

**Institutions and Incentives in Knowledge Production  
and Diffusion:  
From Science to Innovation**

Présentée le 2 octobre 2020

au Collège du management de la technologie  
Chaire en économie et management de l'innovation  
Programme doctoral en management de la technologie

pour l'obtention du grade de Docteur ès Sciences

par

**Charles Chadi AYOUBI**

Acceptée sur proposition du jury

Prof. T. A. Weber, président du jury  
Prof. D. Foray, directeur de thèse  
Prof. R. Veugelers, rapporteuse  
Prof. A. Jaffe, rapporteur  
Prof. G. De Rassenfosse, rapporteur

*A Oum Kalthoum et Sabiha*

# Acknowledgements

Writing a thesis, is a unique experience that shares a lot in a common with the path of an entrepreneur launching a new exciting startup: it is filled with moments of doubt and struggle but all we remember from it eventually is the never-ending learning and the magnificent encounters we build. The PhD has been an occasion for me to learn about theoretical and applied economics but more importantly a beautiful journey during which nothing would have been possible without the support of countless people. I wish first to express my deep gratitude to my advisor, Dominique Foray. Meeting Dominique and starting my PhD with him was an amazing opportunity I am grateful I got. He gave me great autonomy while being constantly supportive in all my endeavors over the last five years. His intuitions and suggestions were crucial for the development of this dissertation. Apart from being a great mentor over the course of my PhD, I had the unique chance of having Dominique as a coauthor of the last paper of this thesis. It was an incredibly enriching experience to learn directly from Dominique's unique economic intuition and his incredibly rich knowledge of the field. All over, the support and the resources he offered were undoubtedly extremely valuable, but his trust and help, especially when I mostly needed it, were priceless. *Merci pour tout.*

I am grateful to all members of my defense committee: Professor Adam Jaffe, Professor Reinhilde Veugelers, Professor Thomas Weber and Professor Gaétan de Rassenfosse. Adam and Reinhilde were daily inspirations for me and their work has inspired many of the reflections that can be found in this manuscript. As for Gaétan, I had the chance to interact with him regularly at EPFL and he has been a great support and an inspiration for me. He was always available and engaged for discussing research ideas and methodologies, I am very grateful for his valuable support. I also feel bound to acknowledge several truly enriching encounters I had with professors in other universities met in various conferences and seminars. Professor Jacques Mairesse who taught me how to keep a high-level perspective on research while keeping a burning passion alive and an inimitable sense of meticulousness. I also thank Ina Ganguli who has been an inspiration with her research papers before being a model for research quality and ethics when I got the chance to meet her in person. I thank Professor Paula Stephan for her constant advice and support, Professor Patrick Gaulé for his inspired comments, Professor Francesco Lissoni for his sense of rigor and passion for research, Professor Annamaria Conti for her inspiring enthusiasm and active support, Professor Mickaël Bikard for his advice and kindness, Professor Stefano Baruffaldi for being an example to follow at every step, and Professor Markus Simeth for his sense of precision and candor. I also thank all the members of the Collegio Carlo Alberto, Max Planck Institute in Munich (MPI), University of Bath, NBER Summer Institute, University of Nice, University of Maastricht, MIT Sloan School of Management, Laboratory of Innovation at Harvard (LISH), Georgia Tech, Exeter Business School, University of Lausanne, University of Geneva, Universidad Pompeu Fabra (UPF), University of Sussex (SPRU), Copenhagen Business School (CBS), and Boston University (BU) for their extremely valuable comments and help. I want to deeply thank all the professors of the Gerzensee Study Center who gave us a beautiful glimpse into what research in economics is about at the very beginning of our PhD experience: Professors Klaus Schmidt, Ricardo Reis, Pietro Gottardi, Sergio Rebelo, Mark Watson, Bo Honoré, Ricardo Alvarez, John Moore, and Jörgen Weibull. Professor Weibull was also always available to help after the end of the courses program, his understanding and explaining of game theory concepts have forever marked the way I reason about these questions and have greatly inspired the fourth chapter of this thesis.

The people I am most indebted however for the completion of this PhD are my coauthors without whom none of the research I conducted would have been possible. When I was still a young and lost PhD student five years ago, I had the unique chance of meeting Fabiana who acted as a true mentor for me. She took the time to involve me in projects she was working on and truly believed in my ability to contribute to them. Fabiana is a model of what every researcher dreams to be, and an amazing supervisor to become for all her future mentees: genuinely helpful, sharp-minded, amazingly organized, hard-working, and kindhearted. None of the projects we have together would have gone this far if it wasn't for Fabiana's sense of balance and intelligence. After almost five years of common

work I am grateful and proud to still be working with her on several exciting research projects. Through Fabiana, I had the chance to meet Michele who a genuinely beautiful encounters of my PhD years. Michele is for sure one of the most brilliant researchers I have met over the last five years. His creativity and practical ingenuity are boundless only rivaling his incredible kindness. I keep on learning every day I interact with him. For both of them, I can confidently say I am grateful to have them as colleagues and now as true friends: Per tutti Grazie mille.

While my relationship with Michele and Fabiana taught me how to become a good researcher, my encounter with Boris is probably the most enjoyably interactive encounter I had during my PhD. I met Boris during the Gerzensee economics coursework program and from that day we developed a deeply beautiful connection and friendship that I had never developed before. I can confidently say that my encounter with Boris transformed my PhD years. I met an amazing coworker, a coauthor, a colleague, a companion of research and maybe more importantly a true friend. I cannot even imagine (and I don't want to) what my PhD years would have been like without my encounter with Boris. He was always there when I had moments of doubt, and when I had almost lost taste for research, he gave me back the passion for questioning things, for digging deeper into the inner mechanisms of human interactions. Our countless debate and discussion added to our perfectly harmonized minds led to several research papers that I am immensely proud of today and I hope for many to come. Un tout grand merci pour tout.

I thank all my colleagues and friends of the EPFL (extended) family. In particular, Yara from bringing joy around her wherever she goes, for always being there for supporting everyone and for reminding me of my oriental vibe, Gabriele Cristelli for always triggering amazingly deep discussions in research and life matters and for being a brilliant colleague and friend, Giada for always being there for me, bringing some colorful vibes to colorless days and for being the older sister I never had, Omar for having the best dark humor there is and for late night discussions the night owls we are had in office, George for showing us that it is possible to play basketball in the office and for his great sympathy, Ling for her continuous kindness and thoughtfulness, Maxime for being a true PhD partner until I realized he was too brilliant to finish his PhD in more than 3 years, Michael for bringing lightness in the offices of CDM, Yu for showing us that skating is possible even for the most hardworking ones, Matthias for showing that you can be the most talkative colleague even by being only once a week in the office, Raphaël for always thriving for things to be better, Paul for being the perfect colleague, Rachel for her sense of humor and kindness, Max for competing with Paul for the title of the perfect colleague and the insightful discussions, Claudia for hosting me when I arrived and always helping me out, Monica for her truthfulness, Abhik for the greatest sarcasm, Corinne for being so sweet, Emilio, Giovanni, Phil and Gabriele Pellegrino for being amazing colleagues as well as great football partners, and Kilian for showing us that angels come by around us to bring light to our lives. All discussions with them were a source of inspiration and motivation. I thank Cyrielle, Amandine, Ema, Ilona, Françoise, Cristina, Céline, Alexandra for making CDM such a joyful place to work in.

I am also grateful to all the great friends I made around Switzerland. The Unil crew that brought laughter and fun my PhD life: Kevin, Elio, Fabio, Silas, Andrea, Valentina, Clem, Vaihbav, Mael, Moritz, Ron, Morten, Fernando, J-live, Pascal, Vani and all the others. And above all Paolita without whom the last five years would have been so much less fun, for all the laughter, the tough moments, the life discussions and the endless absurd debates that I can only hope are the first of many to come: grazie mille. Quentin, for showing me every day how to be beautifully balanced and still amazingly brilliant at all your endeavors, for being a personification of kindness and empathy: je suis on ne peut plus reconnaissant de te compter parmi mes amis. And to the Noëmi, bringing up reality to a world of dreaming researchers and for the sick leaves, danke. To François and Chloé, Merci. Away from Lausanne, I was incredibly lucky to meet wonderful people around the country, the Basel crew, Beau, Andreas, Ali, Manuel, Aimilia, Costanza, Giulia, Taka, Jean-Marc, Martina, Tara, Pedro, Ben, Mauri, Agu, Rahel and Kailin thank you for making every Gerzensee trip a moment I was looking for. Sophia and Noémie thank you for being amazing friends and showing me that Gerzensee friends can work together to organize successful events.

## Acknowledgements

---

To all my one life friends in Paris, Masri, Katra, Attieh, Capo Baroudi, Dib, Safi, Chami, Chacha, Loulou, Tro, Max, Gé, BnB, Bank, Liss, Paki, Markus, Box, JE, Balou, Jack, Guth, VTF, and around the globe, Hakam, Boulos Ghorra, Malak, Yushi, Challoub and a special thank you to the greatest doctor, 7akim el kell Michel: Merci.... thanks especially because it feels like we never got separated when I meet you.

Finally, and most importantly: thanks to my parents, my brother Benjamin and especially to my family here in Switzerland, Mira, Randa, Taha, Carine, Line, Léa and Karim for their endless love, support, for always being by my side, and for making me feel loved and cherished. To Blue, who always reminded me of taking care of myself when I was drown in writing up the pages of this dissertation, danke.

Last but not least, I owe a debt of gratitude to Alice, for her love, for her patience with me in all situations, for her understanding, for helping me out when I needed it the most, for all the moments where she made me feel special, for all the beautiful experiences and the ones to come. You made my PhD years a success at a completely new level and keep on making me every day a better person, there are no words to express how grateful I am to you, for everything: obrigado.

Lausanne, le 19 juin 2020

## Acknowledgements

---

# Abstract

This thesis presents four essays providing novel empirical and theoretical insights on the incentives and institutional structures that favor knowledge production and diffusion. The first two studies analyze these processes in the realm of scientific research, while the last two essays evaluate broader applications with social welfare implications for economists and policymakers alike.

The first essay (chapter 2) of this dissertation, in collaboration with Michele Pezzoni and Fabiana Visentin, exploits a dataset on all applicants to a prestigious Swiss grant to explore a central process in academic life: the application for funds. The results suggest that scientists applying to a grant significantly increase their publications' quality and quantity, learn more, and extend their collaboration network. Beyond the effect of applying, receiving the research funds increases the probability of co-authoring with co-applicants, but it does not have any additional effect on other scientific outcomes. These results justified the title of the chapter since, as it is the case in the Olympics, in research grants, "the important thing is not to win, it is to take part."

The second essay (chapter 3), also in collaboration with Michele Pezzoni and Fabiana Visentin, uses the same empirical context to explore the determinants of knowledge flows among collaborating scientists. The chapter proposes a new methodology based on journal references to track knowledge flows among researchers working together. The results suggest that geographical distance does not significantly affect the knowledge flows between team members, but the cognitive distance separating two members does. More specifically, there is an inverted U-curve effect of cognitive distance on the learning among team members: the higher the distance between two scientists in terms of subjects studied, the more they exchange knowledge, up to the point when the distance becomes detrimental because they have too little common ground to communicate.

The third essay (chapter 4), in collaboration with Boris Thurm, goes beyond the exploration of the determinants for scientists' knowledge production and diffusion to delve into the incentives of all individuals to exchange knowledge. The chapter has two major contributions. First, it acts as a literature review of the empirical evidence on non-financial incentives for knowledge diffusion, such as social recognition, career prospects, and moral considerations. Second, the chapter proposes a simple economic model with heterogeneous agents holding both selfish and moral motives to derive policy implications.

The last essay (chapter 5), in collaboration with Dominique Foray, delves into a specific case of knowledge diffusion, the integration of machine learning technologies in healthcare. The analysis suggests that machine learning has the potential for spurring innovation in healthcare but faces several institutional levers. Collecting quantitative data on patents and publications, and qualitative data on hospitals, the results show that machine learning affects healthcare in different ways than older information and communication technologies. The appearance of new business models encourages tech giants to enter the healthcare sector. These patterns have the potential to increase social welfare by reducing externalities in terms of innovation complementarities, but they pose new challenges such as competition policy and human capital formation.

## Keywords

Knowledge production, knowledge diffusion, incentives, innovation policy, scientific research, social welfare.

# Résumé

Cette thèse présente quatre essais offrant de nouvelles perspectives empiriques et théoriques sur la production et la diffusion du savoir. Les deux premiers essais évaluent ces questions dans le monde de la recherche scientifique, alors que les deux derniers essais explorent ces processus dans un cadre plus large.

Le premier essai utilise les données d'un fond de recherche suisse, pour évaluer l'impact d'une demande de financements de recherche sur la productivité scientifique. Les résultats suggèrent que le simple fait de demander des fonds augmente considérablement la qualité et la quantité des publications, favorise les collaborations scientifiques et permet d'accumuler davantage de connaissances. Au-delà de ces effets, le fait de recevoir des fonds de recherche a un effet positif sur les collaborations scientifiques, mais n'augmente pas significativement la productivité. Ces résultats justifient le titre de l'étude suggérant que, comme aux Jeux Olympiques, pour le financement de la recherche, "le plus important n'est pas de gagner, mais de participer".

Le deuxième essai utilise le même contexte empirique pour évaluer les facteurs favorisant l'échange de connaissances entre chercheurs. L'étude s'appuie sur une nouvelle méthodologie utilisant les références scientifiques pour quantifier les flux de connaissances entre collaborateurs. Les résultats suggèrent que la distance géographique n'a pas d'effet significatif sur l'échange de connaissances. En revanche, la distance cognitive a un effet en U inversé sur le degré d'échange. Autrement dit, plus leurs domaines de connaissances sont différents, plus les chercheurs échangent des connaissances, et ce jusqu'à un point où la distance entre leurs domaines de compétence devient trop grande et leur communication en pâtit.

Le troisième essai va au-delà du monde de la recherche scientifique pour se pencher sur les principes motivant tout un chacun à partager ses connaissances. Cet essai fournit deux contributions. D'une part, grâce à une revue de la littérature, il discute les mécanismes non-économiques favorisant la diffusion du savoir, telles que la reconnaissance sociale, les perspectives de carrière et les considérations morales. D'autre part, le chapitre propose un modèle économique avec des préférences hétérogènes, intégrant le gain personnel et la moralité, qui permet d'évaluer les politiques favorisant le partage de connaissances sous un nouvel angle.

Le dernier essai s'intéresse à un cadre spécifique de diffusion des connaissances : l'intégration des technologies d'intelligence artificielle au domaine de la santé. L'étude suggère que l'intelligence artificielle a le potentiel de stimuler l'innovation dans la santé mais fait face à plusieurs défis institutionnels. L'enquête menée auprès d'hôpitaux ainsi que les données quantitatives sur les brevets et les publications suggèrent que l'intelligence artificielle a un impact différent sur l'innovation dans la santé comparée aux technologies informatiques traditionnelles. De nouveaux modèles commerciaux émergent, avec des régimes d'appropriation basés sur les données, et les grandes entreprises de technologie manifestent leur volonté d'entrer sur le marché de la santé. Ces schémas peuvent favoriser l'innovation en réduisant les externalités, mais posent de nouveaux défis en termes de politique concurrentielle et de formation de capital humain.

## Mots-clés

Production de connaissances, diffusion de connaissances, incitations, politique d'innovation, recherche scientifique, bien-être social.



# Contents

<b>Acknowledgements</b> .....	<b>ii</b>
<b>Abstract</b> .....	<b>vi</b>
<b>Résumé</b> .....	<b>vii</b>
<b>Chapter 1 Introduction</b> .....	<b>1</b>
1.1 Motivation .....	1
1.2 Contribution to the literature .....	2
1.3 Overview of the dissertation essays .....	3
1.4 Policy implications .....	7
<b>Chapter 2 The important thing is not to win, it is to take part: What if scientists benefit from participating in research grant competitions?</b> .....	<b>9</b>
2.1 Introduction.....	9
2.2 Empirical context.....	12
2.3 Data.....	13
2.3.1 Applications and applicants.....	13
2.3.2 Searching for a group of potential applicants .....	17
2.3.3 Scientific outcomes .....	23
2.4 Methodology.....	24
2.4.1 Estimation strategy .....	24
2.4.2 Effect of applying.....	25
2.4.3 Effect of being awarded .....	26
2.5 Results.....	27
2.6 Discussion.....	30
2.7 Conclusion .....	33
<b>Chapter 3 At the origins of learning: Absorbing knowledge flows from within the team</b> .....	<b>35</b>
3.1 Introduction.....	35
3.2 Individual learning: Definition and determinants .....	37
3.2.1 Team characteristics .....	38
3.2.2 Scientist vs. Rest of the team .....	39
3.3 Data.....	41
3.3.1 Team.....	41
3.3.2 Learning.....	41
3.3.3 Empirical setting .....	42
3.3.4 Variables.....	46

---

3.4	Estimation strategy.....	52
3.5	Results.....	52
3.6	Discussion and conclusion .....	58
<b>Chapter 4</b>	<b>Knowledge diffusion and morality: Why do we freely share valuable information with strangers? .....</b>	<b>61</b>
4.1	Introduction.....	61
4.2	Knowledge and data sharing.....	64
4.2.1	The economic properties of shared knowledge.....	64
4.2.2	Non-financial incentives for knowledge sharing.....	66
4.3	Model.....	69
4.3.1	Setting.....	69
4.3.2	Analysis.....	71
4.3.3	Peer pressure.....	76
4.3.4	Financial incentives and morality.....	79
4.4	Discussion.....	82
4.4.1	Contribution of the model.....	82
4.4.2	Policy implications .....	83
4.4.3	Limitations and further work.....	85
4.5	Conclusion.....	86
<b>Chapter 5</b>	<b>Machine learning in healthcare: Mirage or miracle for breaking the costs deadlock? .....</b>	<b>87</b>
5.1	Introduction.....	88
5.2	The innovative significance of ML in healthcare .....	90
5.2.1	Productivity in healthcare .....	90
5.2.2	ML solutions to costly healthcare operations .....	92
5.2.3	Diffusion of ML in Healthcare.....	93
5.3	Supply-side: science, invention and business models .....	96
5.3.1	Evidence on patents and publications .....	97
5.3.2	Main actors and entry of tech firms .....	101
5.3.3	New business models.....	105
5.4	Demand-side: Hospitals transformation and needs .....	107
5.5	Discussion.....	109
5.5.1	The disruptiveness of ML as a solution to spur innovation in healthcare.....	109
5.5.2	Limitations and framework conditions.....	110
5.6	Conclusion.....	112
<b>Chapter 6</b>	<b>Conclusion .....</b>	<b>115</b>
6.1	Scope and findings.....	115
6.2	Future developments.....	115
<b>References.....</b>		<b>117</b>

---

References for chapter 1:.....	117
References for chapter 2:.....	120
References for chapter 3:.....	122
References for chapter 4:.....	125
References for chapter 5:.....	132
References for chapter 6:.....	137
<b>Appendix.....</b>	<b>139</b>
Appendix A2 for chapter 2:.....	139
Appendix A3 for chapter 3:.....	154
Appendix A4 for chapter 4:.....	156
Appendix A5 for chapter 5:.....	161
<b>Curriculum Vitae.....</b>	<b>166</b>

# Chapter 1 Introduction

*“A generation before Nash could have accepted a narrower definition of economics, as a specialized social science concerned with the production and allocation of material goods. [...] But today, economists can define their field more broadly, as being about the analysis of incentives in all social institutions.”* (Myerson, 1999)

During the last couple of months before the finalization of this thesis, our societies have had to face an unforeseen event that distressed our economies and limited several professional activities. Finding effective solutions to save lives and progressively relaunch economic activities depends more than ever on our ability to produce and diffuse knowledge about the disease and potential cures. In this context, better understanding the incentives favoring the process of producing and diffusing knowledge is fundamental. This thesis aims at modestly contributing to this quest.

## 1.1 Motivation

At least since the pioneering works of Romer on endogenous growth (Romer, 1986, 1990), innovation has been recognized by economists as a fundamental driver of economic growth (Grossman and Helpman, 1994). The production of valuable innovations depends on the ability of economic agents to produce and exchange knowledge (Cassiman and Veugelers, 2006; Cohen et al., 1990.; Jaffe, 1986, 1989). Hence, identifying the incentives favoring knowledge production and diffusion and analyzing the institutions governing the process is essential to unravel the mechanisms leading to innovation. However, knowledge suffers from a characteristic that makes it challenging to study: it is hard to measure. As stated by Krugman (1991), “knowledge flows [...] are invisible, they leave no paper trail by which they may be measured and tracked,” which explains our still limited understanding of its underlying mechanisms. This thesis aims at addressing these concerns by investigating empirically and discussing theoretically some of the determinants and implications of knowledge production and diffusion.

The four chapters of this thesis add to the discussion on the mechanisms of innovation by bringing novel insights on the determinants of knowledge production and diffusion in various contexts. It starts by conducting two empirical exercises analyzing the determinants of knowledge production (chapter 2) and diffusion (chapter 3) in the context of scientific research. More specifically, the first essay examines the efficiency of the public funding of scientific research - a traditional policy instrument

- on the quantity and quality of knowledge production, while the second essay studies the determinants of knowledge exchange and diffusion among researchers. Then, the third essay (chapter 4) follows this discussion by proposing a theoretical model that considers the role of non-financial incentives for sharing knowledge. Finally, the fourth essay (chapter 5) wraps up the thesis with an empirical assessment of the patterns of knowledge production and diffusion in the context of machine learning innovation in healthcare.

## 1.2 Contribution to the literature

Knowledge suffers from what economists have defined as a market failure: firms and individuals must invest effort, time, and money to produce new and useful ideas, but everybody can then benefit from their value (Stiglitz, 1999; Foray, 2004). This feature of knowledge, therefore, bears the risk of leading to free-riding problems and an underprovision of the needed knowledge. A classical solution to overcome this risk is the public provision of funds for scientific research to stimulate production and dissemination of new ideas (Arrow, 1962; Dasgupta and David, 1994; Stephan, 2012). Regarding knowledge production, several studies have empirically investigated the efficiency of research funding on scientific production showing that public funding is not as effective in creating incentives to produce knowledge as economic theory would suggest (Jacob and Lefgren, 2011; Gush et al., 2018; Azoulay et al., 2019). Using a unique dataset on Swiss scientists, the first essay shows that the public funding system can also function as an indirect incentive mechanism even for researchers not receiving funds.

Beyond the production of new knowledge, its diffusion is necessary for maximizing social benefits. Although knowledge is a non-rival good - meaning that its exploitation by one agent does not reduce the value for another - the process of knowledge diffusion is not costless (Cowan et al., 2000; Gertler, 2003). Since the seminal work of Jaffe and Trajtenberg (1993) on the impact of geography on knowledge dissemination, several factors affecting the transmission of valuable ideas, such as human interactions and mobility, have been identified (Breschi and Lissoni, 2001; Feldman and Kogler, 2010). The third chapter of this thesis adds to this empirical literature by introducing a measure of knowledge capital for scientists and evaluating the main determinants of knowledge flows among collaborating researchers. The results of the analysis suggest that, in the realm of scientific research, geographical distance and social homophily matter less than the cognitive distance separating researchers.

Following on the determinants of knowledge diffusion, the fourth chapter presents a theoretical framework integrating morality as a motivation for knowledge sharing. In fact, since the seminal work of Arrow (1962) analyzing the optimal incentive system for overcoming the knowledge market failure,

most of the theoretical contribution of economic scholars on the matter have focused on financial instruments (ranging from intellectual property rights to tax credits) to stimulate knowledge production and diffusion (David, 1993; Tirole, 2017). Nevertheless, a growing body of literature has brought empirical evidence on the importance of non-financial incentives for the efficient production and dissemination of knowledge (Xu, 2020; Boudreau and Lakhani, 2015; Gallus, 2017). Chapter 4 of this thesis contributes to the extant theoretical literature on the incentives for knowledge sharing by designing a model consistent with the observations on intrinsic motives and derives novel policy instruments favoring the production and diffusion of knowledge.

The creation of new knowledge and its dissemination among economic agents is most useful when it leads to socially desirable innovations. An area where innovation is direly needed due to alarmingly increasing costs is the healthcare sector (Cutler, 2011; Kocher and Sahni, 2011; Baumol, 2012). Hence, chapter 5 examines the integration of a particular set of knowledge, the promising Machine Learning (ML) technologies, in the healthcare sector, and its potential for increasing the innovative capabilities of the field. Exploiting a recently produced algorithm for the detection of ML patents and publications, the chapter identifies the main institutions governing the demand and supply of this new knowledge and discusses the challenges and opportunities posed by its diffusion.

### 1.3 Overview of the dissertation essays

The first essay exploits a dataset of scientists applying for a leading Swiss funding program to assess the ability of public funding institutions to foster scientific production and collaboration among researchers. The debate in the scientific community on the participation in public funding competitions mainly focuses on the costs they entail for researchers (Ioannidis, 2011; Stephan, 1996). Nonetheless, highly competitive grants require an extensive commitment in the submission phase and strong collaboration among co-applicants. Therefore, this first essay aims at evaluating potential scientific benefits from taking part in research grant competitions. The empirical analysis provides several original findings. First, the comparison of scientific outcomes of applicants to a competitive grant shows little to no significant difference between awarded and non-awarded researchers, which suggests that receiving additional research funds has a limited impact on subsequent research outcomes. Second, considering the global population of all potential applicants to the grant, results point out that, controlling for past trends, applicants to the grant have significantly higher scientific outcomes in terms of quality and quantity. More precisely, scientists participating in the grant application process boost their number of publications, average impact factor, learning, and collaboration network regardless of the success of the application. In other words, these findings suggest that in competitive research funding - like in the Olympics- "the important thing is not to win, but to take part." Finally, citation results indicate that

applicants to the grant witness a decrease in their citations, mainly driven by their exploration of new fields where their reputation is yet to be constructed.

A key mechanism that could explain why applying researchers boost their scientific outcomes is the knowledge they gain by interacting with their co-applicants. Aiming to investigate that hypothesis, the second essay evaluates the determinants of knowledge flows among scientific researchers when working on a common project. Scientific research witnesses a steady and consistent increase in the size of research teams (Wutchy et al., 2007; Jones, 2009), leading to a growing need for efficient knowledge exchange among team members. Using the same empirical setting as the first essay, the second study introduces a new measure of scientists' knowledge stock to assess the determinants of knowledge flows among researchers. The analysis evaluates the effects of three types of distances on the probability of observing knowledge flows among scientists working on a common project: geographical distance, social distance, and cognitive distance. On the one hand, findings suggest that geographical and social distances are not associated with significant differences in knowledge flows. On the other hand, for a scientist applying to a research grant, the probability of learning from a co-applicant is significantly affected by the cognitive distance separating the two. Specifically, the cognitive distance between two co-applicants has an inverted U-shaped effect on the share of knowledge exchanged. Hence, if two researchers share very similar sets of knowledge, then they will have little to learn from each other. Higher cognitive distance is then desirable for increasing knowledge flows but with a certain limit beyond which the two researchers will start to suffer from a lack of common ground to communicate.

Going beyond the scope of scientific research, the third essay delves into the social preferences of individuals to better understand the intrinsic motives that lead them to share valuable knowledge and data even when they incur a personal cost. Information and communication technologies allow individuals, scientists, and organizations to collaborate in new ways (Boh et al. 2007; Williamson et al. 2012). Individuals now share knowledge with unacquainted others and form open collaborations such as open-source software, and firms pool crowds to cultivate better solutions. The classical economic approach with the self-centered *homo oeconomicus* type of preference fails to explain this behavior. If people were simply maximizing their own benefit, they would not put time and effort into sharing knowledge with no insurance of receiving anything in return. This third essay crafts a model of this interaction, precisely tackling the issue of why some individuals are willing to share valuable knowledge at their own cost. The model builds on the existing literature showing that preferences integrating morality are favored by evolution (Alger and Weibull, 2013; Ayoubi and Thurm, 2018) to better portray the behavioral motives of agents. The analysis indicates that it is possible to achieve large welfare increases at a low economic cost by providing easy-access sharing infrastructure and by communicating on the usefulness of the knowledge for other users. This observation suggests that current institutional

systems of knowledge production could be optimized by integrating morality into individual preferences. Discussing non-financial incentives for sharing data and knowledge, this essay contributes to the debate on the design of efficient knowledge dissemination policies.

The fourth essay explores the knowledge diffusion process in the healthcare sector by examining the adoption of Machine Learning (ML), a general-purpose technology, among health experts. The healthcare sector suffers from alarmingly rising costs (Cutler, 2011). The current rapid advances in ML, a subfield of artificial intelligence (AI), offer new automation and prediction capabilities (Brynjolfsson and Mitchell, 2017) that could, if properly integrated, help address the increasing costs' issue. The adoption of ML techniques can, for instance, increase the efficiency and quality of the service and offer new possibilities for prevention, diagnosis, and treatment selection. This essay evaluates to what extent the development of ML-driven solutions can provide relevant opportunities for healthcare innovation. The objective is to improve our understanding of the institutional and organizational conditions required to realize this potential. It provides both qualitative and quantitative evidence on the development of ML in healthcare and discusses the institutional and framework conditions for its successful implementation. Building on a powerful search methodology recently developed by WIPO for patents and publications in ML (WIPO Technology Trends, 2019), the study monitors the production and adoption of ML technology by the healthcare sector using publication and patent data. The chapter proposes two major findings. First, the patenting rate in the field of ML applied to healthcare remains rather low in comparison to a soaring publication rate. This result is mainly driven by new appropriation mechanisms based on the possession of data rather than patenting. Second, ML is allowing the entry of tech giants directly into the healthcare sector. This observation induces both positive externalities (e.g., the reduction of transaction costs) and negative externalities (e.g., a weaker competitive environment due to the increased market power of tech companies).

The following table offers a summarized overview of the chapters of this thesis:



**Table 1.1: Overview of thesis chapters**

	Chapter 2	Chapter 3	Chapter 4	Chapter 5
Title	The important thing is not to win, it is to take part: What if scientists benefit from participating in research grant competitions?	At the origins of learning: Absorbing knowledge from within the team	Knowledge diffusion and morality: Why do we freely share valuable information with strangers?	Machine learning in healthcare: Mirage or miracle for breaking the costs deadlock?
Research Question	What is the impact of applying to a research grant on scientific outcomes?	What are the determinants of knowledge exchange in scientific research team?	Why do individuals share knowledge at their own cost? What are the determinants of that sharing?	What are the patterns of adoption of ML knowledge in healthcare?
Methodology	Empirical econometrics: Propensity score matching and difference-in-differences	Empirical econometrics: Regression with controls	Game theoretical modeling: Utility maximization in a social dilemma	Empirical approach: Quantitative descriptive statistics and qualitative survey data
Data sources	Applicants to SINERGIA + Publication data (Scopus)	Applicants to SINERGIA + Publication data (Scopus)	Descriptive Wikipedia statistics and Scopus publications data	Publication data (Scopus) + Patent data (Patstat) + Survey data
Key findings	Applying to the grant increases publications' quality and quantity, collaboration rate, and learning. Being awarded the grant increases collaborations but has no additional impact on scientific productivity.	Geographical and social distance have little impact on knowledge flows, but cognitive distance has an inverted U-curve effect on knowledge flows.	Morality can explain the sharing behavior of individuals. The perception of the social benefit and the level of sharing in the population affects individual behavior.	Publications soar exponentially, but the patenting rate in ML applied to healthcare is still low. Tech companies are a (growing) key player in the field.

## 1.4 Policy implications

The results exposed in this thesis have direct implications for innovation scholars and policy-makers alike. It adds to our current understanding of the institutions governing knowledge production, as well as of the externalities – both positive and negative - they can produce by intervening in the process of knowledge production and diffusion.

The classical motivation for public funding of science lies in the market failure predicting under-provision of knowledge and insufficient incentives for sharing it (Nelson, 1959; Arrow, 1962; Stephan, 2012). The results of chapters 2 and 5 suggest that the intervention of public institutions can produce unexpected positive externalities beyond merely ensuring the provision of the necessary resources that market mechanisms do not guarantee. Chapter 2 suggests that the public funding of scientific research can create a positive externality by offering an indirect incentive for knowledge production and accumulation, even to the agents not receiving the funds. The process of funding science, therefore, generates “spillovers” to all applicants rather than only benefitting the researchers who got awarded with funds. Similarly, the results of chapter 5 suggest that if the research in a field, namely machine learning, is excessively dominated by private institutions, then the economy might end up leaving the choice of the direction of science in the hands of private interests, not maximizing social welfare. Therefore, public provision of scientific research provides a positive externality to society by ensuring that research production is directed towards socially desirable projects.

By contrast, the results of chapters 3 and 4 indicate that some institutions can also produce undesired negative externalities. Many funding agencies in Europe and the United States have been pushing lately for more interdisciplinary research in the initiatives they support (SINERGIA is an example) intending to stimulate knowledge exchange among researchers. The results of chapter 3 suggest that, while some degree of interdisciplinarity can be desirable for stimulating knowledge flows, there is a limit to the process, and an excessive cognitive distance among collaborators can hinder knowledge diffusion. By assessing the role of non-financial incentives in the production and diffusion of knowledge, chapter 4 suggests that financial rewards might hamper sharing by reducing the intrinsic motives of individuals. As stated by Stephan (1996), scientists have several reasons for producing scientific discoveries, not limited to financial rewards.

The design of optimal public funding policies must account for all these factors to maximize the effective production and diffusion of knowledge. In this sense, the various results exposed in this dissertation have relevant implications for instruments to be implemented by policymakers. For funding agencies, the insights of chapters 2 to 5 suggest at least five different instruments that can lead to more

socially desirable outcomes in terms of knowledge production and diffusion. First, when launching new funding programs, funding agencies should invest in the communication about the call to favor the participation of the largest number of scientists to increase “application spillovers.” Second, the observations of chapter 2 suggest that, for the application effect to function properly, it is preferable to concentrate most of the efforts on the scientific part of the application. For example, scientists applying for SINERGIA insist on how the focus on the scientific part was key to them benefiting from the application phase, as compared to other funding programs. Funding agencies could thus give a higher weight in the evaluation process to the scientific part rather than administrative aspects. Third, the results of chapter 3 suggest that interdisciplinarity is rightfully favored by funding agencies as it boosts knowledge exchange among co-applicants. However, by evaluating applications, funding agencies should give particular care to the degree of interdisciplinarity they encourage, and the tools developed in chapter 3 can provide the proper instrument to estimate the level of diversity of the team in terms of cognitive distance. Fourth, the discussion of chapter 4 suggests that increasing the awareness about the impact of sharing data and knowledge on others can work as a strong incentive for scientists to share both their data and the results of their research, therefore, implying that insisting on that aspect can be an efficient mechanism to increase knowledge diffusion. Finally, the observations of chapter 5 suggest that data functions more and more as an appropriation mechanism for firms, which gives them high market power. Therefore, providing incentives for publicly-sponsored scientists to share their data, and ensuring it respects ethical and confidentiality concerns would have more widespread benefits as the number of individuals capable of using scientific work would rise.

# Chapter 2 The important thing is not to win, it is to take part: What if scientists benefit from participating in research grant competitions?

*Disclaimer: This chapter, written in collaboration with Michele Pezzoni and Fabiana Visentin, is now published in Research Policy, DOI: [10.1016/j.respol.2018.07.021](https://doi.org/10.1016/j.respol.2018.07.021) (Accepted 27 July 2018). Working with Michele and Fabiana has been an amazing learning and collaborating experience. The researcher I am today owes a lot to their support and countless discussions.*

## **Abstract**

“The important thing is not to win, it is to take part,” this famous saying by Pierre de Coubertin asserts that the value athletes draw from Olympic games lies in their participation in the event and not in the gold they collect during it. We find similar evidence for scientists involved in grant competitions. Relying on unique data from a Swiss funding program, we find that scientists taking part in a research grant competition boost their number of publications and average impact factor while extending their knowledge base and their collaboration network regardless of the result of the competition. Receiving the funds increases the probability of co-authoring with co-applicants but has no additional impact on the individual productivity.

## 2.1 Introduction

Throughout the years of economic history, the benefits of competition in terms of social welfare and knowledge production have been debated (Arrow, 1962; Aghion et al. 2005). In contexts where only part of the competitors gets all the monetary reward, competition takes on the characteristics of a race. Research grant competitions offer a stylized example: the scientists who submit the most convincing proposal to the funding agency ‘win’. However, is this competition creating a winners-take-all situation? Or do participants find any benefit in only taking part in the race? Scientific grants are a convenient setting to evaluate the benefits and drawbacks of competition for all contestants, winning or not. We use unique data on a Swiss grant and find that merely taking part in a competition is useful regardless of the result. Specifically, this paper is the first to bring empirical evidence on an overlooked aspect of the research grant process, i.e. the effect of taking part in a grant competition on the scientific productivity,

learning, and collaboration of scientists. Furthermore, we complement the extant literature on the impact of receiving funds from a public funding agency.

When conducting research work, scientists are guided by financial remuneration, puzzle-solving satisfaction, and search for fame and glory (Stephan, 1996). However, nowadays, regardless of their initial motivation, scientists need substantial funding to produce science. Hence, the ability to raise funds is becoming a key skill in managing research laboratories (Etzkowitz, 2003) and a base in the evaluation of scientists' performances along with publication records (Ruben, 2017). Researchers spend an increasing number of hours in writing grant proposals with an uncertain outcome, and, when awarded, in managing the resources they receive. Developing empirical evidence on the benefits and drawbacks of these time-consuming activities would support policymakers and funding agencies in crafting funding systems. However, to our knowledge, extant literature does not include any analysis of the impact of the application process, and studies evaluating funding efficiency are still scarce.

The debate in the scientific community on the opportunity to participate in research grant competitions mainly focuses on the costs they entail. Ioannidis provocatively stated that "the research funding system is broken: researchers don't have time for science anymore. Because they are judged on the amount of money they bring to their institutions, writing, reviewing, and administering grants absorb their efforts" (Ioannidis, 2011). Similarly, Stephan (2010) claimed that "grant applications divert scientists from spending time doing science" and reported an insightful example: "a funded chemist in the U.S. can easily spend 300 hours per year writing proposals". She added that "while some of this effort undoubtedly generates knowledge, much of it is of a 'bean-counting' nature and adds little of social value." (Stephan, 2010). These criticisms are based on the high costs that scientists sustain in applying for competitive grants considered as wasted efforts if the competition turns out to be unsuccessful.

Nonetheless, highly competitive grants require an extensive commitment in the submission phase. Scientists are asked to spend time elaborating an appealing research idea and accurately planning its execution to persuade the evaluators that they will fulfill the promised deliverables. As a matter of fact, as reported by Chubin and Hackett (1990), between 52% and 67% of applicants to NIH and NSF grants pursue the research project they applied for when they did not receive the funds for it suggesting that receiving funds is not the only decisive element in the conduction of research projects. Also, since the grant call is often designed with the requirement of having co-applicants, the application process could be an occasion for scientists to build collaboration linkages (DeFazio et al. 2009).

In this paper, we compare two groups of scientists with the same characteristics differing only in the decision to participate or not in a grant competition. Adopting a difference-in-differences approach, we assess if scientists, who decided to apply, perform differently from the ones who did not. We

use a unique dataset of 775 grant applicants to SINERGIA, a Swiss funding program sponsoring interdisciplinary collaboration where researchers are asked to submit a joint proposal to access funds. We then select a control sample of potential applicants, i.e. scientists with observable characteristics as close as possible to the applicants in our sample using a propensity score matching approach. Since the scientist's observable characteristics used to match applicants with potential applicants might not be perfect proxies for the scientist's quality, commitment, and ability, an instrumental variable approach is added to assure the reliability of our identification strategy.

We find that when applying for a SINERGIA grant, regardless of the result of the application, scientists increase their productivity in terms of number of publications and increase the average impact factor of the journals where they publish. These results suggest that the efforts incurred to apply for the grant pay in the subsequent quality and quantity of the researchers' scientific production. Applicants also expand their collaboration network by co-authoring with their co-applicants. However, in writing multi-disciplinary and long-term projects for a grant like SINERGIA, scientists enter new fields in which they have to acquire new knowledge (Azoulay et al., 2011) and where their reputation requires time to be established. As a result, we observe, for the applicants, a reduction in the average number of citations received per paper.

If on one side, scientists question the utility of participating in a grant competition, on the other side, there is rising attention of researchers for managing their budgets efficiently, partly driven by the growing desire of governments to control public money spending. As an example, illustrating the increasing public pressure on scientists, since the early nineties, the U.S. government is asking funding agencies to report the outcomes of projects publicly supported. The U.S. Government Performance and Results Act of 1993 states that "the Director of the Office of Management and Budget shall require each agency to prepare an annual performance plan covering each program activity set forth in the budget of such agency. Such plan shall "[...] establish performance goals to define the level of performance to be achieved by a program activity; [...] establish performance indicators to be used in measuring or assessing the relevant outputs, service levels, and outcomes of each program activity" (sec. 2803). As an example of an action taken to respond to such government regulatory interventions, one of the largest American evaluation programs assessing the impact of public investment in research, STAR METRICS, was launched to provide taxpayers with precise information on the value of their investments (Lane, 2011).

Despite a growing demand for an evaluation of publicly supported scientific research, extant studies do not provide convergent findings on the effect of receiving funds on researchers' scientific outcomes. The results suggest a limited impact of funding on the main scientific outcomes of scientists,

but the magnitude of the effects and the outcomes analyzed vary across studies (Arora and Gambardella, 2005; Jacob and Lefgren, 2011; Gush et al., 2015; Azoulay et al., 2015; Carayol and Lanoe, 2017).

The disparity in the empirical findings could result from several technical limitations. In a comprehensive review, Jaffe (2002) identifies three main difficulties encountered when attempting to evaluate the effects of research funding. First, there might be information availability issues since it is often difficult to retrieve detailed information about the full sample of scientists applying for a grant, awarded and non-awarded, and demographic information about the studied scientist is sparse. Second, even with the accessibility to such data, the estimation of the funding effect might be biased because most productive scientists are also the ones having a higher probability to be funded. Third, the standard bibliometric measures, such as the number of the publications and the average impact factor, might provide only a partial picture of the effects of being awarded a grant.

Our study proposes a set of solutions to tackle these obstacles. We use a comprehensive dataset of scientists, including both awarded and non-awarded applicants and, following Fox (1983), we exploit the richness of our dataset to include both individual-level variables and environmental characteristics in our analysis. We then introduce new scientific outcomes to capture more extensive aspects of scientific production, such as collaboration and learning.

In our analysis of the impact of funding, we find that receiving funds represents a proper incentive to realize the potential collaborations claimed in the application phase. Precisely, we find that the probability of co-authoring with at least one other scientist listed in the application is higher within awarded applications. However, concerning the productivity of funded researchers, we observe similar results to the ones of the literature. Being awarded has a limited but not significant impact on the quality and quantity of a researcher's scientific productivity.

The organization of this paper is as follows: Section 2.2 sets the empirical context, Section 2.3 describes the data and main variables, Section 2.4 exposes the estimation strategy, Section 2.5 presents the findings, Section 2.6 discusses the results, and Section 2.7 concludes.

## 2.2 Empirical context

The Swiss National Science Foundation (SNSF) is the main national funding agency in Switzerland. It plays in the country the same role of the National Science Foundation (NSF) in the United States or the European Research Council (ERC) in Europe. The SNSF supports researchers' activities and their careers. SINERGIA program is one of the flagship programs in its portfolio. It was launched in 2008 and designed to promote breakthrough research and collaboration of scientists affiliated with different

institutions. As mentioned in the application guidelines, scientists are required to collaborate with colleagues from another institution as a condition for securing research funding, i.e., scientists need to submit a proposal for a “research work carried out collaboratively” (SNSF, 2011). The application process for SINERGIA is similar to the one of NSF and ERC grants. Researchers based in research universities and public research institutions obtain public funds on a competitive basis by submitting their proposal to the selection committee of the SNSF. The committee then selects the most promising projects to which the funds are allocated.

In most cases, a SINERGIA project involves four or five scientists led by a main proponent coordinating the overall project. All disciplines are eligible for funding through the program. Applicants propose interdisciplinary projects or projects where co-applicants belong to the same field but are specialized in different subfields.<sup>1</sup> The criteria considered in evaluating the application are the value added by the joint research approach, the research complementarities of the applying groups, and the coherence of the projected collaboration. The screening of applications is a two-step evaluation process. In the first step, external reviewers assign a provisional score to each application. In the second step, an internal committee of the SNSF, the Specialized Committee for Interdisciplinary Research, based in Bern, assigns a final score to each application using a scale where 6 is the highest score and 1 the lowest. The evaluation process takes six months to be completed. The decision to award projects is also based on the funds available; all awarded projects received funds.

## 2.3 Data

This section describes the characteristics of applications and applicants (paragraph 2.3.1), illustrates the procedure applied to select a control sample of potential applicants (paragraph 2.3.2), and presents the scientific outcomes (paragraph 2.3.3).

### 2.3.1 Applications and applicants

The scarcity of information disclosed by the funding agencies about their application selection process has often limited the capacity of scholars to estimate the effects of public funding activities adopting ideal identification strategies. Our scientific partnership with SNSF provided us with the opportunity to have all – both awarded and non-awarded - grant applications submitted by Swiss researchers applying for the SINERGIA grant in the period 2008-2012.<sup>2</sup> We also have access to the scores

---

<sup>1</sup> An example of a project in different disciplines is one including Math, Hydrology and Geophysics, while an example of a single-discipline project with two sub-disciplines is one with Biochemistry and Genetics.

<sup>2</sup> All concerned applicants were contacted by the SNSF and had the possibility to oppose the transmission of their data.



assigned to the applications, the final funding decisions, and demographic information about applicants. We complement this information with applicants' publication records using the Elsevier's Scopus database. To perform our analysis, we select applications in Engineering and Science & Medicine.<sup>3</sup> Our final sample includes 255 grant applications and 775 distinct applicants. Our unit of analysis is the pair applicant-application. Considering that each applicant can be involved in more than one application, our sample counts 1,060 applicant-application pairs. Precisely, in 22% of the cases, applicants persistently apply by participating in more than one call. However, only 8% of the applicants apply again after having been awarded.

As application characteristics, we consider five sets of variables measuring the application funding decision, the project size, its quality, the applicant team composition, and discipline. We capture the funding decision using a dummy that equals one if the application is awarded, zero otherwise (Awarded). The size of the research project is proxied by two variables, the Amount requested in Swiss Francs (CHF) and the number of co-applicants listed in the application document (N. of co-applicants). The quality of the project is evaluated using a variable that ranges from 1 to 6, according to the grade assigned by the selection committee to the application (Grade). We proxy the ethnic composition of the applicants' team using a dummy that equals one if all the applicants are affiliated with Swiss institutions (Swiss team) and their geographical dispersion using a continuous variable measuring the average distance in terms of travel time between the researcher's affiliation and the co-applicants' affiliations (Distance hours). For gender, we use a dummy that equals one if there is at least one female researcher among the co-applicants (At least one female researcher on the team). Finally, we identify the discipline of the application with a dummy that equals one if the application is in the domain of Science and Medicine and zero if the application is in the domain of Engineering (Science & Medicine).

Table 2.1 reports the key figures describing application characteristics. The applications in our sample are composed, on average, by 4.19 members, with a minimum of 2 and a maximum of 11 members. Concerning team composition, about 13% of the teams have only Swiss members, while the others are multinational teams. When classified by discipline, 36% of the applications are in Engineering, whereas 64% are in Science & Medicine. A SINERGIA grant covers personnel costs, research costs, coordination costs, and, to a limited extent, investment costs. The average amount requested per application is 1.67 million CHF, with a minimum of 0.35 million CHF and a maximum of 6.85 million CHF. Figure 2.1

---

<sup>3</sup> In this study, we exclude from the original sample applications in the Humanities and Social Sciences since book contributions represent a significant part of the field publication outcomes and are not collected with accuracy in the Elsevier's Scopus database. Applications in the Humanities and Social Sciences represent 19% of the total sample.

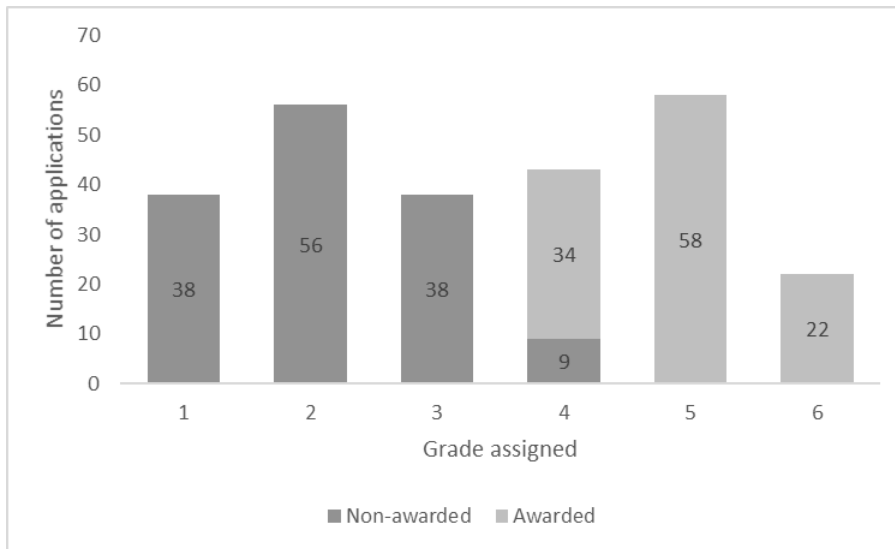
represents the distribution of the number of grant applications by the score assigned and the final funding decision. A total of 9% of the applications obtained the maximum score, 6, and 45% of the applications were awarded.

**Table 2.1: Application characteristics (Number of applications=255).**

	Mean	Std. Dev.	Min	Max
Awarded	0.45	0.50	0	1
Amount requested (in million CHF)	1.67	0.76	0.35	6.85
N. of co-applicants	4.19	1.59	2	11
Grade	3.39	1.60	1	6
Swiss team	0.13	0.33	0	1
At least one female researcher on the team	0.36	0.41	0	1
Distance hours	4.39	3.59	1	19.75
Science & Medicine	0.64	0.48	0	1

The table shows the key figures concerning the 255 grant applications included in our sample. The table reports mean, standard deviation, minimum and maximum values of five sets of variables regarding the application funding decision (*Awarded*), the project size (*Amount requested* and *N. of co-applicants*), its quality (*Grade*), the applicant team composition (*Swiss team*, *At least one female researcher on the team*, and *Distance hours*), and the main discipline of the application (*Science & Medicine*).

**Figure 2.1: Distribution of the number of grant applications by the score assigned and final funding decision.**



The figure shows the distribution of the grant applications according to the grade assigned by the evaluation committee. The grades range from 1 to 6, where 6 is the highest grade.

As applicant characteristics, reported in Table 2.2, we consider two sets of variables, respectively measuring the applicant's demographic and bibliometric characteristics before the application year. The demographic characteristics include a dummy that equals one if the researcher is a female and zero otherwise (*Female*) and the variable *Seniority* that measures the time since the first year of scientific activity of the researcher<sup>4</sup>. We consider a set of bibliometric measures proxying the applicant's publication characteristics before the application year. These measures include a set of variables computed in the five years preceding the application: the publication count (*Publication count pre-application*), the average impact factor of the journals where the applicant published (*Average IF pre-application*), the average number of citations received per paper (*Average citations pre-application*), the yearly average number of authors per paper (*Average authors pre-application*), and the existence of at least one co-authored paper between the scientists and the other applicants (*Co-applicant collaboration pre-application*).

The SINERGIA funding program targets established researchers who demonstrated their ability to conduct excellent quality independent research. In most cases, applicants are associate or full professors with good publication records. The average seniority of the applicants is 18.52 years since the start of their research activity. The average number of applicants' publications is 31.35 in the five

<sup>4</sup> To identify this beginning year, we track the first publication authored by the researcher looking at her self-citations within the publications available in our study sample (2003-2015). We consider this publication year as a proxy for the year when the researcher started her activity.

years preceding the application year. An average applicant has received 4.28 yearly citations per paper at the application time and has published on journals with an average impact factor of 5.59. When looking at gender distribution, 15% of applicants in our study sample are female. In 40% of the cases, the applicant researchers have established collaborations with the other applicants, i.e., they have co-authored at least one article with them in the five years preceding the application.

**Table 2.2: Applicant characteristics (Number of applicants=775).**

	Mean	Std. Dev.	Min	Max
Seniority	18.52	9.37	0.00	53.00
Female	0.15	0.36	0.00	1.00
Publication count pre-application	31.35	27.29	1.00	225.00
Average IF pre-application	5.59	3.60	0.10	28.61
Average citations pre-application	4.28	4.26	0.04	48.62
Average authors pre-application	5.14	1.13	1.00	10.40
Co-applicant collaboration pre-application	0.40	0.47	0.00	1.00

The table shows the key figures concerning the 775 grant applicants included in our sample. It reports the mean, standard deviation, minimum and maximum values of two sets of variables regarding the applicants' demographic characteristics (*Seniority* and *Female*) and the applicants' bibliometric characteristics (*Publication count pre-application*, *Average IF pre-application*, *Average citations pre-application*, *Average authors pre-application*, and *Co-applicant collaboration pre-application*).

### 2.3.2 Searching for a group of potential applicants

To estimate the effect of applying for a grant, we construct a control sample of researchers who would have been eligible to apply but did not apply for a SINERGIA grant. To do so, we retrieve a group of potential applicants with profiles similar to the ones of the applicants of our sample. We find a potential applicant, i.e., a matched control, for each of the 1,060 applicant-application pairs in our sample. To identify potential applicants, we proceed in two steps. First, we define a large pool of scientists eligible to apply for SINERGIA. Second, we extract from this pool of scientists the ones who match the profiles of the applicants using a propensity score matching approach.

We consider as eligible scientists all the publishing scientists affiliated with one of the twelve major Swiss universities<sup>5</sup>. From all the publications of the scientists affiliated with those universities in the period 2003-2015, we retrieve 25,715 authors who were active in the period 2008-2012, i.e., the period during which the SINERGIA grants were awarded. We consider a scientist active in a given year

<sup>5</sup> University of Neuchatel, ETHZ, EPFL, University of Lausanne, University of Fribourg, University of Genève, University of Bern, University of Basel, University of Lugano, University of Zurich, University of Luzern, and University of St. Gallen.

$t$  if she has at least one publication in the time window  $[t-5, t-1]$  and at least one in  $[t, t+4]$ . Each of the 25,715 scientists is observed yearly, leading us to a pool of 86,694 scientist-year pairs.

To extract from the pool of scientist-year pairs the most appropriate control for each applicant-application pair, we use a propensity score matching based on a logit estimation of the probability of applying for SINERGIA<sup>6</sup>. In this estimation, the dependent variable (*Applicant*) equals one for the 1,060 applicant-application pairs and zero for all the remaining scientist-year pairs in the pool. We identify 1,060 controls, one for each of the 1,060 applicant-application pairs. We define the 1,060 controls as the potential applicants for the SINERGIA grant.

To identify potential applicants, we consider as relevant matching characteristics: the researcher's *Seniority*, her fundraising profile, and her bibliometric characteristics before the application year<sup>7</sup>. The fundraising profile of a scientist is captured through two variables: *Other active funding* and *Previous expired funding*. The variable *Other active funding* is a dummy, which equals one if the scientist has at least one active project granted (other than SINERGIA) at the moment of the application and zero otherwise. In contrast, the variable *Previous expired funding* is a dummy that equals one if the scientist has raised funds in the past with a grant that was expired at the moment of the application to SINERGIA and zero otherwise. As funding, we consider the European Union grants<sup>8</sup> and the SNSF grants other than SINERGIA<sup>9</sup>.

The bibliometric characteristics include the variables *Publication count pre-application*, *Average citations pre-application*, *Average IF pre-application*, *Average authors pre-application*<sup>10</sup>. To improve the matching quality, we include in the regression the average yearly variation of the bibliometric characteristics. These additional variables allow us to take also into account the trends of these bibliometric indicators over the 5-year window considered. Precisely, we calculate the average yearly growth (decline) of the publication count over the five years of observation preceding the application year (*Average publication trend*). Similarly, we construct the variable *Average citation*

---

<sup>6</sup> To be conservative, and exclude the possibility that some of our results are driven by some loosely matched controls, we also limit the analysis to the controls having a propensity score different by less than 1% compared to the one of the actual applicant they are matched with. Our results remain stable across the two approaches.

<sup>7</sup> For a potential applicant, we consider as application year the year when her matched applicant submitted her application.

<sup>8</sup> From the European Union we retrieved the data of the grants awarded between 1998 and 2013, namely FP5 (1998-2002), FP6 (2002-2006) and FP7 (2007-2013). We collected this data on the CORDIS platform online (<https://data.europa.eu>).

<sup>9</sup> For Switzerland, we used the P3 database of the Swiss National Science Foundation (SNSF) which makes data available on the projects and people that it has supported (<http://p3.snf.ch>).

<sup>10</sup> As a robustness check, we considered the logarithmic transformations of *Publication count pre-application*, *Average IF pre-application*, *Average citations pre-application*, and *Average authors pre-application* when predicting the propensity scores presented in Table 2.3. The estimation of the coefficients of the interaction term *Applicant\*Post-Application*, in this case, are similar to the ones reported in Table 2.6. The results of this check are available upon request.

*trend*, *Average IF trend*, and *Average authors trend*. To account for the continuity in the productivity of the scientists, we add as further bibliometric characteristic the variable *Productivity break*, which counts the number of years without any publication in the five-year window preceding the application year (Mairesse and Pezzoni, 2015).

We also include as relevant matching characteristic the stock of knowledge of the scientist before the application year. The stock of knowledge, as more extensively described in paragraph 2.3.3, is represented by the number of distinct journals listed in the references of the articles published by the scientist in the five years preceding the application (*N. of journals pre-application*).

Finally, we control for *Application year*, *Affiliation*, and *Discipline*<sup>11</sup> fixed effects. Table 2.3 shows the logit estimation of the regression used to predict the matching probability.

---

<sup>11</sup> See Appendix A2.1 for a detailed description of the attribution of disciplines to scientists.

**Table 2.3: Propensity score matching regression, logit estimation.**

	Logit Applicant
Seniority	0.00015 (-0.0094 ; 0.0097)
Other active funding	1.13*** (0.92 ; 1.33)
Previous expired funding	0.16 (-0.045 ; 0.37)
Publication count pre-application	0.018*** (0.013 ; 0.023)
Average IF pre-application	0.079*** (0.050 ; 0.11)
Average citations pre-application	-0.13*** (-0.17 ; -0.10)
Average authors pre-application	0.71*** (0.66 ; 0.75)
Average publication trend	0.040 (-0.024 ; 0.10)
Average IF trend	-0.044 (-0.12 ; 0.027)
Average citation trend	-0.023 (-0.088 ; 0.042)
Average authors trend	-0.045 (-0.16 ; 0.066)
Productivity break	-1.49*** (-1.63 ; -1.34)
Log(N. of journals pre-application)	0.74*** (0.58 ; 0.91)
Constant	-8.27*** (-9.06 ; -7.48)
Dummy Application year	Yes
Dummy Discipline	Yes
Dummy Affiliation	Yes
N. of Scientists	25,715
Observations	86,694
Pseudo R2	0.60

The table shows the coefficients estimated for the regression predicting the probability of applying for a SINERGIA grant as a function of a scientist's demographic and bibliometric characteristics. The regressors include fixed effects for *Application year*, *Affiliation*, and *Discipline*. The dependent variable is a dummy that equals one if the scientist applies to SINERGIA, zero otherwise. The 86,694 observations refer to the entire pool of scientist-year pairs from which we extract 1,060 potential applicants, one for each 1,060 applicant-application pairs. We estimate a Logit model. In reporting the statistical significance of the coefficients, we apply the standard thresholds i.e. \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ . 95% confidence intervals are reported in parenthesis below each coefficient following Cummings (2013, 2014).

Table 2.4 shows the descriptive statistics for four groups of researchers: awarded applicants, non-awarded applicants, the entire set of applicants, and the potential applicants selected with the propensity score matching approach. Statistical t-tests on the averages of the main bibliometric characteristics of applicants (column 3) and potential applicants (column 4) do not reject the null hypothesis that the two groups have the same averages. This evidence confirms the accuracy in the selection of potential applicants with profiles similar to the applicants. The only case where we reject the null hypothesis of the t-test is for the variable *N. of journals pre-application*<sup>12</sup>.

---

<sup>12</sup> We conducted two additional robustness exercises. One where we restrict the study sample to the applicant-potential applicant pairs having a highly similar propensity score and one where we construct a control sample for which *N. of journals pre-application* is the main matching criterion. Our results remain stable across these two additional control samples. The results of the econometric exercises described in Section 2.4 (Methodology) and conducted using these two alternative control samples are reported in Appendix A2.4.



**Table 2.4: Descriptive statistics for the four groups of scientists before the application year.**

	(1)			(2)		
	Awarded (469 obs.)			Non-Awarded (591 obs.)		
	Mean	Min	Max	Mean	Min	Max
Publication count pre-application	30.23	1.00	225.00	34.34	2.00	183.00
Average IF pre-application	6.14	0.10	19.97	5.19	0.10	28.61
Average citations pre-application	4.61	0.12	38.67	4.03	0.04	48.62
Co-applicant collaboration pre-applica- tion	0.37	0.00	1.00	0.42	0.00	1.00
N. of journals pre-application	125.35	1.00	495.00	135.50	1.00	594.00
Average authors pre-application	5.26	1.00	10.40	5.08	1.33	10.16
Seniority	17.37	0.00	50.00	17.95	0.00	52.00
Other active funding	0.48	0.00	1.00	0.38	0.00	1.00
Previous expired funding	0.37	0.00	1.00	0.35	0.00	1.00
Productivity break	0.50	0.00	4.00	0.41	0.00	4.00
	(3)			(4)		
	Applicants (1,060 obs.)			Potential Applicants (1,060 obs.)		
	Mean	Min	Max	Mean	Min	Max
Publication count pre-application	32.52	1.00	225.00	30.59	1.00	191.00
Average IF pre-application	5.61	0.10	28.61	5.51	0.10	35.21
Average citations pre-application	4.29	0.04	48.62	4.18	0.01	32.24
Co-applicant collaboration pre-applica- tion	0.40	0.00	1.00	0.02	0.00	1.00
N. of journals pre-application	131.01	1.00	594.00	114.50	2.00	515.00
Average authors pre-application	5.16	1.00	10.40	4.87	1.25	15.00
Seniority	17.69	0.00	52.00	17.93	0.00	52.00
Other active funding	0.43	0.00	1.00	0.46	0.00	1.00
Previous expired funding	0.36	0.00	1.00	0.40	0.00	1.00
Productivity break	0.45	0.00	4.00	0.54	0.00	4.00

The table shows the average bibliometric characteristics before the application year for four groups of scientists: Awarded applicants (469 obs.), non-awarded applicants (591 obs.), the whole sample of the applicants (1,060 obs.), and potential applicants selected according to the propensity score matching procedure (1,060 obs.). Comparing the characteristics of Applicants and Potential applicants (Columns 3 and 4), we find no statistical difference, i.e.,  $p\text{-value} > 0.1$ , of the means of the two groups in terms of *Publication count pre-application*, *Average IF pre-application*, *Average citations pre-application*, *Average authors pre-application*, *Seniority*, *Other active funding*, *Previous expired funding*, and *Productivity break*.

### 2.3.3 Scientific outcomes

We measure scientists' outcomes along five dimensions. The first three dimensions are standard bibliometric measures commonly used in the literature. We then add two new measures to draw a more extensive picture of the outcomes that scientists can achieve in their scientific research. All these variables are computed in the five-year window after the SINERGIA application year.

In assessing publication quantity and quality of scientists, we use standard bibliometric measures: publication count, citations, and impact factor. Following the mainstream literature (Arora and Gambardella 2005; Carayol and Lanoe 2017), we account for the publication quantity as the count of papers published by a scientist (*Publication count*). As for publication quality, we use impact factor (*Average IF*) and citations (*Average citations*) as two quality measures capturing different aspects of the impact of a publication. Citations indicate the number of times other scientists refer to a work and could be defined as the "peer assessments of the significance of scientists' contributions" (Cole and Zuckerman, 1984, pp. 231). The impact factor expresses the quality of the journal where a work is published (Long, 1992). For a work, being published in a journal with a high impact factor is commonly seen as a signal of high intrinsic quality.

When measuring the impact of a scientific grant with the standard bibliometric outcomes, as stated by Jaffe (2002), we are only "examining this tiny piece of a very complicated puzzle [...], we are not looking at the spillovers that are perhaps the true reason for these programs." In our case, collaboration and knowledge transfer are stated as "true reasons" for implementing the grant. We consider two novel measures aiming to capture these two core aspects of the SINERGIA grant: successful collaboration and individual learning.

To measure successful scientific collaborations, we consider the work relationship consolidated among SINERGIA co-applicants through a co-authorship. More specifically, we measure a successful collaboration as a dummy (*Co-applicant collaboration*) that equals one if a scientist co-authors at least one paper with her co-applicants and equals zero if she does not.

To measure the knowledge transfer among scientists, we construct a new measure of individual learning. As stated by Huber (1991), "Learning consists in knowledge acquired by any unit of an organization and available for acting upon." Hence, building on Ayoubi et al. (2017), we define individual learning as the increment to the stock of knowledge of a scientist between two periods. Following Uzzi et al. (2013), we use the journals cited as the building blocks of knowledge and thus consider the knowledge stock of a scientist at a given point in time as the journals that she cited in her publications. Our proxy for individual learning is then the difference between the knowledge stock of the scientist before and after the application time. More specifically, we measure individual learning at a given

moment as the number of distinct journals in the references of the publications of the scientist not observed in a previous period (*Learning*).

Table 2.5 shows the descriptive statistics of the scientist's outcomes for the applicants (distinguishing awarded and not awarded) and the potential applicants selected.

**Table 2.5: Descriptive statistics for the scientist's outcomes after the application year.**

	(1)			(2)		
	Awarded (469 obs.)			Non-Awarded (591 obs.)		
	Mean	Min	Max	Mean	Min	Max
Publication count	36.28	1.00	189.00	38.64	2.00	251.00
Average IF	5.95	0.73	20.48	4.97	0.10	19.78
Average citations	2.82	0.03	17.38	2.36	0.06	27.55
Co-applicant collaboration	0.71	0.00	1.00	0.61	0.00	1.00
Learning	95.46	1.00	368.00	101.74	2.00	556.00
	(3)			(4)		
	Applicants (1,060 obs.)			Potential Applicants (1,060 obs.)		
	Mean	Min	Max	Mean	Min	Max
Publication count	37.60	1.00	251.00	30.39	1.00	165.00
Average IF	5.40	0.10	20.48	5.21	0.10	31.96
Average citations	2.56	0.03	27.55	3.93	0.11	78.13
Co-applicant collaboration	0.66	0.00	1.00	0.01	0.00	1.00
Learning	98.96	1.00	556.00	70.22	1.00	318.00

The table shows the descriptive statistics of the scientific outcomes after the application year for four groups of scientists: Awarded applicants (469 obs.), non-awarded applicants (591 obs.), the whole sample of the applicants (1,060 obs.), and potential applicants selected according to the Propensity Score Matching procedure (1,060 obs.). The five variables included in the table, *Publication count*, *Average IF*, *Average citations*, *Co-applicant collaboration*, and *Learning*, are the dependent variables of the regression models estimated in Section 2.5.

## 2.4 Methodology

### 2.4.1 Estimation strategy

To estimate both the effects of applying for a grant and receiving the funds, we rely on a difference-in-differences approach where we compare the changes in scientific outcomes between before and after applying or receiving the funds (Angrist and Pischke, 2008). The identifying assumption is that a

scientist's outcome trends would be the same in the absence of the treatment represented by applying for the grant program and receiving the funds<sup>13</sup>, respectively.

#### 2.4.2 Effect of applying

To formally evaluate whether applying has a significant effect on the subsequent scientist's outcomes, we compare the applying scientists to the matched control sample of potential applicants described in paragraph 2.3.2 above.

For estimating the application effect, we rely on the equivalent formulation of the difference-in-differences reported in Equation 2.1 and estimated with an Ordinary Least Squares (OLS).

$$Scientist's\ outcome_{it} = \beta_0 + \beta_1 Applicant_i + \beta_2 Post-Application_{it} + \beta_3 (Applicant_i * Post-Application_{it}) + (Scientist's\ characteristics)_i' \beta_4 + \varepsilon_{it}$$

(Equation 2.1)

Where  $i$  and  $t$  refer to the scientist  $i$  observed at time  $t$ . We observe the scientist in two periods, before her application ( $t=0$ ) and after her application ( $t=1$ ). To define the two periods for a potential applicant who, by definition, does not apply to SINERGIA, we use the application year of her matched applicant.

The variable *Scientist's outcome<sub>it</sub>* is, in turn, one of the five dependent variables described in paragraph 2.3.3 (*Publication count, Average IF, Average citations, Co-applicant collaboration, Learning*). All the dependent variables, except for *Co-applicant collaboration*, are log-transformed<sup>14</sup>.

The dummy *Applicant<sub>i</sub>* equals one if scientist  $i$  is an applicant to SINERGIA and zero otherwise. *Post-Application<sub>it</sub>* is a time dummy that takes a value of zero if we observe the scientist  $i$ 's outcomes before her application ( $t=0$ ), and a value of one if we observe the scientist's outcomes after ( $t=1$ ). The interaction term *Applicant<sub>i</sub>\*Post-Application<sub>it</sub>* marks a scientist  $i$  who experienced the application and whose outcomes are observed after the application time. The estimated coefficients of this interaction measure the effect of applying for a SINERGIA grant<sup>15</sup>.

---

<sup>13</sup> We formally test for this assumption using the fully flexible model for parallel paths introduced by Mora and Reggio (2017) and find that our sample fulfills this identifying assumption.

<sup>14</sup> Note that this log transformation is possible since all the values of the variables of interest are strictly positive for all the scientists of our sample (See minimum values in Table 2.5).

<sup>15</sup> In comparing applicant versus potential applicant scientists, the heterogeneity of the former category might bias our results. Specifically, applicant scientists include awarded and non-awarded scientists that might behave in a different way. To address this concern, we run a robustness check where we include as controls the variable *Awarded* and the interaction term *Awarded\*Post-Application*. The new results, available in appendix A2.4, are consistent with the results reported in Table 2.6.

The vector *Scientist's characteristics* includes the variable *Seniority* as well as the *Application year*, *Affiliation*, and *Discipline* fixed effects.

A potential issue with our results could come from the endogeneity arising due to applicants' self-selection. The propensity score matching, along with the difference-in-differences approach described in this section, is meant to take into account this concern. In Appendix A2.2, we present an additional estimation exercise using an instrumental variable approach evaluating the impact of applying for a grant (Wooldridge, 2012). Our main results remain stable across the different econometric exercises.

### 2.4.3 Effect of being awarded

After having evaluated the impact of applying, we focus our attention on the subsample of applicant scientists, and we estimate the effect of being awarded a SINERGIA following Equation 2.2.

$$\begin{aligned} \text{Scientist's outcome}_{it} = & \beta_0 + \beta_1 \text{Awarded}_i + \beta_2 \text{Post-Application}_{it} + \beta_3 (\text{Awarded}_i * \text{Post-Application}_{it}) \\ & + (\text{Scientist's characteristics}_i)' \beta_4 + (\text{Application characteristics}_i)' \beta_5 + \varepsilon_{it} \end{aligned}$$

(Equation 2.2)

Where  $i$  and  $t$  refer to the scientist  $i$  observed at time  $t$ . We observe the scientist in two periods, before the application ( $t=0$ ) and after the application ( $t=1$ ). The variable *Scientist's outcome* $_{it}$  is, in turn, one of the five dependent variables described in paragraph 2.3.3 above (*Publication count*, *Average IF*, *Average citations*, *Co-applicant collaboration*, *Learning*). All the dependent variables, with the exception of *Co-applicant collaboration*, are log-transformed<sup>16</sup>.

The dummy  $\text{Awarded}_i$  equals one if scientist  $i$  was awarded a SINERGIA grant and zero otherwise.  $\text{Post-Application}_{it}$  is a time dummy that takes a value of 0 if we observe the scientist's outcomes before the application ( $t=0$ ), and a value of 1 if we observe the scientist's outcomes after ( $t=1$ ). The interaction term  $\text{Awarded}_i * \text{Post-Application}_{it}$  marks a scientist  $i$  who was awarded a grant and whose

---

<sup>16</sup> Although some previous studies apply a Regression Discontinuity Design (RDD) approach (Jacob and Lefgren, 2011), this approach is not suitable in our case. The main reason is that we do not have a ranking of the applications but six grades, from 1 to 6. For us, the threshold is located at grade 4 since all the applications above it are awarded while none of the ones below are. We cannot rank the applications within grade 4 to position correctly the threshold. This limitation is related to the procedure used by SNSF to select the awarded applications graded 4. All the application graded 4 go through a second round of evaluation where the committee revise each application but without producing a ranking. In order to test the robustness of our results with an approach similar to the RDD, we conduct our analysis by restricting our sample to applications of comparable quality, i.e. those graded 3 or 4. The results, available in appendix A2.4, are consistent with the results reported in Table 2.7.

outcomes are observed after the application time. The estimated coefficients of this interaction measure the effect of being awarded a SINERGIA grant.

The vector *Scientist's characteristics* includes the variables *Seniority*, *Female*, *Other active funding*, and *Previous active funding* as well as the *Application year*, *Affiliation*, and *Discipline* fixed effects. The vector *Application characteristics* includes the variable *Grade*, the dummies *Swiss Team*, *At least one female researcher*, *Science & Medicine*, and the continuous variables *Log(Amount requested)*, *Log(N. of co-applicants)*, *Log(N. of disciplines)*, and *Log(1+Distance hours)*. We also control for the presence of previous applications to SINERGIA and previously awarded SINERGIA applications using two dummies, *Previous SINERGIA application* and *Previous SINERGIA awarded*. The first variable equals one if the applicant has at least one previous application to SINERGIA, and the second equals one if the applicant has at least one previous SINERGIA awarded, zero otherwise.

## 2.5 Results

This section summarizes the results of the two main regressions described in Section 2.4. First, we present the main findings for the difference-in-differences regression estimating the effect of applying for the SINERGIA grant (Equation 2.1)<sup>17</sup>. Second, we present the results of the regression estimating the impact of receiving the funds (Equation 2.2). Following the approach of Cummings (2013, 2014), we report confidence intervals in parenthesis below the coefficient estimates in all the regression tables of this paper.

Table 2.6 reports the impact of applying for a SINERGIA grant. Columns 1 to 3 cover the regression results for the three standard bibliometric measures: Publication count in logarithmic terms (*Log(Publication count)*) in column 1, the average impact factor of the journals where the scientist publishes in logarithmic terms (*Log(Average IF)*) in column 2, and the average number of citations received per paper in logarithmic terms (*Log(Average citations)*) in column 3. Columns 4 and 5 report the results for the dummy co-applicants' collaboration realized (*Co-applicant collaboration*) and the individual learning of scientists in logarithmic terms (*Log (Learning)*), respectively.

We find that scientists who applied for a SINERGIA grant are more productive in quantitative terms than scientists who did not apply. Specifically, applicants publish, on average, 43% papers more than non-applicant in the five-year window following the application. We also observe that they increase the average impact factor of the journals where they publish by 7%. However, applicants' articles

---

<sup>17</sup> For the two discrete variables *Publication count* and *Learning* we also performed a Poisson estimation. The results of these estimations, available in appendix A2.4, are consistent with the results reported in Table 2.6.

receive, on average, 33% fewer yearly citations than those of the potential applicants. Applicants have a 19% higher probability of establishing a co-authorship with their co-applicants than potential applicants and learn more on average.

**Table 2.6. Regression results for the estimation of Equation 2.1 comparing applicants to potential applicants.**

	(1) OLS Log(Publication count)	(2) OLS Log(Average IF)	(3) OLS Log(Average cita- tions)	(4) Probit Co-applicant Collaboration	(5) OLS Log(Learning)
Applicant*Post-Application	0.43*** (0.36 ; 0.49)	0.070*** (0.018 ; 0.12)	-0.33*** (-0.39 ; -0.27)	0.19*** (0.12 ; 0.27)	0.36*** (0.29 ; 0.44)
Applicant	0.23*** (0.17 ; 0.29)	-0.030 (-0.095 ; 0.036)	-0.11*** (-0.19 ; -0.030)	0.43*** (0.37 ; 0.50)	0.043 (-0.019 ; 0.11)
Post-Application	-0.25*** (-0.30 ; -0.20)	-0.067*** (-0.11 ; -0.019)	-0.15*** (-0.21 ; -0.096)	-0.013 (-0.073 ; 0.046)	-0.62*** (-0.68 ; -0.57)
Seniority	0.031*** (0.027 ; 0.034)	-0.0019 (-0.0047 ; 0.00084)	-0.00053 (-0.0036 ; 0.0026)	-0.0014 (-0.0031 ; 0.00027)	0.012*** (0.0092 ; 0.015)
Other active funding	0.33*** (0.27 ; 0.40)	0.19*** (0.14 ; 0.24)	0.17*** (0.11 ; 0.23)	-0.0081 (-0.039 ; 0.022)	0.25*** (0.19 ; 0.30)
Previous expired funding	0.22*** (0.15 ; 0.29)	-0.026 (-0.077 ; 0.025)	-0.068** (-0.13 ; -0.0024)	0.030** (0.00098 ; 0.060)	0.098*** (0.041 ; 0.16)
Constant	1.76*** (1.60 ; 1.91)	1.17*** (1.04 ; 1.30)	0.99*** (0.84 ; 1.13)		3.37*** (3.22 ; 3.51)
Dummy Application year	Yes	Yes	Yes	Yes	Yes
Dummy Discipline	Yes	Yes	Yes	Yes	Yes
Dummy Affiliation	Yes	Yes	Yes	Yes	Yes
Appl./Potential appl.	2,120	2,120	2,120	2,120	2,120
Observations	4,240	4,240	4,240	4,240	4,240
R2 / Pseudo R2	0.435	0.261	0.243	0.39	0.466

The table shows a difference-in-differences estimation in the equivalent regression formulation, where the effect of the treatment, i.e., applying for a SINERGIA grant, can be read in the coefficient of the interaction variable *Applicant\*Post-Application*. The five scientific outcomes considered are *Publication count*, *Average IF*, *Average citations*, *Co-applicant collaboration*, and *Learning*. The sample includes 4,240 observations. The 1,060 Applicant-application pairs and the 1,060 Potential applicants, i.e., 2,120 Applicant-application pairs/Potential applicants-matched application pairs, are observed in two periods, namely before and after the application year. Columns 1, 2, 3, and 5 report OLS estimates, whereas Column 4 reports the marginal effect of a Probit estimation that considers the binary nature of the dependent variable *Co-applicant collaboration*. Robust standard errors are clustered around the application. In reporting the statistical significance of the coefficients, we apply the standard thresholds, namely \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ . 95% confidence intervals are reported in parenthesis below each coefficient.

Table 2.7 shows the results of the estimation of Equation 2.2, comparing awarded applicants to non-awarded applicants. On average, awarded applicants do not perform better than non-awarded ones regarding the quantity and quality of their scientific production and learning, but they have a higher chance of establishing a co-authorship with their co-applicants. Specifically, being awarded the funds increases the probability of co-authoring with co-applicants by 17% on average<sup>18</sup>.

<sup>18</sup> In our sample, each applicant can be involved in more than one call. To explore how this persistency in applying could affect our results, we conduct two distinct additional exercises. In the first one, we run a set of regressions limiting our study sample to the first grant call. In the second one, we adjust the control sample of awarded applicants by performing a conditional difference-in-differences estimation matching awarded applicants with similar non-awarded ones. The propensity score

**Table 2.7. Regression results for the estimation of Equation 2.2 comparing awarded to non-awarded scientists.**

	(1) OLS Log(Publication count)	(2) OLS Log(Average IF)	(3) OLS Log(Average citations)	(4) Probit Co-applicant collaboration	(5) OLS Log(Learning)
Awarded*Post-Application	0.034 (-0.050 ; 0.12)	0.027 (-0.024 ; 0.079)	0.086 (-0.020 ; 0.19)	0.17*** (0.086 ; 0.26)	-0.025 (-0.15 ; 0.10)
Awarded	-0.15** (-0.30 ; -0.0031)	-0.015 (-0.18 ; 0.15)	-0.11 (-0.31 ; 0.095)	-0.15** (-0.28 ; -0.0086)	-0.19** (-0.36 ; -0.025)
Post-Application	0.16*** (0.11 ; 0.22)	-0.0087 (-0.042 ; 0.025)	-0.53*** (-0.60 ; -0.46)	0.20*** (0.15 ; 0.26)	-0.25*** (-0.33 ; -0.18)
<b>Scientist's characteristics</b>					
Seniority	0.023*** (0.018 ; 0.028)	0.0039** (0.00039 ; 0.0075)	0.0058** (0.00081 ; 0.011)	-0.0032* (-0.0066 ; 0.000067)	0.012*** (0.0081 ; 0.016)
Female	-0.11** (-0.22 ; -0.012)	0.032 (-0.056 ; 0.12)	0.10* (-0.0032 ; 0.21)	0.0080 (-0.069 ; 0.086)	0.033 (-0.072 ; 0.14)
Other active funding	0.15*** (0.076 ; 0.23)	0.13*** (0.066 ; 0.20)	0.16*** (0.076 ; 0.24)	-0.024 (-0.083 ; 0.034)	0.10*** (0.028 ; 0.18)
Previous expired funding	0.15*** (0.069 ; 0.24)	-0.054 (-0.12 ; 0.012)	-0.083* (-0.17 ; 0.0056)	0.053* (-0.0082 ; 0.11)	0.048 (-0.026 ; 0.12)
<b>Application characteristics</b>					
Grade	0.040 (-0.0078 ; 0.088)	0.048* (-0.0023 ; 0.099)	0.068** (0.0066 ; 0.13)	0.042** (0.00098 ; 0.083)	0.071** (0.015 ; 0.13)
Swiss team	-0.048 (-0.19 ; 0.098)	-0.0071 (-0.13 ; 0.12)	-0.039 (-0.16 ; 0.082)	-0.0028 (-0.13 ; 0.13)	0.095 (-0.022 ; 0.21)
At least one female researcher	-0.047 (-0.13 ; 0.034)	-0.013 (-0.098 ; 0.071)	0.016 (-0.079 ; 0.11)	0.014 (-0.055 ; 0.084)	-0.011 (-0.11 ; 0.087)
Log(Amount Requested)	0.028 (-0.096 ; 0.15)	0.12** (0.000036 ; 0.24)	0.11 (-0.045 ; 0.27)	-0.097* (-0.20 ; 0.0011)	-0.013 (-0.16 ; 0.13)
Log(N. of co-applicants)	0.051 (-0.081 ; 0.18)	-0.037 (-0.19 ; 0.11)	0.015 (-0.17 ; 0.20)	0.30*** (0.17 ; 0.43)	0.035 (-0.16 ; 0.23)
Log(N. of disciplines)	-0.016 (-0.090 ; 0.057)	0.049 (-0.030 ; 0.13)	-0.0081 (-0.10 ; 0.087)	-0.021 (-0.074 ; 0.031)	0.084* (-0.016 ; 0.18)
Science & Medicine	-0.13** (-0.26 ; -0.0019)	0.41*** (0.28 ; 0.55)	0.42*** (0.26 ; 0.59)	-0.095 (-0.21 ; 0.025)	0.57*** (0.41 ; 0.73)
Log(1+Distance hours)	0.039 (-0.013 ; 0.091)	0.082** (0.020 ; 0.14)	0.063* (-0.0090 ; 0.14)	-0.023 (-0.074 ; 0.028)	0.054* (-0.0059 ; 0.11)
Previous SINERGIA application	0.13** (0.029 ; 0.24)	-0.064 (-0.16 ; 0.035)	-0.066 (-0.20 ; 0.070)	0.060 (-0.031 ; 0.15)	0.14** (0.033 ; 0.25)
Previous SINERGIA awarded	0.067 (-0.083 ; 0.22)	-0.025 (-0.18 ; 0.13)	-0.095 (-0.27 ; 0.080)	0.083 (-0.043 ; 0.21)	0.052 (-0.078 ; 0.18)
Constant	1.82** (0.083 ; 3.56)	-0.96 (-2.57 ; 0.65)	-1.22 (-3.47 ; 1.03)		3.18*** (1.25 ; 5.12)
Dummy Application year	Yes	Yes	Yes	Yes	Yes
Dummy Discipline	Yes	Yes	Yes	Yes	Yes
Dummy Affiliation	Yes	Yes	Yes	Yes	Yes
Applicant-application pairs	1,060	1,060	1,060	1,060	1,060
Observations	2,120	2,120	2,120	2,120	2,120
R2 / Pseudo R2	0.373	0.435	0.371	0.12	0.504

The table shows a difference-in-differences estimation in the equivalent regression formulation, where the effect of the treatment, i.e., being awarded a SINERGIA grant, can be read in the coefficient of the interaction variable *Awarded\*Post-Application*. The controls include fixed effects for the *Application year*, *Affiliation*, and *Discipline* of the scientist. The five scientific outcomes considered are the *Publication count*, *Average IF*, *Average citations*, *Co-applicant collaboration*, and *Learning*. The sample includes 1,060 Applicant-application pairs. The 1,060 pairs are observed before and after the treatment year, for a total of 2,120 observations. Columns 1, 2, 3, and 5 report OLS estimates, whereas Column 4 reports the marginal effect of a Probit estimation that considers the binary nature of the dependent variable *Co-applicant collaboration*. Robust standard errors are clustered around the application. In reporting the statistical significance of the coefficients, we apply the standard thresholds, namely \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ . 95% confidence intervals are reported in parenthesis below each coefficient.

matching performed is based on the scientist's application history as additional relevant matching criterion. Results for both exercises, available upon request, are consistent with the ones reported in Table 2.7.



## 2.6 Discussion

In a period of increased selectivity of public grants, scientists are concerned about the utility of spending energy and time in participating in grant competitions where the odds of getting awarded are low. In Section 2.5, we find that scientists who decide to apply for a grant and put the effort necessary to do so, perform differently from the ones who did not. Specifically, we find that participating in a research grant competition per se has a positive effect on the scientist's number of publications, the average impact factor of the journals where she publishes, her probability of collaborating with co-applicants, and her learning. However, applicants receive, on average, fewer citations per paper than non-applicants.

These results suggest that the sunk costs incurred by applicants to craft their application proposal are not fruitless. The application process is time-consuming and requires dedication in writing a valuable proposal. The efforts provided to put together the project and set the research agenda boost the level of advancement and the quality of applicants' research, hence positively stimulating the subsequent number of scientific publications and the average impact factor of the journals where they are published. When applying for grants, scientists design projects and create work ties with co-applicants. We find that they build on these ties afterward regardless of the result of the grant competition. Applicants are exposed to the knowledge of their co-applicants, and spillovers are likely to occur. Hence, the observed positive impact of applying on learning could be explained by the interactions with their co-applicants when crafting the project (Ayoubi et al. 2017).

Interestingly, we observe that applicants to SINERGIA receive fewer citations on average than potential applicants. This observed decrease in citations could be the result of a decline in the quality of the scientific production or a loss of visibility due to entry to new fields of research when scientists are involved in interdisciplinary projects (Azoulay et al. 2014). The observed increase in the average impact factor seems to discard the first interpretation, i.e., a decline in quality. The decline in citations appears hence driven by an entry to research fields with which applicants are not familiar and where they need time to establish a reputation. Entering new research fields, applicants need to acquire new knowledge. In Table 2.6 (Column 5), we find that applying stimulates the tendency of scientists to learn. The lack of visibility hypothesis would then be confirmed if we observe that the knowledge acquired shows a move towards new fields of research.

To give empirical ground to the hypothesis of the decreased visibility when entering new fields, we implement a test to evaluate how different is the newly acquired knowledge compared to the initial stock of knowledge of the applying scientists. Building on Ayoubi et al. (2017), we measure this difference using the cognitive distance separating the scientist's knowledge stock before the application and

the knowledge acquired after the application. Specifically, we estimate the average cognitive distance between the newly acquired knowledge and the original knowledge stock of the scientist by using a measure based on the list of journals the scientist references in her scientific publications. We consider these journals as a proxy for a scientist's knowledge mobilized in her research<sup>19</sup>.

Table 2.8 shows that the knowledge newly acquired by applicants is more diverse in terms of subjects than for potential applicants. In other terms, applying for a SINERGIA grant seems to encourage scientists to enter new fields of knowledge. Interestingly, we find no significant effect of the variable *Awarded* on the entry to new fields, which means that scientists being awarded with funds did not publish research that is more diverse in terms of subjects than non-awarded ones. Consequently, we can assert that SINERGIA awardees are not selected into performing more interdisciplinary research. Two possible explanations can be driving our results. On the one hand, the SINERGIA selection committee could tend to award scientists who are in relatively close disciplines instead of promoting proposals that are more risk-taking and with a higher diversity of disciplines. Recent evidence in research evaluation systems and other formal appraisals of science showing a tendency to favor research within the same field (Rafols et al. 2012; Chavarro et al. 2014) would support this interpretation. On the other hand, this result also suggests that the SNSF is driving the interdisciplinarity of research mainly through an incentive mechanism. The mere act of writing a proposal for an interdisciplinary project seems to be a strong enough incentive to carry on research with a greater diversity of subjects.

---

<sup>19</sup> For a more in-depth discussion of the process and the tools used, see Appendix A2.3.

**Table 2.8: Regression testing entry to new fields of knowledge.**

	OLS Log(1+Journal distance)
Applicant	0.22*** (0.14 ; 0.30)
Awarded	-0.065 (-0.16 ; 0.034)
Seniority	0.0083*** (0.0050 ; 0.012)
Other active funding	0.0047 (-0.082 ; 0.092)
Previous expired fund- ing	0.14*** (0.070 ; 0.21)
Constant	6.17*** (5.96 ; 6.38)
Dummy Application year	Yes
Dummy Discipline	Yes
Dummy Affiliation	Yes
Appl./Potential appl.	2,120
R2	0.224

The table shows the estimated effects of applying (*Applicant*) and being awarded a grant (*Awarded*) on the scientist's entry to new fields of knowledge. The entry to new fields of knowledge is proxied by the average cognitive distance between the newly acquired knowledge and the original knowledge stock of the scientist relying on a measure based on the list of journals cited in the scientist's articles (*Journal distance*). We observe 2,120 Applicant-application pairs/Potential applicants-matched application pairs, of which 1,060 are Applicant-application pairs, and 1,060 are Potential applicants. Robust standard errors are clustered around the application. In reporting the statistical significance of the coefficients, we apply the standard thresholds, namely \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ . 95% confidence intervals are reported in parenthesis below each coefficient.

Having assessed the value of applying for a competitive grant for scientists, we consider the impact of the funding decision for applicants. Doing so, we contribute to the scarce literature on the effect of receiving funds on researchers' scientific outcomes. We find that being awarded has no significant effect on the scientific productivity of grant recipients, regarding the quantity or quality of papers published. However, considering the values of the estimated coefficients and their respective confidence intervals in Table 2.7, we observe a limited but positive impact of funding on publication productivity. As a matter of fact, the first line of Table 2.7 shows that the effects of the interaction *Awarded\*Post-Application* on *Publication count*, *Average IF*, and *Average citations* are not significant because their p-values are greater than the standard significance levels, i.e., 1%, 5%, and 10%. However, these relationships are likely to be significant and positive because the range of their confidence intervals lies mostly in the strictly positive domain, and their coefficients, albeit small, are also positive<sup>20</sup>. This observation is consistent with the literature on the impact of funding on research, finding a little but positive effect of

<sup>20</sup> This tendency of the coefficients of *Average IF* and *Average citations* to be positive is confirmed in the two exercises we performed in the evaluation of the impact of being awarded (see footnote 18). Precisely, we find that the coefficient of *Average citations* becomes significant at the 5% level when we limit the sample to first applications to SINERGIA. As for the second exercise, in the conditional difference-in-differences estimation using propensity score matching, we find that the coefficient of *Average IF* becomes significant at the 1% level.

funding on standard bibliometric outcomes (Arora and Gambardella, 2005; Jacob and Lefgren, 2011; Gush et al., 2015; Azoulay et al., 2015; Carayol and Lanoe, 2017). However, concerning our additional measures, we observe that co-applicants of an awarded project have higher chances to co-author with each other than co-applicants of a non-awarded project. In other words, financial resources have no impact on the individual publication productivity but provide an incentive to consolidate the collaboration started in the application phase. We might explain these results on productivity by the fact that scientists are not strictly dependent on public grants to sponsor their research. They might have access to private funding or use their public bulk funding. However, if a grant is released, scientists are incentivized to finalize the collaborations planned when they submitted their projects.

## 2.7 Conclusion

In this paper, we find that scientists participating in a research grant competition reap the benefits of the efforts spent, even if they do not obtain the desired reward. Applicants for a SINERGIA grant increase their scientific productivity and learn while exploring new fields of knowledge. These positive externalities lead us to claim that “the important thing is not to win; it is to take part.”

This noteworthy result contributes to the debate in the scientific community concerning the utility for researchers to spend time writing proposals to raise money for their research. Our results suggest that scientists derive some benefits from the time spent writing proposals. Hence, scientists should be less reluctant to invest time and effort entering research grant competition since these could be the occasion to launch new trends of research, build working ties with fellow researchers and acquire new knowledge. On the side of funding agencies, our results imply that promoting the calls for proposals as well as encouraging scientists to apply could be as efficient as increasing the funds dedicated to financing research projects.

During the exploratory phase of our work, discussing directly with the SNSF and with actual applicants for SINERGIA, we noted that a peculiarity of the SINERGIA grant is the fact that the administrative requirements of the submitted proposals are limited compared to other grants. In other words, most of the work scientists do when applying for the grant is directly related to the scientific project they are crafting and could, therefore, be useful in further research regardless of the result of the competition. This characteristic of SINERGIA could be partly responsible for our findings, and it, more importantly, suggests that other funding agencies could follow a similar pattern in designing their calls for grant proposals.

Research grants are a relevant setting for analyzing the impact of competition on participating agents. We found that awarded researchers are not the only beneficiaries in a research grant competition since the mere act of being involved in the process is beneficial in many aspects. This finding might be extended to other competitive contexts such as start-ups running for venture capital funding or firms applying for calls of public procurement contracts. When looking for investors to sustain their business, entrepreneurs are asked to write demanding business plans and to have precise strategic planning. Hence, the efforts spent in performing these requirements could be useful to improve the business performance of the start-up regardless of the result of the funding decision of investors. Similarly, firms could boost the efficiency of their projects when working on meeting the requirements of a public procurement call and benefit whether they succeed in winning the contract or not.

To the best of our knowledge, this study is the first to consider the effect of applying on the scientific outcomes of researchers. Our results proved to be robust across several estimation strategies suggesting a treatment effect of the application process. However, one might still claim that our estimation strategy does not entirely account for the selection bias. We believe that this concern could justify a reduction of the magnitude of our results but not their direction. Moreover, we are aware that SINERGIA is a peculiar grant which might attract highly experienced scientists with interdisciplinary profiles, thus limiting the scope of our findings. However, this paper intends to open the way for a better consideration of the effect of participating in competitions, and future studies could further investigate the mechanisms that lead applicants and non-applicants with similar profiles, to perform differently.

# Chapter 3 At the origins of learning: Absorbing knowledge flows from within the team

*Disclaimer: This chapter, written in collaboration with Michele Pezzoni and Fabiana Visentin, is now published in the Journal of Economic Behavior and Organization, DOI: [10.1016/j.jebo.2016.12.020](https://doi.org/10.1016/j.jebo.2016.12.020) (Accepted December 22, 2016).*

## **Abstract**

Empirical studies document a positive effect of collaboration on team productivity. However, little has been done to assess how knowledge flows among team members. Our study addresses this issue by exploring unique rich data on a Swiss funding program promoting research team collaboration. We find that being involved in an established collaboration and team size foster the probability of an individual learning from the other team members. We also find that team members with limited experience are more likely to learn from experienced peers. Moreover, there is an inverted U-shaped effect of cognitive distance on the probability of learning from other team members.

## 3.1 Introduction

This paper assesses the characteristics of a research team that foster the probability of learning from one another. We add to the Science of Team Science (SciTS) literature by investigating research teams from the point of view of team members' learning, an aspect often neglected in favor of team productivity analysis.

Over the past century, the process of scientific knowledge production has fundamentally changed. Nowadays, the teamwork model of conducting science has mainly replaced the single scientist model (Jones et al., 2008; Wuchty et al., 2007). Several reasons explain this trend. First, the cost of scientific instrumentation leads scientists to organize in teams in order to share resources and to avoid cost duplication. Second, lower travel and communication costs increase scientists' mobility and favor the creation of multi-institution teams. Third, certain fields such as physics, chemistry, engineering, and biology are characterized by an increasing level of complexity, which requires the joint effort of specialized scientists. It becomes implausible for a single individual to master all of the technical skills and

knowledge needed to set up a laboratory, run an experiment, analyze the data, and manage the publication process.

There is a general consensus among scholars that a “collaboration [outcome] is greater than the sum of its parts” (Katz and Martin, 1997). Even if some authors present some drawbacks to the collaboration, such as higher coordination costs (Bikard et al., 2015; Mowatt et al., 2002), or ghost and honorary authorships (Mowatt et al., 2002), most empirical studies agree that collaboration has a positive impact on publication productivity. Not only does teamwork have a greater value than solo-author work, but teamwork also positively affects the productivity of each team member (Defazio et al., 2009; Lee and Bozeman, 2005). The most common explanation of the greater value and higher productivity of teamwork is that it allows scientists to combine their knowledge, prompting scientific discoveries (Uzzi et al., 2013).

These studies investigating the dynamics of scientific teams are the building blocks of the SciTS literature (Börner et al., 2010; Whitfield, 2008; Stokols et al., 2008). As described by Börner et al. (2010), SciTS is “an emerging area of research centered on the examination of the processes by which scientific teams organize, communicate, and conduct research.” Our study aims to shed light on a process often neglected in this literature: the exchange of knowledge among members of a scientific research team and their ability to learn from one another.

While the process of learning has been investigated within the organizational literature, starting from the ‘80s with the work of Levitt and March (1988), this process remains largely unexplored in the SciTS literature. The organizational literature claims that learning new skills and using them within a firm is critical to the innovation and productivity of the firm (Argote and Miron-Spektor, 2011). Similarly, we believe that, for scientists, acquiring new knowledge and exploiting it in their research work is key to their productivity and for the novelty of their contribution. Precisely, it has been shown that broadening the researcher horizon and exchanging knowledge within the framework of interdisciplinary teams is key to the production of innovation and high-impact scientific research (Börner et al., 2010). More recently, a study by Misra et al. (2015) suggests that scientists who are more open to other disciplines are more likely to produce higher-quality research. The process of acquiring new diverse knowledge is, therefore, central to the effectiveness of a scientific team.

We identify the factors that promote the learning of an individual from her teammates. We claim that this portion of learning can be affected by the characteristics of the team the individual is working in. As team characteristics, we consider the quality of the research project that the team members are working on, the team size, and the discipline. We also distinguish co-ethnic teams from multi-ethnic teams, and teams with at least one female scientist from only male teams. Moreover, we consider as a

determinant of learning the distance between an individual and her team members along three dimensions, namely, geographical distance, social distance and cognitive distance (Agrawal et al., 2008, 2003; Bercovitz and Feldman, 2011; Jaffe et al., 1993; Nooteboom et al., 2007).

In our analysis, we use a unique dataset of 255 grant applications to a Swiss funding program promoting team collaboration. This rich dataset allows us to clearly define a team, its boundaries, and its creation date. We then use the bibliographical references of each scientist in the dataset to precisely define her knowledge stock and specify her learning from the other team members.

We find that learning from team members is more likely within larger teams. Team members with a limited stock of knowledge are more likely to learn from more experienced team members. Also, having an already established collaboration is correlated with a higher probability of learning from the rest of the team. We compare the knowledge capital shock content of a scientist with her teammates' knowledge capital stocks to measure cognitive distance. We find an inverted U-shaped impact of the computed cognitive distance between the scientist and her team on the probability that learning originates from within the team. An individual with a knowledge stock differing from that of the others guarantees a buffer for learning something new. At the same time, the difference in the knowledge stocks should not be too large so as to avoid obstacles to effective communication between team members, that is, the situation when team members speak different languages and do not understand one another.

The rest of this paper is organized as follows: Section 3.2 illustrates the individual learning determinants. Section 3.3 describes the data and variables. Section 3.4 describes the estimation strategy. Section 3.5 provides the results, and Section 3.6 concludes.

## 3.2 Individual learning: Definition and determinants

As defined by Salomon and Perkins (1998), paraphrasing Huber (1991), "Learning consists in knowledge [...] acquired by any unit of an organization and available for acting upon" (p. 13). Applying this definition of learning in our empirical setting represented by scientists working in research teams, we proceed in two steps. First, we assess the knowledge accumulated by a scientist before entering a team, based on the literature she relied upon in her work. Then, we consider as learning any increment to this initial stock of knowledge. Finally, we focus on the part of this learning that originates from within the scientific team, i.e., the knowledge transmitted by her teammates.

To identify the determinants of this part of learning originating from within the team, we rely on the more extensive literature of SciTS on team productivity. In the current section, we discuss the factors influencing the probability of learning from other team members for an individual scientist. Specifically,



we consider how the probability of an individual to learn from other team members is affected, on the one hand, by the characteristics of the team, and on the other hand, by the individual characteristics of the scientist in comparison to the rest of her team.

### 3.2.1 Team characteristics

As team characteristics, we consider: being endowed with research funds, working on a high- quality project, having various sizes, and having a different ethnic and gender composition.

The creation of collaborative relationships might be facilitated in teams whose projects have been awarded due to the availability of funds for traveling, team-building activities, workshops, and meetings. Moreover, funds might be used to buy research equipment and materials shared among team members. Collaboration activities favored by the availability of funds are expected to foster knowledge flows among team members, and consequently, individual learning.

We distinguish teams with high-quality research projects from teams with lower-quality research projects. High-quality research projects and the promise of making breakthrough scientific discoveries might stimulate scientists' commitment to working actively together. This might foster knowledge flows and learning among team members.

We expect team size to be positively correlated with the probability of learning from other team members: a greater number of individuals with whom to interact should increase the probability of learning.

Several works have investigated the effect of researchers' co-ethnicity on the probability of knowledge flows (Agrawal et al., 2008, 2003; Freeman and Huang, 2015). The prevalent result in the literature is a positive effect of the co-ethnicity of researchers on the probability of observing a knowledge flow. Following the same line of reasoning, we expect co-ethnic teams to favor knowledge flows among team members, and consequently, individual learning.

A large part of the gender literature focuses on the effect of team gender composition on team productivity (Apestegua et al., 2012; Pezzoni et al. 2016). Woolley et al. (2010) investigated the mechanisms behind the ability to accomplish a task within a group where individuals of different genders collaborate. They find that the presence of females improves team performance, and they attribute this productivity premium to the fact that females have higher social sensitivity, i.e., the ability to understand the mental state of another person. Therefore, the presence of women on a team might favor personal interactions among team members, resulting in a positive effect on team learning.

### 3.2.2 Scientist vs. Rest of the team

In order to compare the scientist to the rest of her team, we consider how different she is from her teammates in terms of geographical localization, social attributes, and cognitive diversity. We define social attributes as all of the individual characteristics that determine the stratification of the research community. These attributes might be exogenous, such as gender and age of the scientist, or acquired, such as the scientist's scientific reputation or her collaboration patterns. In contrast to social attributes, cognitive diversity with respect to scientists concerns only the aspect of knowledge differences among individuals, other things being equal (McPherson et al. 2001).

We categorize geographical localization, social attributes, and cognitive diversity into three types of distances separating the scientist from the other members of the team: geographical distance, social distance, and cognitive distance.

Over the last thirty years or so, in parallel to the increase in the average team size (Wutchy et al., 2007), we witness an even higher increase in the geographical dispersion of the team members (Jones et al. 2008). Adams et al. (2005) show that the average geographical distance of collaborations more than doubled in the last twenty years due to improvements in transport and telecommunications. In a sample of French scientists, Mairesse and Turner (2005) show that, except for immediate proximity (i.e., being affiliated with the same unity), geographical distance has no significant impact on collaboration. According to this evidence, we expect a limited effect of geographical distance among team members on the probability of learning from one another.

In our study, we consider four variables measuring the social distance of the individual from her team that might affect her probability of learning: age, reputation, gender, and previous collaborations. First, regarding age, in a mentor-protégé relationship, the young team member is expected to learn, i.e., receive knowledge from the senior team member, who is expected to transmit knowledge (Campbell and Campbell, 1997). Then, we might expect the age difference between an individual and her teammates to be positively correlated with her attitude to engage in knowledge transmission activities. However, Zenger and Lawrence (1989) find that, in a firm environment, individuals of a similar age tend to exchange information more easily. These two competing effects prevent us from formulating a prediction on age difference effects.

Second, we consider the scientific reputation of team members, as proxied by their publication productivity before the team formation. We identify two possible mechanisms at play within the team. On one side, highly productive members might contribute to the team with larger knowledge stocks and might enhance the probability of learning for less productive team members. On the other side, highly

productive scientists might focus on knowledge exchanges with teammates having similar publication stocks from which they can benefit more; thus, they may decide to isolate low-productive scientists, from which they benefit less. As in the case of age differences, we have two competing hypotheses about the possible effect of scientific team reputation diversity on learning from within the team.

Third, we look at the presence of individuals with the same gender in the team. By relying on the concept of “homophily,” we expect that team members of the same gender would be more likely to benefit from reciprocal knowledge flows and learning (McPherson et al., 2001; Cummings and Kiesler, 2008).

Finally, the mechanisms that affect the learning of the scientist from her team might differ as to whether the individual has long-lasting collaborations with her teammates or not. On the one hand, scientists having previous collaborations with their team have a greater level of familiarity (Cummings and Kiesler, 2008; Bercovitz and Feldman, 2011), and thus benefit from the presence of routinized collaboration activities that facilitate the creation of strong relational ties (Porac et al., 2004). Strong relational ties foster knowledge flows among team members and enhance their probability of learning from one another (Granovetter, 1973). On the other hand, having previous professional collaborations might increase the probability that team members share the same knowledge stock. Redundant knowledge decreases the probability that individuals learn from one another (Burt, 2004). The contrasting effects of the mechanisms at work when the scientist has an established collaboration prevent us from making predictions of the impact on her probability of learning.

In the management literature, a major determinant of the knowledge flows and innovative performance of the team is the cognitive distance separating its members (Knoben and Oerlemans, 2006; Nootboom et al., 2007). In the SciTS literature, the question of the disciplinary diversity of the scientific team is a central issue (Fiore, 2008). The tenets of SciTS discuss the impact of the various types of cross-disciplinary teams on the effectiveness of the collaboration (Stokols et al., 2008; Börner et al., 2010). In order to estimate the level of disciplinary diversity in the team, we use a proxy for the cognitive distance between team members, measuring the distance separating their knowledge stocks before the team formation. According to Nootboom et al. (2007), the cognitive distance between the team members has two competing effects on the knowledge production capacity of a team in an organization. On the one hand, the capacity for absorbing new knowledge is higher when the cognitive distance between the members is low since it is easier for the scientist to absorb knowledge similar to what she already has. Hence, the “*speaking the same language*” effect enhances knowledge flows within the team when the cognitive distance is low. On the other hand, having low cognitive distance between individuals implies that their knowledge stocks are very similar, and thus, the probability of observing a knowledge flow

from other team members is low because they have too little novelty to offer to the scientist. Thus, the so-called “*opening new horizons*” effect has a positive impact on knowledge flows within the team as the cognitive distance increases. Therefore, combining the two effects, we expect the global impact of cognitive distance to have an inverted U-shape on knowledge flows within the team. For low cognitive distances, even if the absorptive capacity is very high, the low diversity of knowledge in the team implies that the knowledge flows within the team remain very limited. The probability of having knowledge flows within the team increases when cognitive distance increases until some optimal point. Then, excessively high cognitive distance blocks the understanding between the individuals and negatively affects the knowledge flows between team members.

### 3.3 Data

#### 3.3.1 Team

We rely on the definition of a team as a group of individuals working together for a limited period of time to pursue a circumstantial goal (Katz and Matrin, 1997). Empirically, in the SciTS literature, teams are often reconstructed through co-authorship relationships (Ding et al., 2010; Wuchty et al., 2007). In our study, we refrain from basing our team definition on publication data. Following Cummings and Keisler’s (2008) definition of a team based on grant applications, we consider a team as a group formed by all the scientists who express their willingness to collaborate by submitting a joint grant application. This definition has three main advantages with respect to the one based on co-authorship relationships. First, it fits the definition of a team as a group of individuals working together to achieve a common goal (Katz and Martin, 1997). The members of the team are the scientists who have their names on the grant application, and the goal of the team is explicitly stated on the grant application. Second, contrary to the common definition of a team based on co-authorship, this definition with clear boundaries allows us to capture even teams not producing any publication and the members of the team who are not mentioned in an eventual publication outcome. Finally, we are able to determine the precise time when the team is formed, independently of the time of the first team outcome, i.e., when the first co-authored article is published.

#### 3.3.2 Learning

At the team formation time, each individual is endowed with a knowledge stock represented by the literature she relied upon in her research work. We follow Uzzi et al. (2013) and proxy the knowledge component used by each scientist as the distinct journals cited in her work. Specifically, we use the list

of distinct scientific journals cited in the papers a scientist published before entering the team. Then, we measure the scientist's learning as the citations to new journals added to her knowledge stock after the team formation.

The learning may or may not be attributed to the interaction with other team members. We consider learning from within the team if the new journal citation observed was present in the knowledge stock of another team member before the team formation. If the new citation cannot be attributed to a knowledge flow from another team member, we classify it as not originating from within the team. It could originate from an outside collaboration, or it could be the result of a self-learning process.

### 3.3.3 Empirical setting

Our study is conducted in the context of the SINERGIA Swiss funding program. The program is sponsored by the Swiss National Science Foundation (SNSF), which is the leading Swiss institution supporting national scientific research. It plays the same role in Switzerland as the National Science Foundation (NSF) in the United States. SINERGIA was launched in 2008 and represents a flagship in the SNSF's funding schemes portfolio. It is designed to promote team collaboration. As mentioned in the application guidelines, researchers are required to collaborate as a condition of securing research funding, i.e., researchers need to submit a proposal for a "research work carried out collaboratively" (SNSF, 2011).

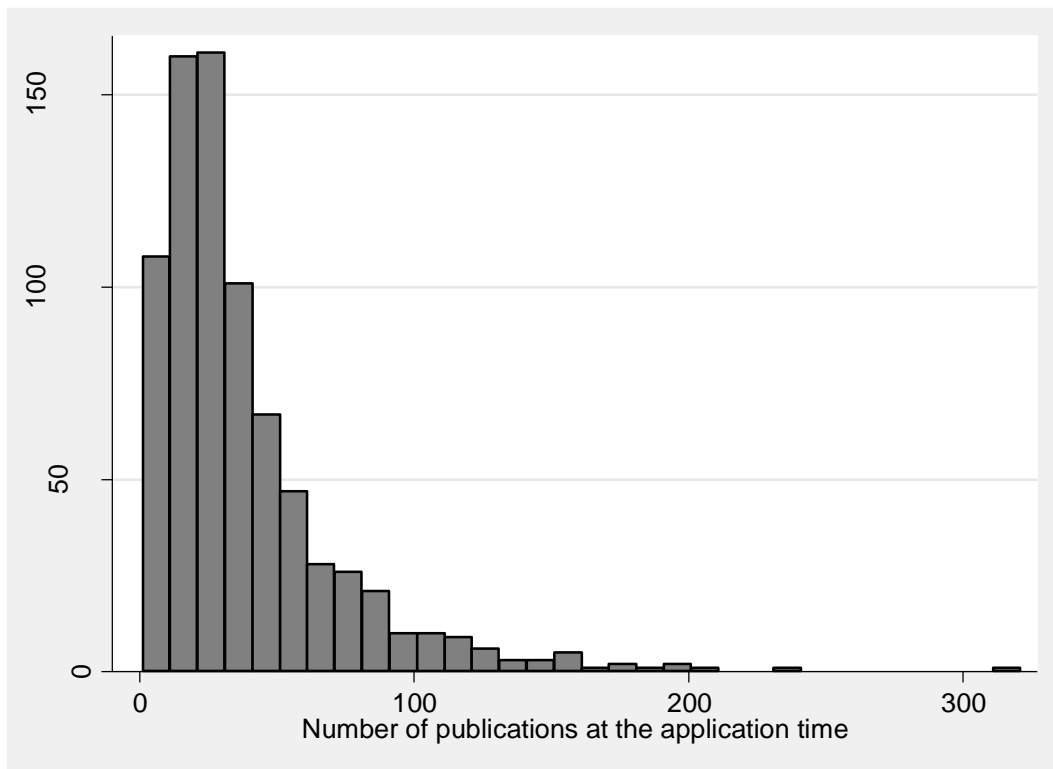
In most cases, a SINERGIA project involves three or four researchers who appear as co-applicants in the grant application. All disciplines are eligible for funding through the program. Applicants propose interdisciplinary projects or projects where co-applicants belong to the same discipline, but are specialized in different sub-fields. The criteria considered in evaluating the application are the value added to the joint research approach, the research complementarities of the applying groups, and the coherence of the projected collaboration. Applications are screened in a two-step evaluation process. In the first step, external reviewers assign a provisional score to each application. In the second step, an internal committee of SNSF, the Specialized Committee for Interdisciplinary Research based in the Swiss capital city Bern, assigns the final score to each application using an alphabetical scale, where A is the highest score, and D is the lowest one. Applications are ranked, and funds are assigned until the annual budget quota is reached. Typically, applications receiving a score below B are not funded.

From all grant applications submitted to the SNSF in the period 2008-2012, we selected applications in Engineering and Science & Medicine. Our final sample is represented by 255 grant applications, which

include 780 unique applicants<sup>21</sup>. The SNSF provided us with grant application data, including final scores assigned and final funding decisions and basic demographic information on applicants (gender, nationality, and birth year)<sup>22</sup>. We matched this information with applicants' publication records using the Scopus database<sup>23</sup>.

The SINERGIA funding program is aimed at established researchers. In the majority of cases, applicants are associate or full professors with good publication records. They have to demonstrate their ability to conduct excellent quality independent research. The average age of an applicant is 47 years old, with a minimum of 30 and a maximum of 69 years old. Figure 3.1 shows the distribution of the count of applicants' publications at the application time. The average number of applicants' publications is 38.

**Figure 3.1: Distribution of the number of scientists' publications at the time of the grant application**



---

<sup>21</sup> In this study, we excluded from the original sample applications in the Humanities and Social Sciences because book contributions represent a large part of the field publication outcomes and are not collected with accuracy in the Scopus database. Applications in the Humanities and Social Sciences represent 19% of the total initial sample.

<sup>22</sup> All concerned applicants were contacted by the SNSF and had the possibility to oppose the transmission of their data.

<sup>23</sup> We match the applicant's surname and the first letter of the name with the author's surname and the first letter of the name. Then, we filter the correct matches by hand-checking the scientist's identity, according to the applicant's characteristics, such as her affiliation, discipline, colleagues' names, and age.

The representative team in our sample is a small one. Ninety percent of the teams have less than six members. A team is composed, on average, of 4 members, with a minimum of 2 and a maximum of 11. Approximately 13% of the teams have only Swiss members, while the others are multi-nationality teams. The average number of nationalities in a team is 2.6, with a maximum of 7 nationalities. The SINERGIA funding program favors inter-institution collaborations. On average, each group has members from 2.8 different affiliations, with a maximum of 6. According to the SNSF's application requirements, a researcher with a foreign affiliation is admitted to apply for the grant only if her competencies and skills are not available in Switzerland. Due to this constraint, when we look at the country affiliations, we note that 66% of the teams include only Swiss affiliations. When classified by discipline, 36% of the applications are in Engineering, whereas the rest are in Science & Medicine. Within the two broad disciplines, each application is classified into sub-disciplines. An application counts, on average, 3.3 sub-disciplines; only 21% of the applications involve only one sub-discipline, while the most diversified application involves 11 sub-disciplines. When we look at the previous collaborations among applicants at the application time, we observe that in 58% of cases, there was at least one co-authorship relationship among the team members. When looking at the applicants' gender distribution, in our sample, women constitute 15% of the total. A SINERGIA grant covers personnel costs, research costs, coordination costs, and, to a limited extent, investment costs. The average amount requested per application is 1,674,320 CHF, with a minimum of 349,901 CHF and a maximum of 6,854,573 CHF.

Figure 3.2 represents the distribution of the number of grant applications by the score assigned and the final funding decision. A total of 8.6% of the applications obtained the maximum score, A, and 45% of the applications were awarded.

**Figure 3.2: Distribution of grant applications by score assigned and final funding decision**

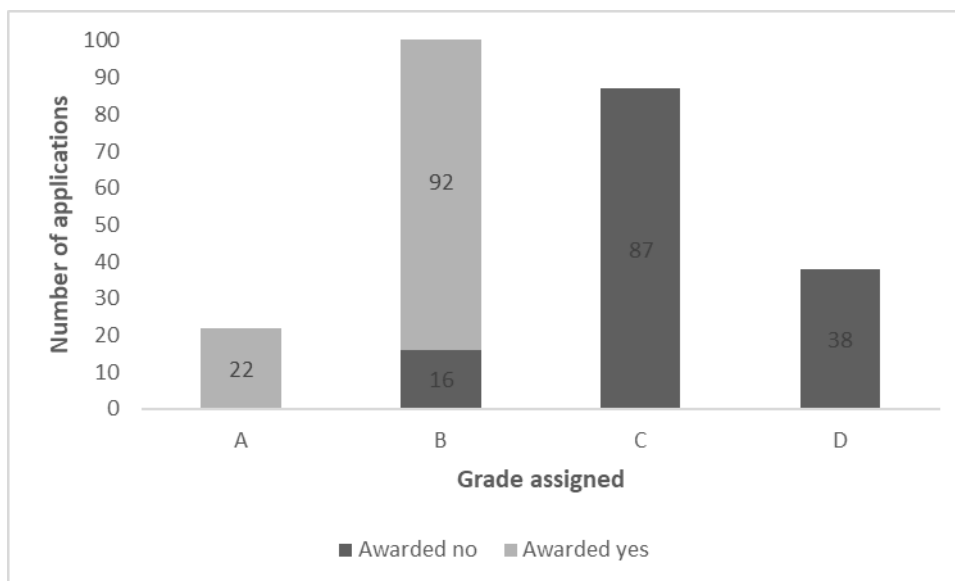


Table 3.1 reports the applicants' characteristics, and Table 3.2 reports the team characteristics.

**Table 3.1: Scientists' characteristics at the time of application (number of applicants=780)**

	Mean	Std. Dev.	Min	Max
Age	47.44	8.07	30	69
Gender (=1 for female, 0 otherwise)	0.16	0.36	0	1
Stock of publications pre-team entry	37.58	34.29	1	318
Stock of journals cited pre-team entry	135.62	102.24	1	644

**Table 3.2: Team characteristics at the time of application (number of teams=255)**

	Mean	Std. Dev	Min	Max
Number of team members	4.19	1.59	2.00	11.00
Number of nationalities represented	2.64	1.08	1.00	7.00
Number of country affiliations	1.38	0.55	1.00	3.00
Number of affiliations	2.79	1.03	1.00	6.00
Number of disciplines	3.30	2.16	1.00	11.00
At least 2 team members already co-authors	0.58	0.49	0.00	1.00
Average team members' age	47.74	4.93	35.09	59.97
Share of women	0.15	0.21	0.00	1.00
Average team members' stock of pubs	43.18	24.75	2.84	153.65
Awarded	0.45	0.50	0.00	1.00
High-quality application (grade A)	0.09	0.28	0.00	1.00
Low-quality application (grade D)	0.15	0.36	0.00	1.00
Amount requested	1674320	764260	349901	6854573
Science & Medicine	0.64	0.48	0.00	1.00
Engineering	0.36	0.48	0.00	1.00



### 3.3.4 Variables

In this section, we describe the dependent variables, the independent variables, and the controls used in the regression exercise. Our independent variables are grouped into two categories: team characteristics and individual characteristics of the scientist in comparison to the rest of her team. For the latter, we consider the three dimensions of distance, namely, geographical distance, social distance, and cognitive distance. Finally, our set of control variables includes individual and journal characteristics.

#### **Dependent variables**

Following Agrawal et al. (2008), we adopt as the unit of analysis the scientist-journal cited pair. For instance, a team composed of four scientists, each of whom cites ten distinct journals, generates forty observations. The scientist-journal cited level of analysis allows us to study the micro-dynamics of the team members' learning processes by isolating each knowledge component and tracing its origins (Börner et al., 2010). For each scientist, we consider two time periods: before and after the team formation. We compare the knowledge stock of each individual in the two periods in order to measure individual learning, namely the new journals cited that appear after the individual enters the team. Hence, we compute our dependent variable using all of the publications available from the constitution of a team until 2015, and the knowledge stocks of the scientists are constructed using all of the publications available before the date of the team constitution<sup>24</sup>.

The dependent variable *Learning from within the team* is a dummy that equals one if the new journal cited by the scientist after team formation is included in the stock of knowledge of at least one other team member before team formation, and zero otherwise. In our study sample, when looking at the origins of learning, 40% of the new journals cited by a team member after team formation originates from within the team. One possible concern in the definition of learning from within the team is that in the specific case of publications co-authored by two members of the same team, it is not possible to disentangle the contribution of each author to the list of references. In other words, we cannot exclude with certainty that the new journal cited by the scientists has been included in the list of references by her co-author teammate and that she actually has not learned from her teammate. In order to avoid this

---

<sup>24</sup> In the main analysis, we do not impose any constraints on the timespan, before or after the time of application. In a robustness check, available upon request, we fix a time window of three years, before and after the time of application, and the results remain stable.

overestimation of individual learning, we define a second variable *Learning from within the team– no co-authored pubs*, for which we exclude the newly cited journals in co-authored articles.

### Team characteristics

We distinguish co-ethnic teams from multi-ethnic teams with a dummy *Co-ethnic team*, which equals one when all of the team members are from the same country of origin and zero otherwise. We use the dummy *At least one female scientist in the team*, which equals one if at least one team member is a female scientist and zero otherwise, to differentiate mixed teams from male-only teams. We take into consideration the quality level of the research project of the team, based on the score assigned by the SNSF evaluation committee. We include in the regression a dummy *High-quality application (grade A)*, which equals one if the application obtains the maximum score, and zero otherwise, and a dummy *Low-quality project (grade D)*, which equals one if the application obtains the minimum score, and zero otherwise. The dummy *Awarded* concerns the final funding decision and is equal to one if the SNSF awards the team of the scientist, and zero otherwise. We use the variables *Amount requested* and *Number of team members* as proxies for the size of the team's project. Our sample includes teams working in two macro-fields: Engineering and Science & Medicine. The dummy *Science & Medicine* is equal to one for Science & Medicine, and zero otherwise. Each team can submit an application that involves one or more sub-fields. Finally, we take into account the number of sub-disciplines listed on the grant application with the variable *Number of disciplines*.

### Scientist vs. Rest of the team

We measure the geographical distance as the average time needed to travel from the affiliation of the scientist to the affiliations of the other team members (*Distance hours*)<sup>25</sup>. The average time needed to reach the other team members for the 780 individuals is 3.4 hours, with a standard deviation of 4 hours.

We measure the social distance over four dimensions. First, the age difference between the scientist and the rest of her team is calculated as the arithmetic difference between the age of the scientist and the average age of her teammates. The average difference is 0.06 years, with a standard deviation of 8.5 years. Then, we standardize the arithmetic difference by subtracting its average and dividing by its standard deviation (*Standardized age difference scientist vs. team*). Second, in a similar fashion, we

---

<sup>25</sup> Another possible variable for measuring the geographical distance could have been the distance in kilometers but, as expected, the time needed to travel and the distance in kilometers are highly correlated (about 0.9); therefore, we include only the connection time variable in the regression.

measure the productivity difference between the individual and the rest of her team, based on the number of publications (*Standardized stock of pubs. Difference scientist vs. team*). The average productivity difference is 6.8 scientific publications, with a standard deviation of 39.3 scientific publications. Third, we introduce the variable *Same-gender scientist vs. team*, which equals one if the team includes at least another team member of the same gender as the scientist, and zero otherwise. This dummy equals one for 92% of the 780 individuals considered. Finally, we consider the *Established collaboration* dummy, which equals one if the scientist has already worked with at least another member of the team in previous joint research projects, and zero otherwise. Previous joint research projects are identified by co-authored scientific articles before the year of the team formation.

We calculate the cognitive distance between the scientist and her team in two steps. First, we extract the references of the publications of all the scientists in our dataset, and we calculate one unique journal distance matrix (D), i.e., a matrix where each cell reports the distance between journal  $i$  and journal  $j$ . This journal distance matrix relies on the references of all articles included in our database, i.e., those published by the 780 scientists included in the analysis. The matrix is based on the assumption that the more the two journals  $i$  and  $j$  are cited together in the references of the same scientific publication, the closer they are. If the two journals are frequently co-cited, we attribute a small distance value to the pair. On the contrary, if the two journals are rarely co-cited, we attribute a large distance value to the pair. In the second step, we use the journal distance matrix (D) to calculate the average cognitive distance between the scientist (S) and her team (T).

In detail, each cell  $D(i, j)$  of the journal distance matrix is calculated as the inverted ratio between the number of publications in which  $i$  and  $j$  are co-cited and the minimum number of publications where  $i$  or  $j$  is cited (Equation 3.1). The denominator of the ratio accounts for the fact that the probability of being co-cited also depends on the number of publications where each of the two journals for which we are measuring the distance appears in the references. The distance measure ranges from one to infinity. In the case of an infinite distance, i.e.,  $i$  and  $j$  are never co-cited in a publication, for computational reasons, we attribute the maximum non-infinite distance of journal  $i$  from all other journals.

$$D(i, j) = \frac{1}{\min(\text{\#pubs where } i \text{ is cited, } \text{\#pubs where } j \text{ is cited})}$$

(Equation 3.1)

We use the journal distance matrix (D) to calculate the average cognitive distance between a scientist (S) and her team members (T). We consider the journals cited by the individual and the journals cited by her team before the team formation; then, we calculate the average distance, as in Equation 3.2.

$$Cognitive\ distance_{S,T} = \frac{\sum_{i=1}^{\#S} \sum_{j=1}^{\#T} D(i,j)}{\#S * \#T}$$

(Equation 3.2)

where #S is the count of journals cited by S, and #T is the count of journals cited by the other team members. We consider the average distance calculated in Equation 3.2 as our measure of the cognitive distance between the scientist S and her team T (*Cognitive distance<sub>S,T</sub>*).

### Individual characteristics

We consider demographic characteristics as determinants of individual learning from within the team, such as age and gender of the scientist. We include in the regression exercise a dummy *Gender* that equals one if the individual is a female and zero otherwise. We include the age of the individual at the time of team formation (*Age*). Since the probability of observing a new citation is correlated with the scientist's ability, we control for the *Stock of publications pre-team entry*, the *Average number of citations per paper*, and her knowledge stock before entering the team (*Stock of journals cited pre-team entry*). As additional controls, we consider the scientist's experience in SINERGIA project applications. We thus include a dummy *Multiple current applications*, which is equal to one if the individual is participating in more than one project at the same time, and zero otherwise. Finally, we take into account the number of previous applications, *Previous applications*, and the number of successful ones, *Previous awarded applications*.

Table 3.3 reports the descriptive statistics concerning the total number of new journals cited after team formation, and the proportion of these citations originating from within the team, according to the characteristics of the applicants. Females tend to have a lower number of new journals cited than males. Young and less experienced scientists learn more than older and more experienced scientists. However, t-tests show that these differences are not significant.

**Table 3.3: Average number of new journals cited post application, and of learning from within the team by gender, age and stock of publications pre-team entry**

	With co-authored pubs.			No co-authored pubs.		
	A. Number of new journals cited post application	B. Learning from within the team	B/A	C. Number of new journals cited post application	D. Learning from within the team	D/C
<i>All scientists</i>	109.20	43.96	0.40	98.09	37.22	0.38
Female scientists	120.73	48.51	0.4	105.34	49.14	0.47
Male scientists	145.08	56.7	0.39	131.39	46.27	0.35
<i>t-test</i>	0.08	0.14		0.05	0.46	
Young scientists (Age<49)	142.41	57.87	0.41	126.29	48.03	0.38
Senior scientists (Age>49)	140.03	52.39	0.37	128.83	46.05	0.36
<i>t-test</i>	0.82	0.17		0.79	0.58	
Large stock of pubs. (Stock>43)	155.2	54.64	0.35	144.92	49.14	0.34
Limited stock of pubs. (Stock<43)	135.18	55.81	0.41	119.59	46.27	0.39
<i>t-test</i>	0.07	0.78		0.02	0.46	

### Journal characteristics

We measure learning by relying on the new journals cited by the scientist. The journal characteristics might affect the number of citations that a journal receives. Hence, in our regression exercise, we control for the following journal characteristics. First, we include in the regression the number of articles where the journal is cited, *Journal frequency*. Second, we control for the fact that the journal is a generalist journal, *Generalist*. As generalists, we consider the following journals: Nature, Science, PNAS, and PlosOne. Finally, we control for the length of the history of the journal proxying its foundation year by the year when the first article published in the journal appears in our database, *History of journal*. For approximately 10% of the journals, we are not able to identify the founding year; in the case when such information is missing, we control with a dummy, *Unknown history*.

In Table 3.4, we consider descriptive statistics at the scientist-journal cited level of analysis adopted in the regression exercise.

**Table 3.4: Regression descriptive statistics considering the study sample used in Table 3.5, Column 4 (106,898 scientist-journal cited observations)**

VARIABLES	Mean	Std.	Min	Max
Learning from within the team- no co-authored pubs.	0.37	0.48	0.00	1.00
<i>Individual characteristics</i>				
Age	47.80	7.76	30.00	71.00
Gender (=1 for female, 0 otherwise)	0.13	0.33	0.00	1.00
Stock of pubs. pre-team entry	43.91	37.81	1.00	318.00
Stock journals cited pre-team entry	182.07	114.13	1.00	644.00
Average number of citations per paper	5.54	4.68	0.05	51.05
Multiple current applications	0.15	0.36	0.00	1.00
Previous awarded applications	0.11	0.31	0.00	1.00
Previous applications	0.33	0.47	0.00	1.00
<i>Team characteristics</i>				
Co-ethnic team	0.13	0.34	0.00	1.00
At least one female scientist on the team	0.40	0.49	0.00	1.00
Awarded	0.44	0.50	0.00	1.00
High-quality application (grade A)	0.10	0.30	0.00	1.00
Low-quality application (grade D)	0.18	0.38	0.00	1.00
Amount requested	1796821	795390	349901	6854573
Number of team members	4.51	1.65	2.00	11.00
Number of sub-disciplines	3.46	2.22	1.00	11.00
Science & Medicine	0.72	0.45	0.00	1.00
<i>Geographical distance</i>				
Distance hours	3.37	4.02	0.00	35.00
<i>Social distance</i>				
Same gender scientist vs. team	0.92	0.28	0.00	1.00
Standardized stock pubs. difference scientist vs. team	-0.04	0.92	-4.08	6.35
Standardized age difference scientist vs. team	-0.10	1.00	-2.77	3.13
Established collaboration	0.34	0.47	0.00	1.00
<i>Cognitive distance</i>				
Log(Cognitive distances <sub>S,T</sub> )	4.70	0.43	1.05	6.14
<i>Journal characteristics</i>				
Journal frequency	480.94	1026.06	51.00	17438.00
Generalists (NATURE,SCIENCE,PNAS,PLOS)	0.01	0.08	0.00	1.00
History of journal (Obs. 96,370)	34.81	18.34	3.00	86.00
Unknown history	0.10	0.30	0.00	1.00

### 3.4 Estimation strategy

We estimate with two Probit models the probability that a new journal cited by the scientist, after team formation, is the result of a process of learning from other team members. In the first model, we consider *Learning from within the team* as a dependent variable, while in the second model, we consider *Learning from within the team – no co-authored pubs* as a dependent variable.

In the two Probit models, we maintain the same set of explanatory variables. We group the explanatory variables into four vectors, namely *team characteristics*, *scientist vs. team*, *individual characteristics*, and *journal characteristics* (Equation 3.3).

$$\begin{aligned}
 P(\text{Learning from within the team} = 1 | \mathbf{x}) = \\
 G(\beta_0 + \mathbf{team} * \beta_1 + \mathbf{scientist vs. team} * \beta_2 \\
 + \mathbf{individual} * \beta_3 + \mathbf{journal} * \beta_4 + \mathbf{controls} * \beta_5)
 \end{aligned}$$

(Equation 3.3)

where G is the standard normal cumulative distribution function. In our estimations, we clustered standard errors at the scientist level.

### 3.5 Results

Table 3.5 reports the results of the estimation of the probability of learning from within the team. In the regressions of columns 1 and 3, we adopt *Learning from within the team* (118,602 scientist-journal cited observations) as a dependent variable, while in the regressions of column 2 and 4 we adopt *Learning from within the team – no co-authored pubs* (106,898 scientist-journal cited observations) as a dependent variable. Columns 1 and 2 consider only the controls, individual and journal characteristics, while columns 3 and 4 add the team characteristics and the scientist vs. team measures.

In columns 1 and 2, we find a weak correlation between the individual characteristics and the probability of learning from within the team. In particular, age (*Age*), gender (*Gender*) of the scientist, the stock of publications (*Stock of publications pre-team entry*), and the average number of citations per paper (*Average number of citations per paper*) have no impact on the probability of learning from the other team members. The stock of journals cited (*Stock of journals pre-team entry*) and having multiple current applications (*Multiple current applications*) have a negative impact on the probability of learning from within the team. After controlling for the team characteristics, i.e., columns 3 and 4, we find the same effect as in columns 1 and 2 for the individual characteristics' variables. The only exception is the significant positive effect of the variable *Stock of publications pre-team entry*. This change is due to the

introduction of the stock of publications difference between the scientist and her team variable (*Standardized stock pub. difference scientist vs. team*).

The journal characteristics are significantly correlated with the probability of learning from other team members. In particular, when journals are frequently cited in the bibliographies of the articles in our dataset (*Journal frequency*) and have a long history (*History of journal*), it is more likely that the citation to the corresponding journal originates from other team members. If the citation refers to a generalist journal (*Generalist*), it is less likely to originate from other team members.

When adding the team characteristics, we find that having at least one female member on the team (*At least one female scientist on the team*) and being member of a co-ethnic team (*co-ethnic team*) have no significant effect on the probability of learning from other team members. The coefficient of the dummy *Awarded* is negative and barely significant when we consider estimations in column 4, but we refrain in interpreting this result since it is not robust across different model specifications. In fact, in a regression available upon request, where we do not control for the quality of the application (*High-quality application (grade A)* and *Low-quality application (grade D)*), the dummy *Awarded* loses its statistical significance.

Finally, in larger teams, the scientist has greater chances to learn from her teammates, since the coefficients accounting for the size of the team (*Amount requested* and *number of team members*) are both significantly positive. The number of sub-disciplines involved in the application has a positive and significant effect (*Number of sub-disciplines*). Scientists in the fields of Science & Medicine (i.e., dummy *Science & Medicine* equal to one) are associated with higher chances of learning from other team members than teams in Engineering (dummy equal to zero).

Looking at the individual characteristics of the scientist in comparison to the rest of her team, we find that the coefficient of the variable capturing the geographical distance (*Distance hours*) is not significant. This result is in line with the findings of Mairesse and Turner (2005), who state that there is no impact of geographical distance on collaboration. An alternative explanation might be that scientists applying for a common grant are already committed to overcoming the costs induced by being located in different places. However, most of the geographical distance values obtained in our sample are relatively low (average of 3.4 hours), since all applications in our dataset need to have at least one Swiss affiliation. Hence, we limit the scope of our analysis on geographical distance to scientific collaborations with small geographical dispersion. As for the social distance between the scientist and her team, we observe that when scientific reputation (*Standardized stock pub. difference scientist vs. team*) increases, the probability of learning within the team decreases. This means that team members with limited scientific reputations benefit more from the learning of other team members than experienced



scientists. We observe that the age difference between the scientist and her team (*Standardized age difference scientist vs. team*) and matching with teammates of the same gender (i.e., dummy *Same gender scientist vs. team* equal to one) are not significantly correlated with the probability of learning from within the team. Finally, we find that if the scientist has an established collaboration with at least one member of the team (*Established collaboration*), it is significantly more likely that she will learn from her teammates.

**Table 3.5: Probit estimation results for the probability of learning from within the team**

VARIABLES	(1) Learning from within the team	(2) Learning from within the team no co-authored pubs.	(3) Learning from within the team	(4) Learning from within the team no co-authored pubs.
<i>Individual characteristics</i>				
Age	-0.00086 (0.0011)	-0.00035 (0.0010)	0.0009 (0.0014)	0.00087 (0.0014)
Gender (=1 for female, 0 otherwise)	-0.017 (0.026)	-0.025 (0.026)	-0.0017 (0.018)	-0.0047 (0.018)
Log(Stock of pubs. pre-team entry)	-0.023 (0.016)	-0.014 (0.016)	0.070*** (0.015)	0.077*** (0.016)
Log(Stock journals cited pre-team entry)	-0.045*** (0.014)	-0.050*** (0.015)	-0.088*** (0.014)	-0.086*** (0.015)
Log(Average number of citations per paper)	-0.0018 (0.011)	0.0069 (0.011)	0.0034 (0.0093)	0.0055 (0.0097)
Multiple current applications	-0.055** (0.025)	-0.041 (0.025)	-0.040** (0.020)	-0.03 (0.021)
Previous awarded applications	0.03 (0.022)	0.021 (0.023)	0.015 (0.020)	0.0032 (0.020)
Previous applications	0.015 (0.017)	0.016 (0.018)	0.00055 (0.012)	0.005 (0.013)
<i>Team characteristics</i>				
Co-ethnic team			0.012 (0.018)	-0.0084 (0.019)
At least one female scientist on the team			-0.012 (0.013)	-0.017 (0.013)
Awarded			-0.018 (0.012)	-0.023* (0.012)
High-quality application (grade A)			0.032 (0.021)	0.022 (0.020)
Low-quality application (grade D)			-0.032* (0.018)	-0.026 (0.019)
Log(Amount requested)			0.075*** (0.015)	0.067*** (0.016)
Log(Number of team members)			0.26*** (0.020)	0.26*** (0.022)
Log(Number of sub-disciplines)			0.024*** (0.0092)	0.020** (0.0093)
Science & Medicine			0.066*** (0.017)	0.066*** (0.018)
<i>Continued next page</i>				
Pseudo R2	0.07	0.07	0.11	0.11
Observations	118,602	106,898	118,602	106,898

Note: Coefficients are marginal effects. Robust standard errors are clustered around PIs. \*\*\*, \*\*, \*: Significantly different from zero at the 1%, 5%, 10% confidence levels.

**Table 3.5 (Continued) Probit estimation results for the probability of learning from within the team**

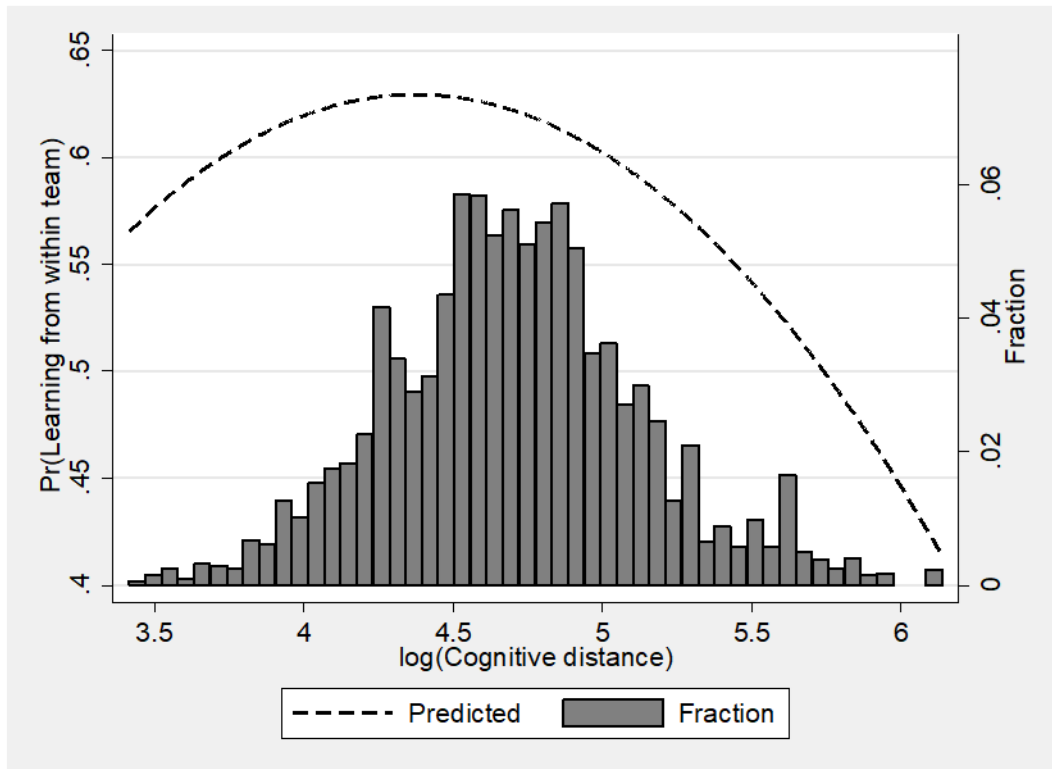
VARIABLES	(1) Learning from within the team	(2) Learning from within the team no co-authored pubs.	(3) Learning from within the team	(4) Learning from within the team no co-authored pubs.
<i>Geographical distance</i>				
Log(1+Distance hours)			0.0082 (0.0081)	0.0086 (0.0084)
<i>Social distance</i>				
Same gender scientist vs. team			0.029 (0.030)	0.031 (0.031)
Standardized stock pub. difference scientist vs. team			-0.100*** (0.011)	-0.10*** (0.011)
Standardized age difference scientist vs. team			-0.013 (0.011)	-0.014 (0.011)
Established collaboration			0.058*** (0.012)	0.046*** (0.013)
<i>Cognitive distance</i>				
Log(cognitive distance <sub>s,t</sub> )			0.55*** (0.20)	0.45** (0.19)
Log(cognitive distance <sub>s,t</sub> ) <sup>2</sup>			-0.064*** (0.021)	-0.055*** (0.020)
<i>Journal characteristics</i>				
Log(Journal frequency)	0.11*** (0.0045)	0.11*** (0.0046)	0.12*** (0.0045)	0.12*** (0.0046)
Generalists (NATURE, SCIENCE, PNAS, PLOS)	-0.083*** (0.021)	-0.057*** (0.022)	-0.096*** (0.021)	-0.069*** (0.022)
Log(History of journal)	0.085*** (0.0042)	0.083*** (0.0042)	0.085*** (0.0043)	0.083*** (0.0042)
Unknown history	0.30*** (0.017)	0.29*** (0.018)	0.29*** (0.018)	0.29*** (0.018)
Pseudo R2	0.07	0.07	0.11	0.11
Observations	118,602	106,898	118,602	106,898

Note: Coefficients are marginal effects. Robust standard errors are clustered around PIs. \*\*\*, \*\*, \*: Significantly different from zero at the 1%, 5%, 10% confidence levels.

Concerning the cognitive distance between the scientist and her team, we find an inverted U-shaped curve. This relationship means that when the cognitive distance between the scientist and the rest of her team is excessively high or low, the scientist is less likely to learn from within the team. A medium level of cognitive distance maximizes the probability of learning from within the team. More precisely, we find a statistically significant and positive coefficient for the linear term of the variable *Cognitive*

*distance* and a statistically significant and negative coefficient for its quadratic term<sup>26</sup>. We estimate the curve plotted in Figure 3.3 below, based on a simplified linear probability model, adopting the same specification of the Probit model. These results confirm the findings of Nooteboom et al. (2007) concerning the two opposing effects of the cognitive distance on the ability to learn something new.

**Figure 3.3: Predicted probability of learning from within the team for the representative individual\* vs. cognitive distance**



\*For us, the representative individual is a male scientist of average age, publication stock, and citations received. He applied for the first time to a SINERGIA. His team is made only of Swiss scientists, has no female members, has been awarded with an amount requested of 1,500,000 CHF. The team is composed of 4 members who work in 3 sub-disciplines in medicine. The average distance from the other team members is 3h traveling. There is at least another individual of the same gender within the team. He has the same stock of publications and seniority as his teammates. He does not have previous collaborations with his teammates.

According to our definition of a team, individuals work together for a limited period of time, pursuing a circumstantial goal. This definition focuses our analysis on scientists applying for a grant. In order to extend the scope of our analysis, we tested for the possibility of selection bias for this category of scientists. These scientists could have different characteristics from the rest of the population, such as higher

---

<sup>26</sup> Appendix A3 reports an extensive discussion on the statistical significance of this result in relation to the sample size.

ability, which might influence the probability of learning from their teammates. We test the presence of selection bias by implementing a Heckman analysis. We find no evidence of the presence of selection bias<sup>27</sup>.

### 3.6 Discussion and conclusion

This paper contributes to the SciTS literature by identifying the factors that promote the learning of individuals from their team members when working in a scientific research team. Unique to our study is the fact that we measure the basic component of the knowledge stock, and we keep track, within the team, of the knowledge flows from one individual to another.

When an individual enters the team, she contributes to it with her knowledge stock, and, at the same time, she has the occasion to learn from other team members. We find that team characteristics and the characteristics of the scientists, in comparison to the rest of her team, affect her probability of receiving knowledge from the members of the team. Precisely, in terms of social distance, we find that having an established relationship with at least one team member and working with scientists with higher scientific reputations, enhance the probability of learning from within the team. As for the cognitive distance between the individual and the other team members, we find that it generates an inverted U-shaped relationship with the probability of observing knowledge flows from within the team. This result suggests that there is an optimal level of cognitive distance that favors learning. An individual should have a knowledge stock that differs from that of the others in order to guarantee a buffer for learning something new. At the same time, the knowledge stock difference should not be too large so as to avoid 'speaking a different language' obstacles, which can hinder the effective flow of knowledge among team members.

Our results have a direct implication for the SciTS literature. While a large part of this literature shows that researchers working in teams have higher productivity, we focus on the individual and team characteristics that stimulate team members' learning from their colleagues. Nowadays, funding agencies are increasingly promoting the constitution of interdisciplinary teams for conducting research. For

---

<sup>27</sup> The ideal selection equation of the Heckman procedure would include all of the scientists involved in all teams within the whole scientific community. Given that this information is not available, we constructed a sample of non-applicant experienced scientists with similar characteristics as the applicants who are eligible to apply to SINERGIA, but they did not. We used the sample, including both our applicants and the non-applicant experienced scientists to implement the selection equation. The identification variable used in the selection equation is a dummy that equals one if another scientist previously applied for the SINERGIA grant in the institution of affiliation of the focal scientist. This variable is positively and significantly correlated with the probability of applying. In the main equation, we included the inverse Mills ratio computed from the selection equation and, reassuringly, we noticed that its coefficient is not significant.

instance, the latest SINERGIA call for grants states as a new requirement for applicants the need to prove the interdisciplinary composition of the team. In this context, understanding the micro-dynamics of research teams is crucial. Our findings on the social distance and on the optimal level of cognitive distance among team members suggest that, in promoting interdisciplinary teamwork, particular attention should be devoted to team composition. While geographical distance has little impact on the knowledge flows among team members, social aspects should be taken into account. Previous experience in joint research work favors team members' learning. Moreover, scientific reputational differences of team members direct the knowledge flows from productive members to scientists with less experience in research. Finally, while it is common wisdom to promote cross-disciplinary research in order to stimulate creativity, this could have unexpected consequences. It is important to maintain a common knowledge base among team members in order to guarantee knowledge flow absorption.



# Chapter 4 Knowledge diffusion and morality: Why do we freely share valuable information with strangers?

*Disclaimer: This chapter, draws from a working paper written in collaboration with Boris Thurm. Boris was one of the greatest encounters of my PhD. It has been an intense and enriching journey that led us to countless discussions, fruitful collaborations, and ideas. One of them led to the model developed in this chapter, that builds on previous collaborations featured in Boris' dissertation.*

## **Abstract**

Technology enables individuals, scientists, and organizations to share valuable data and knowledge in new ways, not possible before. Scholars are divided on how this phenomenon emerges, especially among strangers. The classical *homo oeconomicus* type of preference does not explain this behavior. If individuals were simply self-centered, they would choose to keep for themselves the valuable information they hold, especially in the absence of any contract or guarantee of reciprocity. In this paper, we explain why some individuals are willing to share valuable knowledge at their own cost by crafting a model with heterogeneously-moral individuals involved in a sharing social dilemma. Our model builds on the recent literature showing that moral incentives are favored by evolution theoretically and have a strong explanatory power empirically. Our analysis highlights the limit of financial incentives and the importance of promoting a sharing culture by enhancing awareness. Shedding light on how people respond not only to financial but also moral incentives, we contribute to the ongoing policy debate on the design of effective open science policies.

## 4.1 Introduction

Regularly, when confronted with a new concept or when we want to verify a piece of information, we head to consult the corresponding Wikipedia page. As teachers, we consult available online classes to get inspiration in the design of our courses. As programmers, we seek precious debugging advice on Stack Overflow. Also, as researchers, we benefit from the freely-available research and voluntarily-shared data. All these informational public goods bring great value in terms of social welfare. However, the availability of this information largely depends on the willingness of contributors who most often voluntarily share without any financial compensation. In this context, it is particularly relevant to better understand the incentives that motivate individuals to freely share knowledge and data to the largest audience. This paper aims to address this question.



The act of sharing knowledge and data has the properties of a social dilemma (Olson, 1965; Lichbach, 1996). Contributing to publicly available knowledge maximizes social welfare, but the time and effort required and the lack of compensation drive individuals out of it. With the common assumption of self-centered *homo oeconomicus* agents, individuals restrain from contributing, and the level of sharing is sub-optimal. To address this issue, economists have developed various incentive mechanisms such as taxes, intellectual property rights, and public subsidies (Samuelson, 1954; Arrow, 1962; Stiglitz, 1989). However, a significant share of the publicly available knowledge lies outside of these designed incentive mechanisms. Although often acknowledged by economists working on knowledge (Arrow, 1962; Stephan, 1996; Dasgupta and David, 2002), few models integrate non-self-centered motives. The classical economic approach thus fails to explain the contribution to crowdsourced initiatives such as open-source software and public repositories.

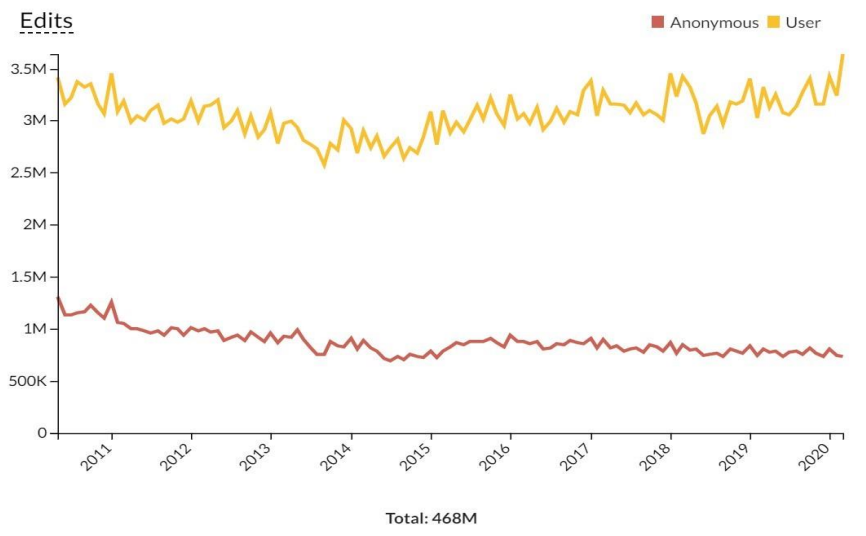
Voluntary contribution to public goods is not limited to knowledge and data sharing. In fact, whether by financing charity actions that provide no direct personal benefit (Carpenter and Matthews, 2017; Andreoni, 1988) or by performing environmentally friendly actions that are time- and effort-consuming (Bruvoll et al., 2002; Sundt and Rehdanz, 2015), selfless contributions to public goods are quite common.<sup>28</sup> Therefore, economists have introduced several preferences challenging the classical self-centered *homo oeconomicus* approach, such as reciprocity (Fehr and Gächter, 1998), social status (Auriol and Renault, 2008; Besley and Ghatak, 2008) and warm glow (Andreoni, 1990), to name a few.

However, these prosocial motives are arguably not sufficient in the particular case of knowledge and data sharing. First, individuals sharing have little guarantee of receiving something in return from people using their contribution, which discounts the reciprocity option. Second, looking at the contributors to the Wikipedia pages over the last ten years (see Figure 4.1), around a quarter were anonymous, while non-anonymous contributors have usernames that rarely match their real names (Gallus, 2017).<sup>29</sup> Hence, it seems rather unlikely that the prosocial actions of sharing individuals arise solely from their desire for status or their pursuit of a warm glow since many cannot even boast about their contribution.

---

<sup>28</sup>While classical models with a *homo oeconomicus* type of preference would suggest no contribution to a public good, empirical evidence from field experiments shows that the level of contribution to the public good is consistently higher than 40% even with varying conditions (Andreoni, 1995; Ockenfels, 1993; Fischbacher et al., 2001; Burton-Chellew et al., 2016).

<sup>29</sup> For more details see <https://stats.wikimedia.org>

**Figure 4.1: Monthly edits on the English Wikipedia between April 2010 and April 2020.**

In this paper, we integrate a specific type of personal motivation: morality. In the sense developed by Kant (1870), morality consists in accounting for the potential outcome of one's action if all others acted similarly.<sup>30</sup> In our model, agents have *homo moralis* preferences. They maximize a weighted average of their selfish payoff and of the payoff that they would get if all individuals act in the same way. We build on the recent theoretical contributions of Alger and Weibull (2013, 2016) and Roemer (2015), showing that *homo moralis* preferences have an evolutionary advantage. Recent empirical evidence also suggests that among six different candidate preferences, the *homo moralis* type of preference offers the highest explanatory power to predict individuals' contribution in public good games (Miettinen et al., 2020).

The *homo moralis* preference offers a solid explanation for non-selfish behaviors. However, empirical evidence suggests that the propensity to share knowledge varies greatly among individuals (Andreoli-Versbach and Mueller-Langer, 2014; Hergueux et al., 2015). This empirical observation is backed up by theoretical findings exhibiting the evolutionary stability of a heterogeneous population of *homo moralis* individuals (Ayoubi and Thurm, 2018). We, therefore, introduce in the model a diverse population of *homo moralis* agents having different levels of morality.

The model consists of a simple framework embedding both the morality of individuals and the heterogeneity in their preferences. The setting we develop explains why individuals share valuable and costly knowledge, the observed level of their contribution, as well as the determinants of this

<sup>30</sup> "Act only according to that maxim whereby you can, at the same time, will that it should become a universal law." (Kant, 1870).

contribution in the larger population. Our work contributes to the extant literature in three ways. First, it adds to the scant theoretical literature on the prosocial incentives for knowledge and data sharing with Stephan (2012), calling for research in this area. Second, while most analyses on the economics of knowledge focus on the financial incentives for the production and diffusion of knowledge, the introduction of moral preferences allows us to discuss innovative incentive mechanisms. Finally, when most theoretical models of sharing use unique representative agents, we account for the empirically observed diversity of preferences (Falk et al., 2018) by proposing a model with heterogeneous agents.

The rest of the paper is organized as follows: Section 4.2 discusses the theoretical setting and reviews the literature on non-pecuniary incentives for knowledge sharing. Section 4.3 presents the model, establishing the conditions for individual cooperation and the average share of cooperation in the population, Section 4.4 discusses the main findings and derives some policy implications, and Section 4.5 concludes.

## 4.2 Knowledge and data sharing

### 4.2.1 The economic properties of shared knowledge

When James Clerk Maxwell, building on the experimental work of Michael Faraday, laid the foundation of the electromagnetic theory, he probably did not foresee that his research would be at the cornerstone of modern societies. Without electromagnetism, there would be no electric power distribution system and even less internet.<sup>31</sup> This little story highlights two key features of knowledge. First, it holds substantial social benefits as it is at the core of economic development (De la Croix et al., 2018). Second, knowledge production is a cumulative process (Furman and Stern, 2011): the creation of new knowledge depends on the existing knowledge pool. In research and development also, innovative ideas draw from the extant literature, and scientific results are the source of numerous innovations (Debackere and Veugelers, 2005; Ahmadpoor and Jones, 2017; Fleming et al., 2019). These two characteristics render the availability of knowledge a powerful way of increasing social welfare.

Among the pool of existing knowledge, the freely available share plays a crucial role. Bell et al. (2019) suggest that a considerable part of the social inequalities is driven by uneven access to valuable knowledge among individuals. Therefore, increasing the availability of knowledge could yield large welfare gains, not only by stimulating research and innovation but also by helping to reduce inequalities.

---

<sup>31</sup> For a history of electromagnetic theory, you can refer to the following Wikipedia page: [https://en.wikipedia.org/wiki/History\\_of\\_electromagnetic\\_theory](https://en.wikipedia.org/wiki/History_of_electromagnetic_theory)

Notwithstanding the numerous social benefits of publicly shared knowledge, individuals have little economic incentives to contribute to the pool of available knowledge. The process of sharing is costly. Contributors have to allocate time, effort, and sometimes money to produce useful knowledge and to store it in the most appropriate format (Smith et al., 2017). In addition to the individual costs born by the contributors, the nature of knowledge makes it hard for them to internalize the social benefits of the knowledge they shared.<sup>32</sup> In academic research and innovation competitions, for instance, the person sharing information, ideas, or data incurs the cost of a potential future scooping by competitors working on a similar project (Thursby et al., 2018). The reason is that publicly shared knowledge is a *public good*: it is non-rival by nature (its use by an individual does not diminish its value for others) and non-excludable by the decision of the contributor (Samuelson, 1954). These characteristics imply that the process of knowledge sharing takes the form of a social dilemma (Kollock, 1998): individuals' economic incentives are not aligned with social needs.<sup>33</sup>

Standard economic theory assuming self-centered *homo oeconomicus* preferences predicts that like other public goods creating a positive externality (individuals do not internalize the cost of producing it), publicly available knowledge will be under-supplied compared to optimal social welfare levels (Stiglitz, 1999). The classical response to this market failure is the implementation of policy tools such as the public funding of research centers and universities, or the design of intellectual property rights (IPRs) (Arrow, 1962; Dasgupta and David, 2002; Stiglitz, 2014).

These policy instruments help in countering the market failure by increasing the widespread availability of valuable knowledge. More specifically, the public funding of scientific research accounts for the absence of economic incentives to invest in risky and long-term research projects that do not generate direct commercial benefits (Stephan, 2012). The projects' results are then expected to be available in public reports and scientific publications (Merton, 1973). Similarly, IPRs are intended to increase the financial incentives for developing innovation and for ensuring the diffusion of knowledge by imposing disclosure (De Rassenfosse et al., 2016; De Hopenhayn and Squintani, 2016).

The classical economic approach of knowledge sharing as a market failure and the tools put in place to counter it have proven effective in providing financial incentives to produce and diffuse knowledge at more socially desirable levels (Singh, 2005; Gangopadhyay and Mondal, 2012; Poege et al., 2019). However, these policy instruments are costly and come with several inherent inefficiencies.

---

<sup>32</sup> Since individuals are unable to internalize its benefits, knowledge creates what we often call a positive *externality*.

<sup>33</sup> This situation is often called a market failure as market incentives fail to lead to a socially optimal situation.

For instance, in the case of IPRs, the incentives to disclose come at the cost of a static inefficiency arising from the monopoly power granted to the holder of the patent (Stiglitz, 2007; Lerner, 2009). As for public funding, several studies show that the allocation system does not necessarily reward the most deserving researchers (e.g., Merton, 1988; Rigney, 2010; Ayoubi et al., 2019). More importantly, the publicly funded academic system has not yet proven effective in ensuring proper dissemination of knowledge (McKiernan et al., 2016).

Furthermore, many individuals decide to contribute to the diffusion of knowledge and data without being directly affected by any of these instruments (Moser, 2013). The main drivers for these actions are non-financial incentives for sharing.

#### 4.2.2 Non-financial incentives for knowledge sharing

In various contexts, the level of contribution to producing and disseminating knowledge is higher than what the classical economic theory would suggest. Numerous studies in the literature have qualitatively and quantitatively accounted for the determinants of knowledge sharing. These studies can be classified into three environments where non-financially motivated sharing of knowledge and data has played a central role: Open source software (OSS), online public repositories, and scientific research.

##### **Open-source software**

One of the most impressive and surprising successes of the last 25 years was the development of high-quality software that was not produced by any firm but simply by the voluntary contribution of millions of individuals online. Open-source software includes very successful products such as Linux, Apache, or Python, to name a few. The OSS platform Sourceforge hosts more than half a million projects, with more than 35 million monthly users worldwide and over four million downloads a day.<sup>34</sup> The OSS community relies on the dedication of users to improve the software and to freely share their upgrades with others, who then build on it to improve and share results again.

Since its development, the operating mode of OSS, and in particular, the absence of financial compensation for the contributors, has intrigued economists (Lerner and Tirole, 2002, 2005). Many qualitative and quantitative studies aimed at better understanding the dynamics of sharing in OSS. The pioneering work of Lee and Cole (2000), analyzing contributors to Linux, identified a few compelling findings. First, the development of OSS was contingent on the progress of web-based tools, making the communication and sharing of information easier and more efficient. Second, the majority of individuals

---

<sup>34</sup> See [sourceforge.net/about](https://sourceforge.net/about).

participating in the Linux project (around 64%) worked on the software during their leisure time. Third, as main reasons for participating in the project, developers mentioned the social usefulness of the product, the recognition by their peers and sense of identity, and the anticipated reciprocity of their actions. Beyond the case of Linux, a comprehensive survey by David et al. (2003) on 2784 developers found that the motivation to commit to OSS projects being rather heterogeneous over the sample, but a majority of contributors do not get any kind of payment for their work. In essence, while some developers emphasized that indirect financial motives such as future career prospects are important (Fershtman and Gandai, 2007), most contributors put forward intrinsic motives as the main reason for participating in this collaborative effort (Lakhani and Wolf, 2003). The most prominent non-financial motives collaborating are the learning benefits of coding with experts (Von Krogh et al., 2003; Hippel and Krogh, 2003; Lakhani and Von Hippel, 2004), the sense of belonging to a community (Lee and Cole, 2003), and the generalized access to a better quality product (Gächter et al., 2010). Many of these motivations are also relevant when contributing to public repositories of knowledge.

### **Public repositories**

A more recent phenomenon of free sharing of valuable knowledge is the contribution to public repositories such as Wikipedia, Github, or Stack Overflow. With more than 50 million articles in around 300 different languages, Wikipedia is the world's largest encyclopedia.<sup>35</sup> Although most of its articles are written and edited by willingly-contributing unpaid users, Wikipedia has managed to achieve a remarkable level of quality (Giles, 2005; Liu and Ram, 2018) and unbiasedness (Greenstein and Zhu, 2012). However, Wikipedia is not a unique success story as it is possible today for any person with an internet connection to access reliable information provided by volunteers on subjects as diverse as debugging code, preparing a recipe, gardening, and finding the solution of a riddle. With their wide accessibility, these platforms bring considerable social value, especially for individuals that did not previously have access to high-quality knowledge sources (Teplitskiy et al., 2017).

Considering the growing importance of public repositories as a source of valuable knowledge, several scholars have investigated the motives for editors to contribute. In a recent article, Xu et al. (2020) study the behavior of more than sixty thousand North American contributors to Stack Overflow. By tracking the career choices of users, they show that indirect career benefits influence the tendency to contribute to the pool of answers on the platform: Stack Overflow partly functions as a signaling device for potential recruiters. However, this argument does not hold for many contributors to public

---

<sup>35</sup> See [meta.wikimedia.org/wiki/List\\_of\\_Wikipedias](https://meta.wikimedia.org/wiki/List_of_Wikipedias).

repositories. In an experiment on the Swiss Wikipedia, Gallus (2017) finds that contributors are sensitive to symbolic rewards that increase their sense of belonging and usefulness, even when the rewards grant no career or financial benefit. Using as natural experiment the blocked access to Wikipedia on mainland China, Zhang and Zhu (2011) show that the contribution of users outside of mainland China (who experienced no change in their access) decreased by 43% on average during the blockage period. This result suggests that voluntary contributors are sensitive to the social effects of their efforts since their motivation to contribute goes down when the audience is reduced.

## Scientific research

In OSS and public repositories, the incentives to produce and to share knowledge are usually aligned, as both processes usually happen hand in hand. However, the situation is different in research. The funding of scientific projects depends on the ability of researchers to produce new ideas and results. Hence, before publishing and setting priority on their findings, the incentives to share data and preliminary results are very low (Dasgupta and David, 2002; Stephan, 2012). More specifically, even after publishing, the current scientific publication system leads to less than 13% of papers in Scopus being listed as open access.<sup>36</sup> This observation fuels a growing debate in the scientific community on the necessity to establish an open science paradigm, i.e., an environment of active public sharing of data and scientific articles (Gewin, 2016). While some researchers have expressed their concerns that open science could encourage free-riding behaviors (Longo and Drazen, 2016), public institutions in Europe and in Switzerland have recently implemented regulations favoring open access to scientific articles and data (Guedj and Ramjoué, 2015; Spichiger, 2018).<sup>37</sup> These decisions are motivated by the potential of open science to increase the diffusion of knowledge, improve the reliability of results, avoid excessive duplication, and help early-career researchers with less visibility (David and Foray, 1996; Boudreau and Lakhani, 2015; Munafò et al., 2017; Farnham et al., 2017). Nonetheless, the implementation of a widespread open-science attitude remains limited so far (McKiernan et al., 2016). Therefore, it is essential to improve our understanding of what drives researchers to share their work.

As already stated by Arrow (1962) and more recently discussed by Dasgupta and David (2002) and David (2004), non-pecuniary motivations are a key reason leading scientists to contribute and share

---

<sup>36</sup> Around 8 million open access articles out of 64 million found on the Scopus database. The trend is however increasing over the last two decades with 20% since 2001 (7M over a total of 35M article) and around 28% since 2011 (5.5M out of 20M).

<sup>37</sup> For instance, the European Commission H2020 Programme requires an open access to all the results, data, and peerreviewed scientific publications produced by the projects funded by the Program: [https://ec.europa.eu/research/participants/data/ref/h2020/grants\\_manual/hi/oa\\_pilot/h2020-hi-oa-pilot-guide\\_en.pdf](https://ec.europa.eu/research/participants/data/ref/h2020/grants_manual/hi/oa_pilot/h2020-hi-oa-pilot-guide_en.pdf)

knowledge, even in the absence of financial incentives. Stern (2004) shows that researchers incur financial sacrifices to pursue their career paths, and thus have strong intrinsic motivations. In a study including more than 1600 scientists, Haeussler (2011) finds that researchers are willing to share even at a personal cost, but there exists substantial heterogeneity depending on the fields. Indeed, the decision to share is affected by the expected reciprocity. Since some communities of researchers are more inclined to share, their actions reinforce the sharing culture. These results are confirmed by Thursby et al. (2018). Thanks to a survey involving scientists in nine different disciplines, they show that a majority of researchers chose to disclose their results before publication, their behaviors depending on the number of researchers working on similar questions and on the sharing culture in the field. In a study of peer-reviewing, another key activity of scientists, Squazzoni et al. (2013) find that paying scientists actually reduces the quality of the review because referees are more sensitive to moral motives when evaluating the work of fellow researchers. Accounting for these empirical findings, the following section introduces a model integrating moral motivations for sharing knowledge.

## 4.3 Model

This section presents a novel model with individuals involved in a knowledge-sharing dilemma. We first describe the setting and the main definitions (section 4.3.1). We then analyze the level of contribution in the population (section 4.3.2). Finally, the last two subsections (4.3.3 and 4.3.4) introduce two extensions to the model in line with empirical observations.

### 4.3.1 Setting

#### Social dilemma

We consider a continuum of individuals  $i \in I = [0,1]$  involved in a knowledge-sharing social dilemma: each individual has a personal incentive to abstain from sharing, but the social welfare is higher when individuals fully share. For each individual  $i \in I$ , we note  $x_i \in [0,1]$  the (individual) degree of sharing. The degree of sharing here encompasses both the quality and the quantity of the contribution. The average level of sharing in the population is then defined as  $\bar{x} = \int_0^1 x_i d\mu(i)$  where  $\mu$  is a density measure for the population.

The payoff received by each individual  $\pi_i$  depends on her own strategy  $x_i$  as well as on the level of sharing in the population  $\bar{x}$ . We assume that  $\pi_i: [0,1] \times [0,1] \rightarrow \mathbb{R}$  is continuous and differentiable in  $x_i$  and in  $\bar{x}$  for all individuals  $i \in I$ . The social dilemma setting then implies that, for all individuals  $i \in I$ :  $\pi_i$  is strictly decreasing in  $x_i$  (there is a cost incurred to share), and strictly increasing in  $\bar{x}$  (individuals



benefit from the higher available knowledge). Moreover, we assume that there is a social benefit of contributing, i.e., individuals are better off when all share than when nobody does.

Formally, we have:

$$\forall i \in I, \text{ and } \bar{x} \in [0,1], \text{ if } x_i^1 > x_i^2 \text{ then } \pi_i(x_i^1, \bar{x}) < \pi_i(x_i^2, \bar{x}) \quad (4.1)$$

$$\forall x_i \in [0,1], \text{ if } \bar{x}^1 > \bar{x}^2 \text{ then } \pi_i(x_i, \bar{x}^1) > \pi_i(x_i, \bar{x}^2) \quad (4.2)$$

$$\forall i \in I, \pi_i(1,1) > \pi_i(0,0) \quad (4.3)$$

Furthermore, for each individual  $i$ , we call  $C(x_i)$  the individual cost associated with sharing knowledge such that  $C: [0,1] \rightarrow \mathbb{R}_+$  is continuous, differentiable and strictly increasing in  $x_i$  (i.e., spending time and effort to contribute is a cost, as specified in equation 4.1). Similarly, let  $\xi(\bar{x})$  be the positive externality of knowledge such that  $\xi: [0,1] \rightarrow \mathbb{R}_+$  is continuous, differentiable and strictly increasing in  $\bar{x}$  (i.e., more contribution increases the positive externality, as specified in equation 4.2). Assuming separability, for a given cooperation share  $\bar{x}$ , the individual payoff for individual  $i \in I$  playing  $x_i$  can then be written as follows:

$$\pi_i(x_i, \bar{x}) = \xi(\bar{x}) - C(x_i) \quad (4.4)$$

Since we assume that the payoff when everybody cooperates is higher than when everybody defects (equation 4.3), we have, for all  $i \in I$ ,  $\xi(1) - C(1) > \xi(0) - C(0)$ . Setting the value of the cost and the externality to zero in zero,<sup>38</sup> we have:  $\xi(1) > C(1)$ .

### Distribution of preferences

The decision of individuals derives from the maximization of their utility. We consider a population with *homo moralis* preferences. The utility of a *homo moralis* individual is a weighted average between her classic "selfish" payoff ( $\pi_i(x_i, \bar{x})$ ) and a "moral" payoff accounting for the payoff she would get if all the rest of the population acted like her ( $\pi_i(x_i, x_i)$ ):

**Definition 1 (Homo moralis utility).** An individual  $i$  is said to have a homo moralis type of preference with a degree of morality  $\kappa_i \in [0,1]$  when her utility follows:

$$u_{\kappa_i}(x_i, \bar{x}) = (1 - \kappa_i) \cdot \pi_i(x_i, \bar{x}) + \kappa_i \cdot \pi_i(x_i, x_i)$$

---

<sup>38</sup> There is indeed no cost when not sharing, and there is no positive externality if nobody shares.

The case where the degree of morality is equal to zero is the classical *homo oeconomicus* utility, while the case where  $\kappa_i = 1$  is called *homo kantientis*<sup>39</sup> (Alger and Weibull, 2013) or "fully moral" individual. A growing empirical and theoretical literature suggests that *homo moralis* preferences are a more accurate representation of human preferences than most other used utilities (Alger and Weibull, 2016; Capraro and Rand, 2018; Miettinen et al., 2020).

However, recent empirical evidence also suggests that prosocial preferences greatly vary across individuals (Falk et al., 2018; Alger et al., 2019; Awad et al., 2020). Consequently, we consider a heterogeneous population of individuals by varying their degrees of morality. More specifically, for each individual  $i \in I$ , the degree of morality  $\kappa_i \in [0,1]$  is independently drawn from a random distribution over  $[0,1]$  with a cumulative distribution function (CDF)  $F_\kappa$  and density  $f_\kappa$ . For instance, the beta distribution with two parameters  $(a, b) \in \mathbb{R}^*_+$  ( $\kappa \sim \text{Beta}(a,b)$ ) is an example of distribution over  $[0,1]$  offering great flexibility (Gupta and Nadarajah, 2004) that we often use as an illustration in what follows.

### 4.3.2 Analysis

#### Sharing behavior

For a given level of cooperation in the population, each individual  $i \in I$  decides on her degree of sharing  $x_i$  based on the following program:

$$x_i \in \operatorname{argmax}_{x \in [0,1]} [(1 - \kappa_i) \cdot \pi_i(x, \bar{x}) + \kappa_i \cdot \pi_i(x, x)]$$

Alternatively, when integrating equation 4.4, and knowing that the argmax function is unchanged with the application of a strictly increasing function, we have:

$$x_i \in \operatorname{argmax}_{x \in [0,1]} [\kappa_i \cdot \xi(x) - C(x)] \quad (4.5)$$

First, note that in a population made solely of individuals with the classical *homo oeconomicus* type of preference,<sup>40</sup> the solution to the program above is simply a corner solution where nobody shares ( $x_i = 0$  for all  $i \in I$ ) because  $(-C)$  is maximized in 0. For the more general case of a population of *homo moralis* individuals, we have the following theorem:

---

<sup>39</sup> The name kantientis is a tribute to the German philosopher Immanuel Kant whose categorical imperative inspired the formal formulation of morality used in this paper.

<sup>40</sup> This situation is a specific example of our setting with the  $(\kappa_i)_{i \in I}$  following a degenerate distribution equal to zero.

**Theorem 1 (Sharing strategy)**

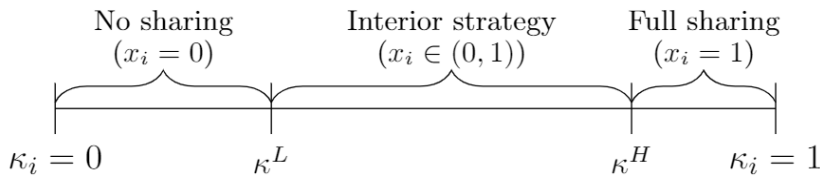
A *homo moralis* individual with a degree of morality  $\kappa_i$  involved in a sharing social dilemma plays a pure strategy ( $x_i \in \{0,1\}$ ) if and only if:  $\kappa_i \leq \kappa^L$  ( $x_i = 0$ ) or  $\kappa_i \geq \kappa^H$  ( $x_i = 1$ ).

Where,  $\kappa^L = \inf_{(0,1]}[\psi_0]$ ,  $\kappa^H = \sup_{[0,1)}[\psi_1]$ , and  $\forall u \in \{0,1\}, \psi_u(x) = \frac{c(x)-c(u)}{\xi(x)-\xi(u)}$

*Proof.* In Appendix A4.1.  $\square$

Theorem 1 implies that the population is divided into three sets. Individuals are in one of the three sets defining their strategy  $x_i$  based on their level of morality  $\kappa_i$  as illustrated in Figure 4.2:

**Figure 4.2: Strategy played by *homo moralis* individuals in a sharing social dilemma depending on their degree of morality  $\kappa_i$**



**Pure strategies**

The result of Theorem 1 can lead to a peculiar case with no interior equilibrium under some simple conditions on the functions  $(\psi_u)_{u \in \{0,1\}}$ . We characterize it in the following corollary:

**Corollary 1 (Population with pure strategies only)**

If the functions  $\psi_0$  and  $\psi_1$  are decreasing on  $[0,1]$ ,<sup>41</sup> then all the individuals in the population play pure strategies, and we have, for all  $i \in I$ :  $x_i = 0$  if  $\kappa_i \leq \psi_0(1)$  and  $x_i = 1$  otherwise.

*Proof.* In Appendix A4.1.  $\square$

<sup>41</sup> Note that the corollary includes the case of non-strictly decreasing functions such as constant functions for instance.

**Application 1: Linear individual cost and externality**

In order to illustrate this situation, we consider a linear form for the cost function  $C(\cdot)$  and the externality  $\xi(\cdot)$ , i.e., we have:

$$\forall x \in [0,1]: C(x) = \gamma x$$

$$\text{And } \forall x \in [0,1]: \xi(x) = \beta_n x$$

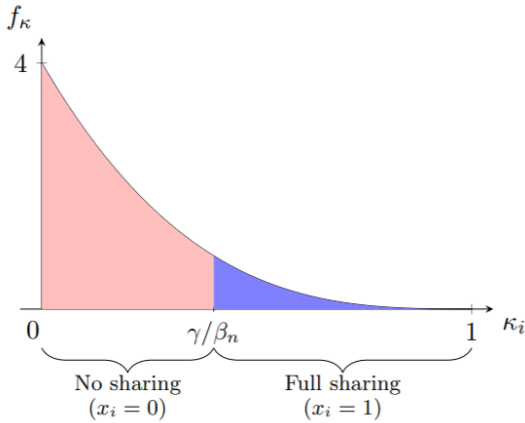
Where  $\gamma > 0$  and  $\beta_n > 0$  represent respectively, the marginal cost of sharing and the marginal benefit produced by the externality. We note the marginal benefit  $\beta_n$  where  $n$  is the size of the population because the positive externality depends on the size of the population.<sup>42</sup> A larger population leads to a higher externality, i.e., we have:  $\frac{d\beta_n}{dn} > 0$ . Note that necessarily  $\gamma < \beta_n$  in order to satisfy the third condition of the sharing social dilemma (equation 4.3).

The functions  $\psi_0$  and  $\psi_1$  are then constant, equal to  $\frac{\gamma}{\beta_n}$ , and we can apply Corollary 1 giving us:

$$\text{For all individuals } i \in I: \quad \text{if } \kappa_i \leq \frac{\gamma}{\beta_n}, \text{ then } x_i = 0, \text{ otherwise } x_i = 1$$

This illustration offers a situation where only two attitudes are possible (either sharing or not). This type of setting can be often encountered by a person holding a piece of knowledge or data and where the question is whether to share this piece or not (i.e., sharing part of the content makes little sense). The population is then divided into two sets (the "contributors" and the "non-contributors") and, for a given distribution of the degree of morality in the population, the size of each group depends on the values of  $\gamma$  and  $\beta_n$  as illustrated in Figure 4.3:

**Figure 4.3: Share of *homo moralis* individuals that contribute in a sharing social dilemma with linear individual cost and externality (illustration with  $\kappa \sim \text{Beta}(1,4)$ )**



<sup>42</sup> An example of  $\beta_n$  taking into account the cumulative property of knowledge discussed in section 4.2.1 is  $\beta_n = \beta e^n$  where an increase of the size of the population increases the pool of people taking advantage of it and the number of people capable of sharing exponentially.

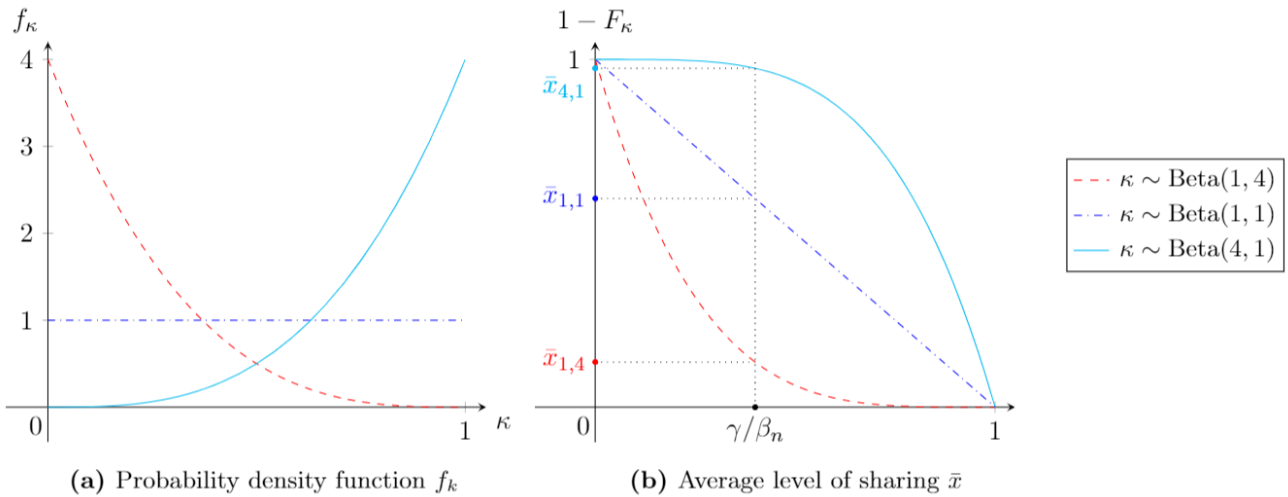
As for the average level of sharing in the population, we have:

$$\bar{x} = \int_{i \in I} x_i d\mu = \int_{\kappa \in [\frac{\gamma}{\beta_n}, 1]} dF_\kappa = 1 - F_\kappa\left(\frac{\gamma}{\beta_n}\right)$$

Hence, regardless of the distribution of morality in the population, since  $F_\kappa(\cdot)$  is increasing, the average degree of sharing in the population is decreasing in the marginal cost of sharing  $\gamma$  and increasing in the externality marginal benefit  $\beta_n$ . This finding is consistent with empirical evidence. Indeed, as shown by Lee and Cole (2000), lowering the individual cost of sharing increases the contribution in the population considerably. Moreover, as suggested by Zhang and Zhu (2011), when the population that benefits from sharing ( $n$ ) increases (and thus when  $\beta_n$  increases), so does the level of contribution in the population.

Furthermore, even with a given cost and size of the population, the distribution of the degrees of morality has a strong impact on the level of sharing in the population. If the distribution is skewed towards higher degrees of morality, the average sharing in the population is higher and vice versa, as illustrated in Figure 4.4.

**Figure 4.4:** Level of sharing in the population under linear individual cost and externality for various distributions of the degrees of morality. When  $\kappa \sim \text{Beta}(4,1)$ ,  $\bar{x}_{4,1} = 1 - (\gamma/\beta_n)^4$ . When  $\kappa \sim \text{Beta}(1,1)$  (Uniform distribution),  $\bar{x}_{1,1} = 1 - \gamma/\beta_n$ . When  $\kappa \sim \text{Beta}(1,4)$ ,  $\bar{x}_{1,4} = (1 - \gamma/\beta_n)^4$ .



## Interior strategies

Following the result of Theorem 1, we can characterize the behavior of individuals playing interior strategies in the Theorem below:

**Theorem 2** (Interior strategy)

A *homo moralis* individual with a degree of morality  $\kappa_i$  involved in a sharing social dilemma has an interior strategy if and only if:  $\kappa^L < \kappa_i < \kappa^H$ . And the degree of sharing  $x_i \in (0,1)$  is then solution to:

$$\frac{\partial C(x_i)}{\partial x_i} = \kappa_i \frac{\partial \xi(x_i)}{\partial x_i}$$

*Proof.* In Appendix A4.1.  $\square$

Application 2: Quadratic individual cost

In order to illustrate a case with interior strategies, we take a quadratic form for the cost function  $C(\cdot)$ , i.e., we have:

$$\forall x \in [0,1]: C(x) = \gamma x^2 \quad (4.6)$$

$$\text{And } \forall x \in [0,1]: \xi(x) = \beta_n x \quad (4.7)$$

Where  $\gamma > 0$ ,  $\beta_n > 0$  and  $\gamma < \beta_n$ . In this case  $\kappa^L = 0$  and  $\kappa^H = \frac{2\gamma}{\beta_n}$ .<sup>43</sup>

Therefore, using Theorem 1, we know that:

- Individuals do not share ( $x_i = 0$ ) if and only if they are *homo oeconomicus*.
- If  $2\gamma \leq \beta_n$ , non-*homo oeconomicus* individuals play an interior strategy if  $\kappa_i < 2\gamma/\beta_n$ , and they fully share ( $x_i = 1$ ) otherwise.
- If  $2\gamma > \beta_n$ , then all non-*homo oeconomicus* individuals play an interior strategy. Moreover, the interior strategy is characterized by Theorem 2:<sup>44</sup>

$$x_i = \frac{\beta_n}{2\gamma} \kappa_i$$

We can then compute the average degree of sharing in the population as follows:

$$\bar{x} = \frac{\beta_n}{2\gamma} \int_{\kappa \in [0, \frac{2\gamma}{\beta_n}]} \kappa dF_\kappa + 1 - F_\kappa(2\gamma/\beta_n) = \frac{\beta_n}{2\gamma} \mathbb{E}[\kappa \mid \kappa < 2\gamma/\beta_n] + 1 - F_\kappa(2\gamma/\beta_n) \quad (4.8)$$

<sup>43</sup> The detailed calculations behind the application are available in Appendix A4.2.

<sup>44</sup> Note that this equation is also valid for *homo oeconomicus* since  $x_i = 0$  when  $\kappa_i = 0$ .

When  $2\gamma \geq \beta_n$ , the average degree of sharing in the population simplifies to:  $\bar{x} = \frac{\beta_n}{2\gamma} E(\kappa)$ . The average level of sharing in the population is thus equal to the average degree of morality in the population weighted by the ratio between the marginal gain from the externality ( $\beta_n$ ) and the cost factor ( $\gamma$ ).

As before, independently of the distribution of morality in the population, the average level of sharing in the population is decreasing in the cost factor  $\gamma$  and increasing in the size of the population  $n$ :

$$\frac{d\bar{x}}{d\gamma} = -\frac{\beta_n}{2\gamma^2} E[\kappa \mid \kappa < 2\gamma/\beta_n] < 0$$

$$\frac{d\bar{x}}{dn} = \frac{d\bar{x}}{d\beta_n} \frac{d\beta_n}{dn} = \frac{d\beta_n}{dn} \frac{1}{2\gamma} E[\kappa \mid \kappa < 2\gamma/\beta_n] > 0$$

Similarly, when the distribution of the degrees of morality is skewed towards higher values of  $\kappa$ , the level of sharing in the population increases. For instance, for  $\gamma/\beta_n = 0.8$ ,  $\bar{x} = 0.5$  when  $\kappa \sim \text{Beta}(4,1)$  (high level of morality in the population), while  $\bar{x} = 0.125$  when  $\kappa \sim \text{Beta}(1,4)$  (low level of morality).

To summarize, the model analyzed in this section puts forward two channels influencing the level of sharing in the population. The first channel is the relative weight of the cost of sharing ( $C(\cdot)$ ) compared to the benefit driven by the externality ( $\xi(\cdot)$ ). Both the individual cost of sharing and the externality benefit depend on the characteristics of the population (such as level of education, population size, access to communication, and sharing technologies). They can also be affected by the policies in place, a topic we discuss in section 4.4. The second channel is the distribution of morality in the population. The level of morality determines the proportion of individuals willing to contribute to the public good as well as their level of contribution. The distribution of the degree of morality in a population mainly depends on cultural and geographical factors (Ayoubi and Thurm, 2018; Alger et al., 2019). The next section (4.3.3) delves into a particular aspect that affects the first channel: peer pressure. Section 4.3.4 considers the impact of financial incentives on the second channel, i.e., on the distribution of morality.

### 4.3.3 Peer pressure

As discussed in section 4.2.2, knowledge sharing is often a matter of communities with strong social influence and peer pressure. When making decisions, individuals are incentivized to harmonize their strategy with the dominant paradigm (Falk and Ichino, 2003). Hence, whether for contributing to an online repository or sharing data among scientists, the level of sharing in the population sets a cultural standard that can push contribution upwards if the norm is to share (Owens, 2016) or downwards if not (Gould and Kaplan, 2011).

In this section, we extend the previous model by integrating peer pressure in the payoff function as an additional cost. Precisely, when an individual plays  $x_i$ , we call  $P : [0,1] \times [0,1] \rightarrow \mathbb{R}_+$  the function capturing the intensity of the peer pressure. The value of  $P(x_i, \bar{x})$  is constructed as to increase when the distance between the individual degree of sharing ( $x_i$ ) and the average degree of sharing in the population ( $\bar{x}$ ) increases. We also set  $P(x_i, \bar{x})$  to zero when the individual and societal decisions are aligned, i.e., we have, for all  $x \in [0,1]$ ,  $P(x, x) = 0$ . This characterization implies that, for each individual, the closer her strategy from the average population strategy, the higher her payoff.

Formally we have:  $\forall i \in I, \pi_i(x_i, \bar{x}) = \xi(\bar{x}) - C(x_i) - P(x_i, \bar{x})$

Therefore, integrating the peer pressure function in the individuals' maximization program, for an individual with a degree of morality  $\kappa_i \in [0,1]$ , the optimal strategy  $x_i$  satisfies:<sup>45</sup>

$$x_i \in \operatorname{argmax}_{x \in [0,1]} [\kappa_i \cdot \xi(x) - C(x) - (1 - \kappa_i)P(x, \bar{x})] \quad (4.9)$$

With the integration of peer pressure, individuals consider not only their own strategy but also the level of sharing in the population when making a decision. In the rest of this section, we explore the application with quadratic cost (application 2 of section 4.3.2) and discuss the impact of peer pressure in that context.

### Application 3: Quadratic individual cost and peer pressure

Taking the functions of cost and externality defined in equations (4.6) and (4.7), and setting, for  $x \in [0,1]$ ,  $P(x, \bar{x}) = \delta(x - \bar{x})^2$ , the program defined in (4.9) becomes:

$$x_i \in \operatorname{argmax}_{x \in [0,1]} [x(c - dx)]$$

$$\text{where } c = \beta_n \kappa_i + 2\bar{x}\delta(1 - \kappa_i)$$

$$\text{and } d = \gamma + \delta(1 - \kappa_i)$$

A *homo moralis* individual, therefore, plays the strategy:<sup>46</sup>

- $x_i = 0$  if and only if  $\kappa_i = 0 = \kappa^L$  and  $\bar{x} = 0$
- $x_i = 1$  if and only if  $\kappa_i \geq \kappa^H = \frac{2\gamma + 2(1-\bar{x})\delta}{\beta_n + 2(1-\bar{x})\delta}$
- Otherwise  $x_i = \frac{2\bar{x}\delta(1-\kappa_i) + \beta_n \kappa_i}{2\delta(1-\kappa_i) + 2\gamma}$ .

<sup>45</sup> We use the fact that the function  $\operatorname{argmax}$  is invariant in strictly increasing transformations and that  $\forall x \in [0,1], P(x, x) = 0$ .

<sup>46</sup> The detailed calculations are available in Appendix A4.2.



We first observe that no one refrains from sharing, even *homo oeconomicus* individuals (i.e.,  $x_i > 0$  for all individuals), except if (almost) all individuals are *homo oeconomicus*. Indeed, individuals with a strictly positive degree of morality always contribute, at least a little. Thus, as long as a little share of the population is made of non-*homo oeconomicus* individuals, the average degree of sharing in the population is strictly positive, and even *homo oeconomicus* individuals contribute.<sup>47</sup>

Second, the effects of changes in cost and externality are similar with and without peer pressure.<sup>48</sup> More precisely, computing the comparative statics for  $\kappa^H$  and  $x_i$ , we observe that  $\gamma$  has a detrimental effect on sharing and  $\beta_n$  a positive one.<sup>49</sup>

The effect of peer pressure is more ambiguous. Considering the strategy of an individual playing an interior strategy  $x_i \in (0,1)$ , a variation of her peer pressure cost factor ( $\delta$ ) has the following influence on her sharing strategy:

$$\frac{dx_i}{d\delta|_i} = \frac{(1 - \kappa_i)(2\bar{x}\gamma - \beta_n\kappa_i)}{2(\delta(1 - \kappa_i) + \gamma)^2}$$

Hence, stronger peer pressure increases the individual degree of sharing if and only if  $\kappa_i < \frac{2\bar{x}\gamma}{\beta_n}$ .

This result suggests that peer pressure increases the sharing propensity only under certain conditions on the degree of morality, on the average sharing in the population ( $\bar{x}$ ), on the cost factor  $\gamma$ , and on the marginal externality benefit  $\beta_n$ . While peer pressure is effective in increasing the sharing behavior on lower morality individuals, it can actually have a detrimental effect for highly moral individuals. As a consequence, the degree of sharing tends to homogenize between individuals, as highlighted by Falk and Ichino (2003).

Finally, the beneficial effect of peer pressure on sharing depends on the initial level of sharing in the population. We can evaluate the impact of the level of sharing in the population on the propensity to share by individual  $i$  as follows:

$$\frac{dx_i}{d\bar{x}} = \frac{2\delta(1 - \kappa_i)}{2(\delta(1 - \kappa_i) + \gamma)} > 0$$

This last result shows that, in the presence of peer pressure, an increase in the average level of sharing in the population also increases the propensity to share of each individual with a degree of

<sup>47</sup> For another example where even *homo oeconomicus* individuals have a strictly positive degree of sharing, see Application 5 in Appendix A4.2.

<sup>48</sup> Note that if  $\delta = 0$  we end up with the same results as in Application 2.

<sup>49</sup> See Appendix A4.2 for the detailed calculations of the comparative statics.

morality  $\kappa_i \in [0,1]$ . This observation contrasts with the situation without peer pressure studied in section 4.3.2 where the individuals' decisions were independent of  $\bar{x}$ , but is more in line with empirical evidence (Owens, 2016; Gould and Kaplan, 2011). Interestingly, since individuals' decisions depend on their perception of  $\bar{x}$  rather than on the actual value, it is essential that individuals perceive that the level of sharing in the population is high to increase their contribution. We discuss the implications of this result in section 4.4.

#### 4.3.4 Financial incentives and morality

Even when considering a population of *homo moralis* individuals, the introduction of financial incentives seems like a desirable approach to increase the level of sharing in the population. The presence of financial subsidies would, in fact, reduce the individual cost of sharing by introducing a financial compensation for it.

Financial incentives can take the form of a subsidy  $\sigma$  paying individuals proportionally to their contribution.<sup>50</sup> More precisely, when an individual plays  $x_i$ , we call  $\sigma(x_i)$  the financial reward for sharing, with  $\sigma: [0,1] \rightarrow \mathbb{R}_+$  a continuous, differentiable, and increasing function. We also impose that, for all  $x \in [0,1]$ ,  $\sigma(x) < C(x)$ , i.e., the financial compensation never exceeds the individual cost (otherwise money would be wasted needlessly). This characterization implies that, for each individual, the higher her sharing, the higher the financial subsidy she receives. The payoff can then be expressed as follows:  $\forall i \in I, \pi_i(x_i, \bar{x}) = \xi(\bar{x}) - C(x_i) + \sigma(x_i)$

Therefore, integrating the financial reward in the individuals' maximization program, for an individual with a degree of morality  $\kappa_i \in [0,1]$ , the optimal strategy  $x_i$  satisfies:<sup>51</sup>

$$x_i \in \operatorname{argmax}_{x \in [0,1]} [\kappa_i \cdot \xi(x) - C(x) + \sigma(x)] \quad (4.10)$$

The financial reward  $\sigma(\cdot)$  acts as a reduction of the individual cost function  $C(\cdot)$ . We can thus reformulate the program above by introducing  $\hat{C}$  such that for all  $x \in [0,1]$ ,  $\hat{C}(x) = C(x) - \sigma(x)$ . Noting that  $\hat{C}(\cdot)$  has the same properties as the function  $C(\cdot)$  defined in section 4.3.1,<sup>52</sup> the solutions to the program (4.10) are given by Theorems 1 and 2. Consequently, since  $\hat{C}(\cdot) < C(\cdot)$  by construction, financial incentives are effective for increasing the level of sharing in the population.

Nevertheless, the effectiveness of the financial subsidy to increase the level of sharing in the population relies on the fact that everything else remains equal when the subsidy is introduced. In

<sup>50</sup> We thus have  $\sigma(0) = 0$ .

<sup>51</sup> We use the fact that the function  $\operatorname{argmax}$  is invariant in strictly increasing transformations.

<sup>52</sup> In particular,  $\hat{C}(\cdot)$  is differentiable, positive and respecting  $\hat{C}(0) = 0$  and  $\hat{C}(1) < \xi(1)$ .

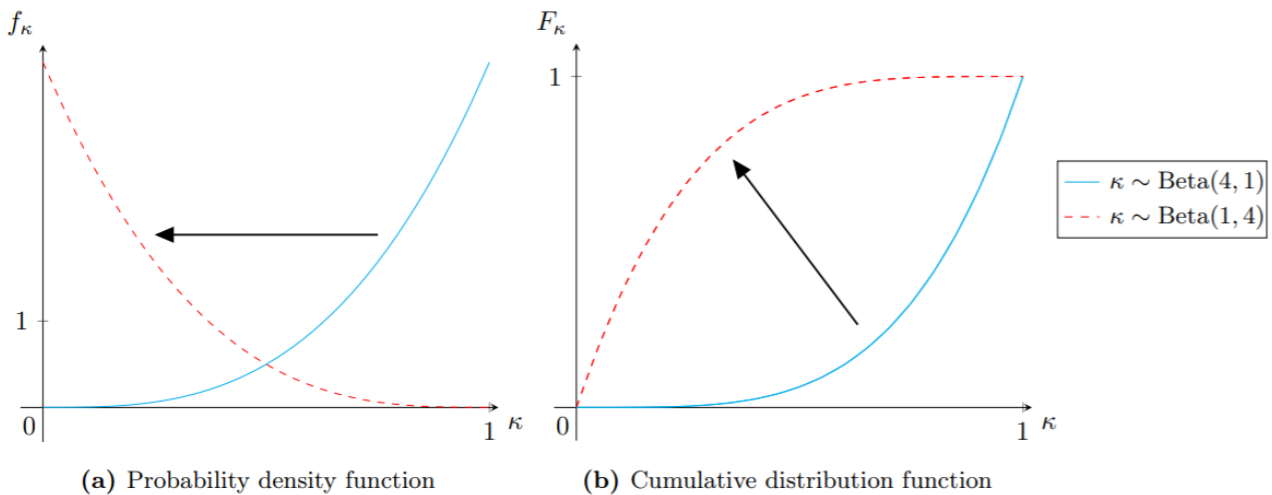
particular, financial compensation works if the distribution of morality in the population remains stable. This situation is, however, rarely met in practice: numerous empirical studies have observed that the morality in a population is highly dependent on the presence of financial compensations.<sup>53</sup> In a leading study, Gneezy and Rustichini (2000) show, for instance, that the introduction of fines for late parents in a day-care unexpectedly increased their tendency to be late, in contradiction to what economic theory would suggest. As extensively discussed by Sandel (2013), "market-based instruments are not inert," they alter the non-financial incentives of individuals. Putting a price tag on an action can push individuals out of the "moral sphere" and into economic calculations.

In our context, the effect of the introduction of financial mechanisms can be modeled by integrating a shock on the distribution of the degrees of morality. The shock then depends on whether financial incentives (or disincentives) are applied. Formally, if we suppose that the distribution of degrees of morality in the population is following a Beta distribution of parameters  $a > 0$  and  $b > 0$ , we model the alteration induced by the introduction of financial motives as follows:

$$\kappa | \lambda_F \sim \text{Beta}(a - \lambda_F, b + \lambda_F)$$

Where  $\lambda_F$  represents the shock induced by the introduction of financial incentives. Figure 4.5 illustrates the effect on the distribution of the degrees of morality, using as an example  $a = 4$ ,  $b = 1$ , and  $\lambda_F = 3$ .

**Figure 4.5: Effect of the introduction of a financial mechanism on the distribution of degrees of morality in the population ( $\lambda_F=3$ )**



In the rest of the section, we illustrate the effect of financial incentives using the second application of section 4.3.2 with quadratic costs.

<sup>53</sup> See Gneezy et al. (2011) for a review of the literature on the subject.

Application 4: Quadratic individual cost and financial incentives

Taking the functions of cost and externality defined in equations (4.6) and (4.7), and setting, for  $x$  in  $[0,1]$ ,  $\sigma(x) = \nu x^2$ , we have for all  $x$  in  $[0,1]$ :  $\hat{C}(x) = (\gamma - \nu)x^2$ . Hence, since  $\gamma > \nu$  by the construction of  $\sigma$ , the setting with the functions  $\xi(\cdot)$  and  $\hat{C}(\cdot)$  is the same as the one in application 2 of section 4.3.2.

Therefore, using Theorems 1 and 2, a *homo moralis* individual plays the strategy:

- $x_i = 0$  if and only if  $\kappa_i = 0 = \kappa^L$ .
- $x_i = 1$  if and only if  $\kappa_i \geq \frac{2(\gamma-\nu)}{\beta_n} = \kappa^H$ .
- $x_i = \frac{\beta_n}{2(\gamma-\nu)} \kappa_i$  if and only if  $0 \leq \kappa_i \leq \kappa^H$ .

Then, using the result of equation (4.8), the average degree of sharing in the population follows:

$$\bar{x} = \frac{1}{\kappa^H} \mathbb{E}[\kappa \mid \kappa < \kappa^H] + 1 - F_\kappa(\kappa^H)$$

where  $\kappa^H = \frac{2(\gamma-\nu)}{\beta_n}$ .

For instance, taking back the example of application 2 in section 4.3.2, we had  $\gamma/\beta_n = 0.8$ , and, with  $\kappa \sim \text{Beta}(4,1)$ , without financial incentives, the level of sharing in the population was  $\bar{x} = 0.5$ . If we consider that the financial incentives reduce the individual cost by a third (i.e.,  $\nu = \gamma/3$ ), then  $\kappa^H = \frac{4\gamma}{3\beta_n}$  and the average level of sharing in the population goes up to  $\bar{x} = 0.75$ . Financial incentives are then beneficial to increase social welfare.

However, if we integrate the effect of financial incentives on the distribution of morality, the conditional distribution of  $\kappa$  follows:  $\kappa \mid \lambda_F \sim \text{Beta}(a - \lambda_F, b + \lambda_F)$ . Thus, with  $\lambda_F = 3$ , we get that  $\kappa \sim \text{Beta}(1,4)$ . Hence, taking the same values for  $\beta_n$ ,  $\gamma$  and  $\nu$ , the average level of sharing in the population becomes  $\bar{x} = 0.1875$ , which is lower than what it was before the introduction of financial incentives.

This simple application illustrates how financial incentives can sometimes have counter-intuitive effects on the contribution to a public good. Moreover, financial incentives also need to be provided by public funds and imply an additional cost for society. One must bear in mind this cost when evaluating the adequate policy to implement by comparing the costs with the benefits retrieved from the introduction of the policy. We discuss these welfare considerations in the next section.

## 4.4 Discussion

This paper adds to the literature on knowledge sharing by crafting a simple model integrating a key feature of human decision making, morality. We designed a model following the construction of homo moralis utility that was proven to have an evolutionary advantage theoretically (Alger and Weibull, 2013; Ayoubi and Thurm, 2018) and a strong explanatory power empirically (Capraro and Rand, 2018; Miettinen et al., 2020). The setting provides a tool to assess the underlying motivation behind the observed willful sharing and to determine the effectiveness of policies in place. In this section, we discuss how the model echoes the literature on knowledge sharing (section 4.4.1), its policy implications (section 4.4.2), and its limitation (section 4.4.3).

### 4.4.1 Contribution of the model

As discussed in section 4.2, most of the theoretical economic literature on knowledge sharing uses representative agents with utility functions in which individuals maximize solely their personal payoff. The model we propose, with heterogeneously-moral individuals, allows reconciling theoretical considerations with empirically observed behaviors. More specifically, while most economic models represent the freely available knowledge as an externality that individuals do not account for in their decision making, our model suggests that individuals actually internalize (at least partly) this externality. This feature of the model has the potential to change our consideration of the knowledge externality and the most effective approach to maximize it. The framework we designed also offers considerable flexibility, which allows to model sharing behaviors in various contexts.

First, by allowing for interior strategies, we account for the variability in the quality and quantity of sharing. Researchers in information science often use three interrelated concepts to categorize the building blocks of the discipline: data, information, and knowledge (Zins, 2007; Badia, 2014). Data is the raw material that can be exploited and polished to produce information that can be analyzed. Information can then be organized and assembled in a way that renders it useful and "actionable," which is what is called knowledge. This relationship can be seen as the degree of refinement of a piece: data is raw and unpolished, information is organized, and knowledge is useful. In our model, the level of sharing ( $x$ ) embeds both the quantity and the quality of the shared element. If one makes a lower value contribution, then the effort is lower, but the usefulness for others is also lower (low individual cost ( $C$ ) and low externality ( $\xi$ ) for a low level of sharing ( $x$ )).

Second, including morality as a motivation for actions, we observe that individuals are willing to engage in individually costly actions because of the potential positive compensation they would get if all others acted similarly. In this sense, our setting accounts for the cumulativeness and social impact of

knowledge discussed in section 4.2 since individuals share because they deem the knowledge others hold as valuable. Moreover, the motivational factors identified in the empirical literature on the subject, such as conditional reciprocity and the access to better general information (see section 4.2.2), are in line with the construction of the model.

Third, while most of the empirical literature on knowledge sharing focuses on the dissemination of quantifiable and codified knowledge, our model is also applicable to the diffusion of tacit forms of knowledge. The examples of sharing loci discussed in section 4.2.2, such as OSS, public repositories, and scientific research, all imply a standardized codification of knowledge before being shared. However, as suggested by Polanyi (1967, 1958), a large part of the knowledge we convey is not translated into tangible codified forms (Cowan et al., 2000). Moreover, as shown by De la Croix et al. (2018), the diffusion of tacit knowledge is at the core of economic development. Hence, since the morality embedded in our model accounts for the sharing of codified or tacit knowledge, it offers a wide scope of application.

Finally, the two extensions described in sections 5.3.3 and 5.3.4, allow the integration of two important empirical observations about individual behavior (peer pressure and the effect of financial instruments on morality) while keeping the main features of the model untouched. These characteristics of the model induce interesting takeaways in terms of policy implications that we discuss in the next section.

#### 4.4.2 Policy implications

The core interest of the model we presented in this paper is to inform policymaking to help in the implementation of effective policies maximizing welfare. By moving away from the classical modeling of individual behavior in economic literature, we bring novel elements to assess the current economic policies and for suggesting new ones.

#### **The limit of financial instruments**

The classical economic approach to tackle market failures and maximize social welfare consists in introducing financial instruments aligning the individuals' interests with social needs. Examples of such policies include public subsidies for producing and disseminating knowledge and intellectual property rights offering the perspective of future revenues (Stiglitz, 2007). While these instruments are effective to reduce the individual cost associated with the production and diffusion of knowledge, the results of section 4.3.4 suggest that they could have limited effects and might even be counterproductive. Indeed, intrinsic motives are a powerful driver of knowledge sharing among individuals, and, as suggested by Gneezy et al. (2011), the introduction of financial instruments might "crowd out" these

intrinsic motives in favor of financial calculations. The randomized experiment conducted by Squazzoni et al. (2013) on peer-reviewing is very telling in this context. The authors test the effectiveness of paying referees for doing reviews and find that the paid referees did a worse job than the ones who were not paid. In other words, while the financial compensation reduces the cost ( $C$ ) to perform a socially desirable action, it also reduces the moral motivation ( $\kappa$ ) and thus the value of the contribution ( $x$ ).

The low efficiency of market-based instruments poses the question of alternative mechanisms to increase the diffusion of knowledge. One option consists in reducing the individual cost of sharing without introducing financial considerations. For instance, as suggested by Cowan and Foray (1997) and confirmed by the survey results of Lee and Cole (2000), reducing the individual costs with better communication infrastructure can promote knowledge diffusion. In practice, facilitating the access to platforms where sharing knowledge is easy and cheap in time and effort is an effective way of reducing the individual cost without having to pay or tax anyone. Similarly, creating public data warehouses with well-designed instructions would greatly reduce the cost of sharing and increase contributions, as shown by the survey results of Kim and Stanton (2016). Finally, offering online tutorials explaining how to draft a Wikipedia article or how to organize and share databases is a relatively inexpensive approach to increase both the quality and the quantity of shared knowledge and data.

### **The role of awareness**

With the classical self-centered *homo oeconomicus* approach, since all individuals' decisions are taken based on cost and benefit considerations, being aware of the social impact of one's action and of the level of sharing in the population holds little importance. With *homo moralis* preferences, however, the perception each individual has of the size of the positive externality and of the level of sharing in the population greatly affects the sharing strategy. More specifically, as discussed in sections 3.2 and 3.3, if an individual perceives the social benefit as low (modeled as a low  $\xi$  in our setting), then her level of sharing will be lower than what it could be if she knew more about the impact of her actions. Similarly, the desire to share among the population would increase if individuals realize that the level of sharing in the population is higher than what they had previously envisioned.

Consequently, a relatively cheap and effective policy is to increase the perceived social benefit of sharing and to publicize sharing actions. For instance, putting banners on websites such as Wikipedia or data warehouses stating the number of daily contributors to the website would increase the awareness of the actual level of sharing in the population. Similarly, as shown by the experiment of Chen et al. (2018) on scientific experts contributing to Wikipedia, publicizing the impact of the shared knowledge on others increases the perception of the positive externality, and therefore enhances the sharing attitude.

The promotion of open source software and data would also benefit from a better understanding of licensing issues, especially among researchers. For example, training individuals and encouraging them to use copyleft<sup>54</sup> conditions could help the diffusion of knowledge in two ways. First, it would decrease the individual cost by curtailing fears of free-riding behaviors. Second, it would increase the perceived social benefit by ensuring that users are using it for personal benefits only. Finally, it can also increase the perceived level of sharing by reinforcing the visibility of the sharing community.

Furthermore, the use of nudges, i.e., suggestive messages without any financial or legal intervention (Thaler and Sunstein, 2008), putting forward the social benefit of sharing behavior and the level of sharing in the population can be a cheap but effective approach to increase the contribution level in the population. While classical economic policy would suggest that awareness on the social benefit is neutral to individual behavior, our model suggests it can be decisive to promote knowledge diffusion.

#### 4.4.3 Limitations and further work

While offering new insights on a central intrinsic motive of individuals to share knowledge and data, our model does not cover all the incentives that one might have for sharing. For instance, the model does not include motives such as status and recognition by peers (Gallus, 2017) and contribution to the direction of knowledge (Thompson and Hanley, 2018). Unfortunately, no model can integrate all empirically observed preferences, but we are confident that our setting offers a broad enough scope and a high explanatory power (Miettinen et al., 2020).

Although our analysis is focusing on individuals' sharing behaviors, our model could easily be applied to employees within organizations. Even if financial competition is the rule of functioning in the context of firms, "moral" actions could be an optimal solution to thrive in business. In fact, in a survey study on firm employees, Wasko and Faraj (2000) show that intrinsic motives in general and morality, in particular, have a key role to play in stimulating knowledge exchange both within and outside of the organization. At the level of the firm itself, having a "moral" approach to business integrating the outcome if all other firms acted similarly can provide a competitive advantage (Kopel et al., 2014; Kurataa and Van Longb, 2019). In practice, many firms in several sectors share their data with suppliers and competitors alike (Garry, 2009; Ghoshal et al., 2018). Some tech giants such as Facebook and Google are also sharing their source code and application programming interfaces (API) to facilitate the

---

<sup>54</sup> Copyleft is the practice of offering people the right to freely distribute and modify a piece of work, at the condition that the same rights are preserved in derivative works. See <https://www.gnu.org/copyleft/>.



deployment of their technologies and increase their quality (Bodle, 2011). An extension of our approach could be to evaluate the business success of firms if their optimal strategy relies on "moral" optimization instead of profit-maximization.

Finally, our model would benefit from empirical investigations estimating the various parameters used, such as the values of individual cost, positive externality, peer pressure, and morality. Obtaining precise estimates in different contexts is necessary to calibrate the model, and, in turn, to provide useful and specific policy recommendations.

## 4.5 Conclusion

Benefiting from effective knowledge sharing among individuals has massive economic consequences. For instance, when a pandemic, such as the Covid-19 outbreak, shakes our economies, fast and effective knowledge and data sharing among scientists all around the world can have a tremendous impact by reducing casualties and accelerating the way towards a vaccine (Johansson et al., 2018). More generally, increasing the pool of available knowledge can spur innovation and improve economic conditions in numerous ways. For instance, favoring knowledge diffusion is fundamental to enhance economic development and to reduce inequalities (De la Croix et al., 2018; Bell et al., 2019), while the development of artificial intelligence technologies can be much more socially beneficial if a widespread data-sharing culture is set (Cockburn et al., 2019). In this context, identifying the determinants of knowledge and data sharing is fundamental. This paper complements the existing literature on the subject by explaining why individuals share without any financial incentives and suggests some policies that would effectively enhance sharing.

Our study is particularly relevant in the ongoing debate about the implementation of open science policies. Several funding agencies in Europe and in the United States have recently expressed their desire to make all publications resulting from their funding freely available (Guedj and Ramjoué, 2015; Spichiger, 2018; Luna-Reyes and Najafabadi, 2019). Evaluating the incentives that most efficiently encourage scientists to share their work is essential for a faster transition towards more accessible knowledge.

This paper aims at bringing a new perspective on the individuals' motives behind knowledge and data sharing. We hope that it will open the way for a better consideration of morality as a factor influencing decision making, and help better understand the mechanisms to increase social welfare by favoring knowledge diffusion.

# Chapter 5 Machine learning in healthcare: Mirage or miracle for breaking the costs dead- lock?

*Disclaimer: This chapter draws on a working paper written in collaboration with my supervisor Dominique Foray. Dominique has been a mentor, an inspiration and an endless source of knowledge and ideas for me. It has been an honor and a true pleasure to collaborate with him, I hope it is only the first of many to come.*

## **Abstract**

The ageing population in all developed economies and the limited productivity characterizing the healthcare sector are leading to alarmingly increasing costs. The current rapid advances in machine learning (ML), a subfield of artificial intelligence (AI), offer new automation and prediction capabilities that could, if properly integrated, help address the healthcare costs deadlock. Are ML-driven solutions the appropriate ingredient to produce this necessary transformation, or are they condemned to face the same destiny as previous attempts to remodel healthcare delivery? This paper aims at bringing first elements to answer this question by providing both qualitative and quantitative evidence on the development of ML in healthcare and discussing the organizational and institutional conditions for the ML potential to be realized. Building on a novel search methodology for publications and patents in ML and on hospital surveys, our results reveal two major observations. On the one hand, while the publication rate in the field has tripled in the last decade, the level of patenting in ML applied to healthcare has so far been relatively low. This result has several potential explanations, such as the early stage of the technology, its rapid growth, and the emergence of new business models based on data accumulation and appropriation rather than patenting. On the other hand, the bulk of firms' publications are produced by IT firms rather than by companies in healthcare. This last observation seems to be driven by the disruptiveness of the new ML technology allowing the entry of new actors in healthcare. The technology producers benefit from their mastery of ML and the lack of investment and capabilities among health experts.

## 5.1 Introduction

Healthcare services are a vital component of all economies (Cutler and Richardson, 1998).<sup>55</sup> Nonetheless, partly due to an ageing population, the healthcare sector is currently facing numerous challenges such as growing care needs, higher societal expectations, and multimorbidity (Atun, 2015), leading to a significant surge in health expenditures<sup>56</sup> (Aizcorbe and Nestoriak, 2011; Cutler, 2017). The increase in costs is all the more alarming, considering that the sector has struggled to raise its productivity in the past (Kocher and Sahni, 2011; Baumol, 1993, 2012).<sup>57</sup> The increasing costs with little perspective of facing them with higher productivity lead the healthcare sector to a cost deadlock. In parallel, the development of Machine Learning (ML) technologies brings highly effective prediction capacities complementing human labor and creating new business opportunities with the potential to increase productivity (Brynjolfsson and Mitchell, 2017; Taddy, 2018). This powerful innovation is, therefore, naturally seen as a candidate to revolutionize healthcare practices. ML algorithms are already starting to provide medical applications in specialties as diverse as radiology, cancer research, and dermatology (Agrawal et al., 2019; Bibault et al., 2018; Haenssle et al., 2018). Numerous startups are building their business models around these applications, and even tech giants are showing a growing interest in applying their technologies in healthcare.<sup>58</sup> Historically, the integration of information and communication technologies (ICTs) in healthcare had a modest impact on the productivity of the sector (Lee et al., 2013), questioning the potential of ML to have a more impactful fate. However, ML has different characteristics, and the ICT revolution targeted operations that are very different from the ones ML technologies can impact (Webb, 2019). In this paper, we thus ask whether -contrasting with the weak impact of previous ICT applications- ML solutions can help address the healthcare cost deadlock and the framework conditions for achieving it.

ML has the properties of a general-purpose technology (GPT) with the potential to diffuse widely in various application sectors (Brynjolfsson et al., 2018). However, this potential can only be realized under a number of conditions on the providers of the technology (supply-side) and on the adopting sector (demand-side). Therefore, in our analysis, we use the GPT framework, which puts much emphasis on the centrality of co-invention, innovation complementarities, and externalities to explain the pattern of

---

<sup>55</sup> The healthcare sector is a major spending area for most developed economies accounting for 9% of the GDP of OECD countries on average in 2017 and 18% in the United States (OECD Health Statistics 2018). [https://stats.oecd.org/Index.aspx?DataSetCode=HEALTH\\_STAT](https://stats.oecd.org/Index.aspx?DataSetCode=HEALTH_STAT)

<sup>56</sup> See Cutler (2017) for a detailed discussion of healthcare costs growth in the United States since 1970.

<sup>57</sup> Identifying a lack of productivity growth characterizing the economics of healthcare, we do not include in our analysis the bio-tech and pharmaceutical industries, which constantly exhibit high productivity growth rates.

<sup>58</sup> <https://www.forbes.com/sites/forbestechcouncil/2020/01/30/big-techs-brewing-battle-over-healthcare-data/#29158c5d6b48>

diffusion in healthcare. Then, building on a powerful search methodology on patents and publications, we monitor the production and adoption of ML technologies by the healthcare sector. To have a more general picture, we complement these quantitative observations with survey data (on hospital transition towards digitalization) as well as qualitative case studies. These pieces of evidence suggest that ML can enhance productivity and offer new types of services and products for healthcare delivery. However, as a GPT, ML exhibits several particular features that makes this potential hard to realize.

First, as for any historical GPT, the bidirectional externalities between ML inventions and the development of applications (also called “co-invention”) as well as the externalities between early users and the next adopters within the specific application sector (e.g., healthcare) are considerable, which makes the economy far from a socially optimal rate of invention and co-invention.

Second, innovational complementarities seem to be particularly difficult to build between ML innovation and healthcare business model applications. Our analysis shows that the patenting rate in the field of ML remains rather low in comparison to a soaring publication rate. This result reflects the emergence of a new appropriability regime – in which building and preserving data advantages become more central than acquiring exclusive rights on technologies - which creates new business challenges for inventors in ML applications (e.g., in healthcare).

Third, the diffusion of ML innovation in healthcare is driven by the way innovational complementarities are built between the GPT and the sector-specific technologies, organizations, business models, and human capital. Healthcare is not an “easy” application sector for effective and rapid deployment of ML innovation, as shown by our survey on hospital digitalization.

However, a fourth feature is observable. It is relatively specific to ML as compared with other historical GPTs and deals with the fact that the GPT inventors in ML (i.e., tech companies) are very active not only in advancing the body of knowledge about the fundamental inventions in ML but also in developing new knowledge in the application sectors. This fact is probably the most interesting findings of our data analysis in terms of scientific publications. In other words, these big companies represent institutions, which are capable of internalizing the externalities from both GPT invention and the development of applications, reducing thereby the size of externalities and making the economy closer to an optimal rate of invention and co-invention. This trend seems to happen not only in the “easiest” application sectors (such as marketing and advertisement) but also in markets that matter for growth and social development– such as healthcare.

This situation is likely to entail many advantages: the entry of the big companies into the healthcare industry does represent an efficient way to internalize externalities, minimize market failures and generate a high rate of innovation in a socially desirable direction. However, this could mean a radical change in the division of inventive activities between the GPT inventors, the co-inventors in application sectors, and

the academic research. These institutional changes can have substantial economic and social effects. Finally, there are two issues raised by the entry of the GPT inventors in a few crucial sectors such as healthcare – the issue of concentration and competition, on the one hand, the issue of privacy on the other. Both issues are addressed at the end of the paper.

## 5.2 The innovative significance of ML in healthcare

This paper focuses on a particular aspect of innovation in healthcare: the development and deployment of ML applications to radically transform the processes of healthcare production and coordination and offer new types of services. As recently shown by Webb (2019), unlike software and robotic technologies, ML solutions are mainly directed at qualified tasks. This characteristic of ML gives it a higher potential than previous technologies to reduce the reliance on labor in healthcare. Moreover, provided that access to large amounts of data is possible, ML algorithms can produce high-quality predictions allowing new applications and business opportunities in various fields, including healthcare (Agrawal et al., 2018). ML has the properties of a general-purpose technology (GPT) with the potential to drive transformational changes in the hospitals and all healthcare services – to an extent not reached during the first phase of ICTs and computerization in healthcare (Sahni et al., 2017; Trajtenberg, 2018, Brynjolffson, Rock and Syverson, 2019; Cockburn et al., 2019). During the first phase of ICTs penetration in healthcare, the new technology was mainly used as transactional tools for billing, monitoring, and error checking and pressed against the existing old infrastructure – which limited the full realization of ICTs potentials and created new costs (Hendrich et al., 2008). Conversely, the second phase, the one of ML, can help transform the very way hospitals deliver medical care.

### 5.2.1 Productivity in healthcare

Many service sectors have witnessed continuous growth in their productivity levels, with output rising per person-hour and leading to a reduction of the labor share through technological changes and IT capital accumulation. Healthcare, however, has lagged behind (Chandra and Skinner, 2012; Bojke et al., 2017). The limited productivity growth in the sector does not imply that it is technologically inert. For several technological metrics such as the share of knowledge workers and capital renewal, the healthcare sector was at least as technologically active as, for instance, manufacturing good producers and progressive services sectors (Feldstein, 2017). The observable poorer productivity performance of healthcare is rather driven by two major economic characteristics of the sector: the centrality of labor and innovation market failures.

The first factor explaining the slow productivity growth in healthcare is the centrality of human labor in all operations. A study by the U.S. Bureau of Labor Statistics (BLS)<sup>59</sup> projects that, although the labor force participation rate will decline in the 2015-2024 decade, around 40% of the newly created jobs will be in healthcare. This trend is partly driven by increasing demand for healthcare services due to the accelerated ageing of the population<sup>60</sup> but, more importantly, to the hardly replaceable nature of human labor in healthcare services. In most sectors (both in industry and in services), human labor is primarily an instrument – “an incidental requisite for the attainment of the final product” (Baumol, 1967). Hence, the fact that the part of human labor in the final product or service (labor input coefficient) is decreasing does not change the evaluation of the quality of the goods or services by consumers. Most industrial sectors build on this property (“human labor is an instrument only”) to sharply increase labor productivity through technological progress and labor-capital substitution (Autor and Salomons, 2018). Conversely, in healthcare, human labor is not only a factor of production. Labor is an end in itself, and the quality of services is judged directly in terms of the amount of labor involved (Baumol, 1993). This situation does not necessarily imply that the potential for productivity increase is not significant. For instance, a considerable part of doctors’ and nurses’ working time is not spent with the patients but instead with documentation and other administrative tasks (Rao et al., 2017; Gardner et al., 2019). There is, therefore, a broad scope for productivity improvement through the integration of new processes and technologies and better organization (Rajkomar et al., 2019). Nevertheless, part of human labor is irreducible in healthcare, implying that the sector, by its very nature, will always exhibit a relatively high labor to capital ratio. One must bear this argument in mind when evaluating the role and place of ML in healthcare – as complementing (H-enhancing) rather than substituting (H-replacing) human labor (Trajtenberg, 2018; Agrawal et al., 2019).

The second reason behind the productivity gap in healthcare lies in the limited incentive for innovation. As discussed by Cutler (2011), the “lack of information and poor incentives” for entrepreneurs in the healthcare sector limits the development of new business models and organizational structures. So far, technological solutions have failed to fix these inefficiencies, and as a consequence, productivity growth has been hampered. There are, of course, massive innovation activities in healthcare as in any sector where intelligent people are learning by doing and thriving to solve problems through « user innovation» and numerous innovation activities are performed<sup>61</sup>. However, the sector seems to lack the right incentive

---

<sup>59</sup> See <https://www.bls.gov/opub/mlr/2015/article/overview-of-projections-to-2024.htm>

<sup>60</sup> See United Nations (2009), World Population Prospects.

<sup>61</sup> See for instance DeMonaco, Ali & Von Hippel (2006) and Oliveira et al. (2015) on user innovation in hospitals and other healthcare organizations.

structure to attract entrepreneurs with the ability to identify business opportunities and to try new business models. Usually, entrepreneurs are attracted by the expectation of capturing a significant fraction of the social value of the innovation (Chesbrough and Rosenbloom, 2002; James et al., 2013). So far, this condition has rarely been met in healthcare. The data-driven revolution of ML can offer new business models for startups and tech companies alike to help in reducing these market failures.

### 5.2.2 ML solutions to costly healthcare operations

The recent advances in ML open new promising possibilities in terms of automation and prediction. As defined by Brynjolfsson et al. (2018), ML allows to continuously learn from previously collected data to establish better quality and lower price predictions. This prediction specificity of ML makes it an ideal candidate for applications outside of the core field of computer science, as shown by Cockburn et al. (2019). Specifically, in the case of healthcare, ML could significantly reduce the cost of several services performed by physicians and thereby to curb health expenses by reducing processing time and increasing the availability of the medical cast for the patients (Yu et al. 2018; Rajkomar et al., 2019).

Agrawal et al. (2019) identify four direct effects through which prediction technologies like ML can affect labor: “Substituting labor by capital (H-replacing), automating decision tasks, enhancing labor tasks (H-enhancing) and creating new decision tasks.” Recent empirical results suggest that ML has the potential to augment labor productivity rather than replacing it (Bessen et al., 2018; Webb, 2019). This property of ML makes it a suitable candidate to improve quality and cut down operating costs. The exploitation of ML capabilities in healthcare has in fact already shown promising results for performing faster high-quality predictions of diseases, for establishing a precise diagnosis in a limited time and for efficiently selecting the optimal treatment for the patient:

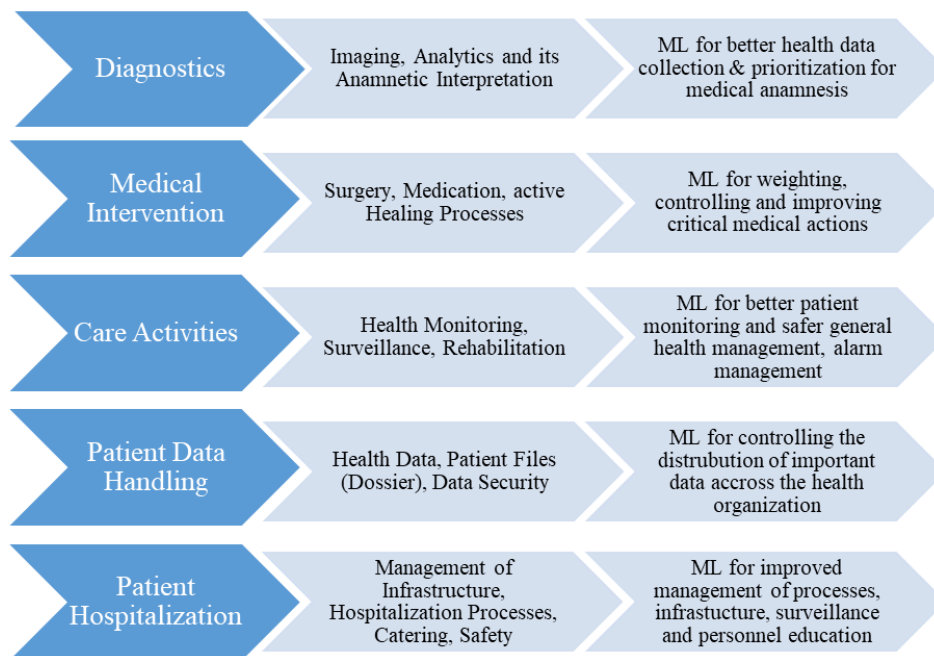
- In prediction: Haenssle et al. (2018) and Akselrod-Ballin et al. (2019) show that advances in deep learning (DL) -a subfield of ML- can sharply increase the rate of identification of cancers (skin cancer and breast cancer respectively) years before they become visible to a human eye.

- In diagnosis: De Fauw et al. (2018) train an algorithm on millions of medical images for fifty different retinal diseases and show that it reduces the time needed by the ophthalmologist to perform the diagnosis while increasing the success rate.

- In treatment selection: Bibault et al. (2018) and Sun et al. (2018) use ML algorithms to predict the response to different cancer treatments based on the characteristics of the patient and imagery results. The procedure drastically reduces wrong choices and costly, inefficient treatments while increasing the patient’s healing chances.

In other words, whether before a patient develops a disease (prevention), when feeling the first symptoms (diagnosis), or when choosing the way to treat it (treatment selection), ML techniques can reduce the time needed to perform medical actions and limit the error rate. Scott (2009) estimated that medical errors cost to around \$30 billion annually to the US economy, which represents a significant share of healthcare expenditures, notwithstanding health consequences for patients. Therefore, the first results already exhibited by ML algorithms in healthcare suggest that daily care operations performed by the medical staff can become faster and cheaper. More generally, beyond the direct provision of health services, the emerging technology revolution – including ML and big data - can offer a myriad of applications in production processes and care coordination. The figure below shows the potential of ML-based transformations and improvements on the map of hospital tasks and processes (figure 5.1).

**Figure 5.1: Hospital tasks and processes – where can ML help?**



Source – N.Bühler (Industry Relations Manager at ETH Zurich) – Elaborated for this project.

### 5.2.3 Diffusion of ML in Healthcare

As suggested by Brynjolfsson et al. (2018), ML “is potentially pervasive, improves over time, and can spawn complementary innovation,” meeting the Bresnahan and Trajtenberg (1995) criteria for a GPT. Empirical evidence about ML as a GPT is growing (Mihet and Philippon, 2019). For example, Cockburn et al. (2019) register the publication trend over time for three different AI fields: learning, robotics, and symbolic logic. For each field, they separate publications in computer science from publications in application fields. For the field of learning, their results suggest a sharp increase in publications that use ML in scientific fields outside computer science. This first evidence suggests that ML is a GPT, with the potential to



have a widespread impact on the economy, accelerating growth. Therefore, as a GPT (Bresnahan and Greenstein, 1997), ML is subject to two classical observations discussed in the following two paragraphs.

### **Characterizing healthcare as an application market for GPTs**

The success of ML in spurring innovation and producing economic value depends on four characteristics of the application sector (healthcare in our case): the aggregate demand, the benefits of the GPT compared to current solutions, the development cost, and the complementarities between the GPT and the application sector (Helpman and Trajtenberg, 1996). This framework developed by Helpman and Trajtenberg (1996) allows identifying the relative easiness of adoption of a GPT by a given sector. Based on this framework, how “easy” is healthcare as an application field for ML? First, the discussion above (sections 5.2.1 and 5.2.2) suggests that the demand for ML in healthcare is potentially very high (fulfilling the first characteristic) and that the solutions it offers should bring a significant added value compared to state of the art (fulfilling the second characteristic). As for the third characteristic, the development cost of ML solutions applied to healthcare will significantly depend on the accumulation of large training datasets allowing increasing returns to scale and reducing the average cost. Finally, the fourth characteristic, the complementarity of the ML GPT with the healthcare sector, will depend on the interaction between the producers of the technology and actors of healthcare adopting the new solution. More specifically, ML can only generate significant innovation in healthcare if entrepreneurs exploit the new technology, reconfigure it to create new business processes and co-invent applications with actors of the adopting sector. Hence, to profoundly transform the healthcare sector, ML technologies should enhance innovation complementarities and co-invention of solutions between both producers of ML technologies and final users in the healthcare sector. This approach is needed to create positive feedback loops and produce increasing returns to scale for both sectors (Bresnahan, 2010). The recently published *Nature Medicine* article (De Fauw et al. 2018) based on a joint research effort by teams from an emblematic Californian software firm and clinicians from a renowned British hospital gives an example of a co-invention offering cheaper, faster and better predictions of retinal diseases compared to well-established ophthalmologists.

### **Externalities in GPTs invention and diffusion**

Vertical improvement and horizontal diffusion of a GPT are governed by feedback and externalities (Bresnahan and Trajtenberg, 1995). The externalities among co-inventors are considerable. For any GPT, early adopters within one sector of application generate learning and informational effects and a range of effects on a better provision of many specific inputs, such as skills and specialized services (Goolsbee and Klenow, 2002; Goldfarb, 2005). In the case of ML, the learning effect is even more substantial.

Precisely, since the quality of the prediction depends on the amount of data generated, algorithms improve exponentially and generate increasing returns to scale to the producer even giving the possibility to “experiment” the technology at a lower cost (Hendel and Spiegel, 2014; Iansiti and Lakhani, 2020). Furthermore, the externalities between inventors of the GPT and co-inventors of applications are substantial. They can be both positive (network effect, learning benefits) and negative (coordination and transaction costs, human capital needs). Since not all these externalities are internalized by contract, the economy is far from a socially optimal rate of invention and co-invention. But here comes a « new » empirical observation: the new tech firms are active in advancing both the body of fundamental ML knowledge and the body of applied ML knowledge – providing a mechanism to internalize some of these externalities.

### **The disruptiveness of ML as a GPT**

ML, unlike previous GPTs, is highly disruptive. Precedent GPTs such as the computer were radical but not disruptive. By disruption, we mean that the current GPT inventors, use it to try to enter all application sectors markets.<sup>62</sup> While computer and software producers as GPT inventors at the computer age did not try (and had little opportunity) to enter healthcare, Alphabet (the parent company of Google) as a GPT inventor at the AI age is doing it (see section 5.3 for a more detailed discussion). The providers of ML technology have proven their ability to enter application fields that appeared initially out of their scope: Alphabet is already testing an autonomous car<sup>63</sup> (transport), Amazon has developed its own streaming service<sup>64</sup> (media), and Facebook has announced its desire to create a new currency<sup>65</sup> (finance). The intrinsic nature of data and the expertise tech companies have developed in training ML algorithms allow increasing returns to scale with the potential to overcome the classically limiting organizational costs. This new situation creates both advantages and disadvantages from a social welfare point of view. On the one hand, it might be an efficient way to internalize externalities (within the company producing the GPT and its applications), minimize market failures, and generate a rate of innovation closer to optimality. On the other hand, it raises concerns in terms of the direction of inventive activities – who is deciding about the direction within crucial application sectors such as healthcare?

---

<sup>62</sup> We follow here the approach described by Christensen (1997) suggesting that new *entrant* firms often use their mastery of a radically new technology to become active on an *established* market. In our case, tech companies are the entrant firms, using ML technology to be active on the healthcare market.

<sup>63</sup> <https://techcrunch.com/2019/08/23/watch-a-waymo-self-driving-car-test-its-sensors-in-a-haboob/>

<sup>64</sup> <https://www.amazon.com/primeinsider/video/prime-video-qa.html>

<sup>65</sup> <https://www.forbes.com/sites/bernardmarr/2019/10/07/facebook-blockchain-based-cryptocurrency-libra-everything-you-need-to-know/#33203f2d4d7a>

Having discussed the diffusion potential of ML, the goal of the following two sections is to bring new empirical evidence on the development of ML solutions in healthcare. Precisely, we provide metrics on ML-based innovations in healthcare, observe the entry of new actors, assess capacities and capabilities of the large healthcare organizations (hospitals) to learn and absorb these innovations, and assess the innovation capacity in this domain. To this end, while section 5.3 investigates the supply side (basic research, inventions, and the centrality of the search for new business models to capture the value of innovations), section 5.4 analyzes the demand side with the capacity of hospitals to generate and use digital innovations.

### 5.3 Supply-side: science, invention and business models

Building on a powerful search methodology developed by WIPO for patents and publications in ML (WIPO Technology Trends 2019), we monitor the production and adoption of ML technology by the healthcare sector using publication<sup>66</sup> and patent<sup>67</sup> data related to the subject of ML in healthcare. The search methodology we use presents several advantages. First, being developed in collaboration with AI experts, it is based on a large corpus of ML related terms capturing very specific ML technologies. More precisely, the search methodology we use includes practical ML techniques terminology such as, for instance, “AdaBoost,” “covariate shift,” and “perceptron”; and specific IPC codes that allow a right balance between exhaustivity and precision. With this, we add to the search method currently used in the literature on the subject, mainly using broader labels such as “deep learning,” “neural networks,” and “unsupervised learning” (Cockburn et al., 2019; Webb, 2019). Second, the search algorithm we use enables us to target a specific subfield of AI (ML) and apply it to a particular application sector (healthcare) with a high degree of precision. Finally, our approach was developed with the concern of capturing the time variation of the terminology, including the recent advances in the field of ML and excluding the terms that were considered as AI at some point in time and are somewhat obsolete nowadays. We thus collect data on the most recent advances in the field of ML, namely the technologies that are most relevant for healthcare applications. This data collection then allows us to estimate the production of ML scientific publications and patents in healthcare over time, by country, and by major actors such as firms and universities.

---

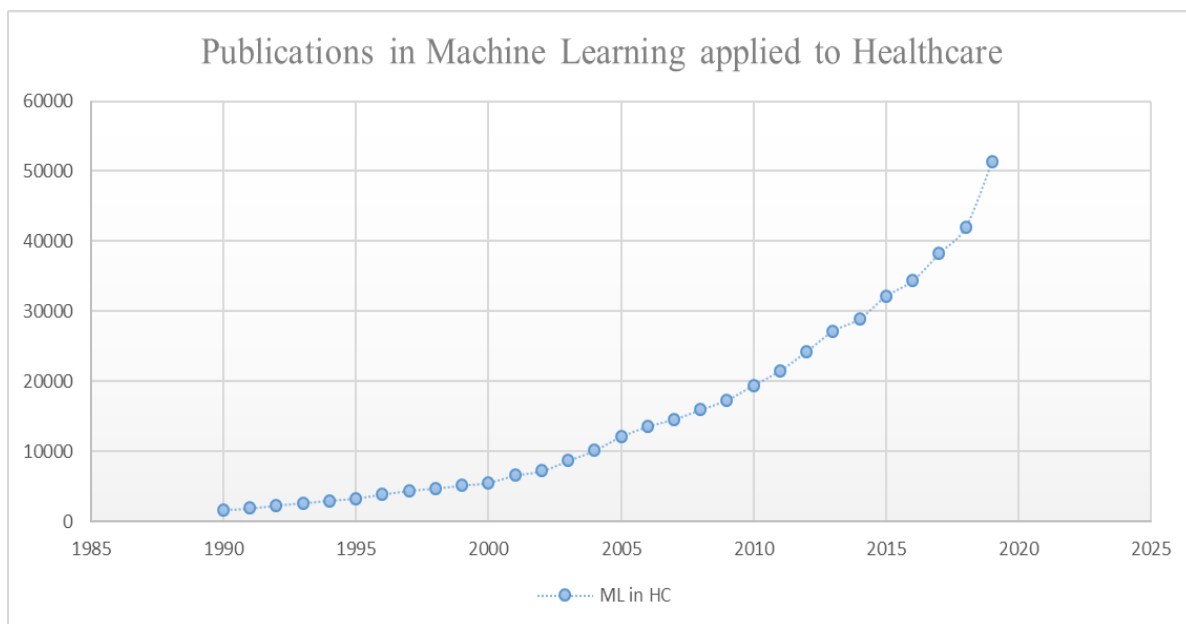
<sup>66</sup> We use a query searching for typical combinations of words associated with ML in the titles, keywords and abstracts of scientific articles based on the work of WIPO in their AI patent landscape (2019). See appendix A5.1 for more details.

<sup>67</sup> We use the WIPO query selecting patents based on the IPC codes associated with ML and then limit to healthcare related IPC codes. See appendix A5.1 for more details.

### 5.3.1 Evidence on patents and publications

The main finding concerning scientific articles is that the publication rate in the area is witnessing a consistent surge in the last ten years with a tripling of the number of publications between 2009 and 2019, which represents an average of 12% yearly over the period (see figure 5.2). During the same time span, the average growth of all scientific publications was three times smaller, with around 4% growth per year. More interestingly, in line with the findings of Cockburn et al. (2019), we observe that while the production of ML articles in computer science has slowed down around 2010, the publication rate in the application sector of healthcare has kept its exponential growth at a rate of around 12% yearly since 1990.

**Figure 5.2 - Evolution of ML publications in Healthcare 1990-2019 (World)**



Source – Scopus, authors' calculation (see Appendix A5.1)

Concerning the distribution of scientific publications in ML applied to healthcare, we observe that the trend is similar among the top five countries in terms of publications<sup>68</sup> and in most European countries (see figures 5.3 and 5.4). The most notable member of the top five is China with a scientific production that went from almost absent in the late 1990s to the second position in the world just behind the United States confirming the country's strategy aiming to be the world leader in AI by 2030 (Roberts et al., 2019).

---

<sup>68</sup> From largest to smallest: United States, China, United Kingdom, Canada, and Germany.

Figure 5.3 - ML publications in healthcare by country (World), 1990 – 2018

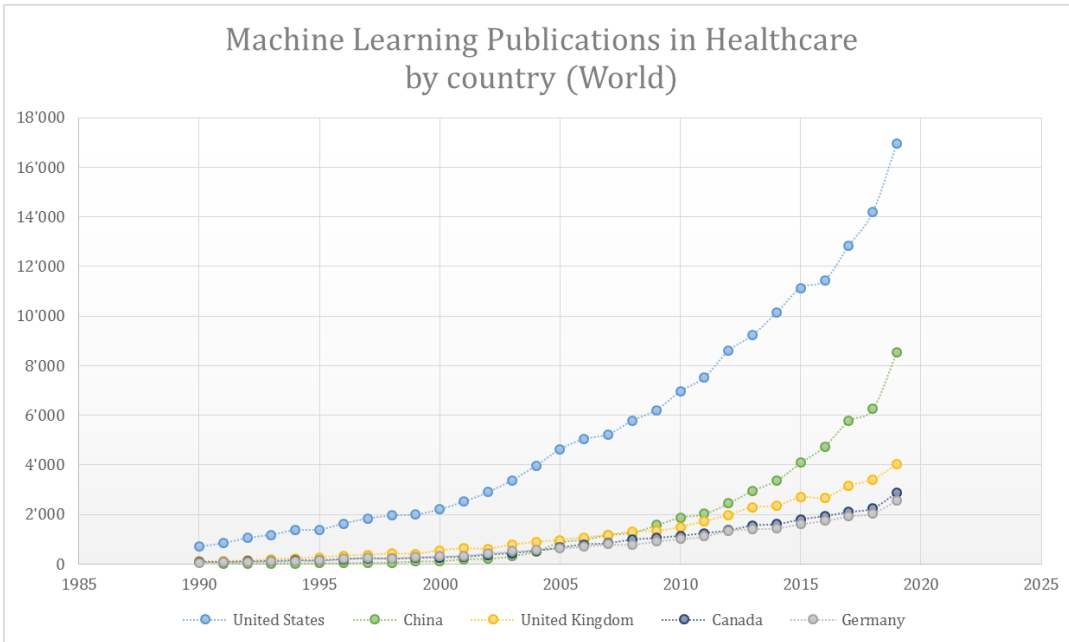
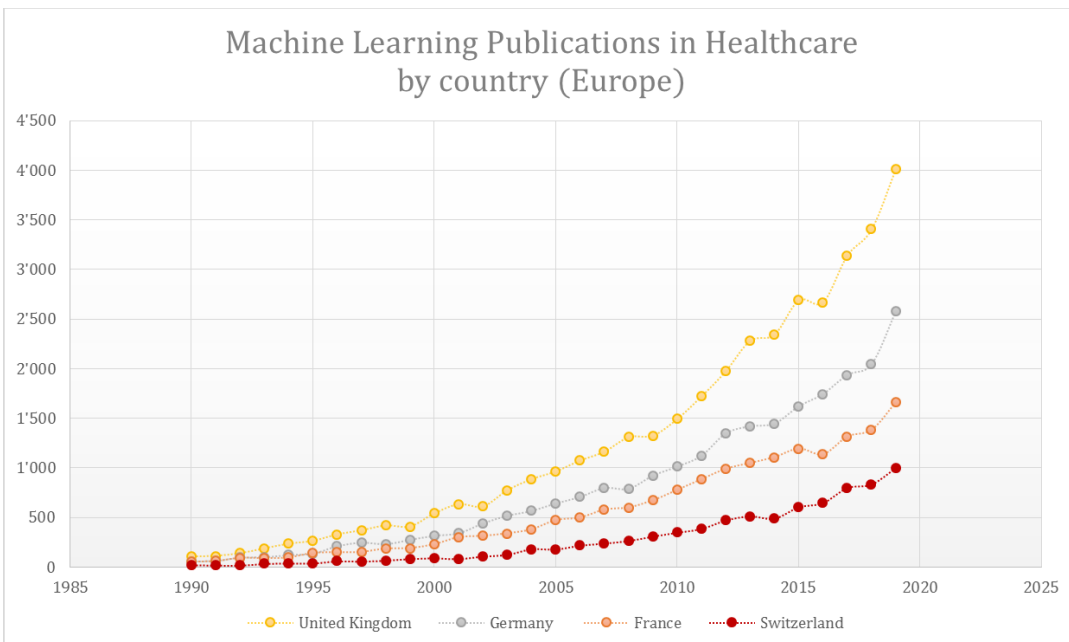


Figure 5.4 - ML publications in healthcare by country (Europe), 1990 – 2018



Source – Scopus, authors' calculation (see Appendix A5.1)

Looking at patents, we observe that the level of patenting in ML with an application in healthcare has so far been relatively low in all countries (see figures 5.5 and 5.6). This trend can be imputed to at least four reasons related to the technology as well as the institutional characteristics: i) the early stage of the technology, ii) the fast-changing technological environment, iii) the patenting legislation for algorithms, and iv) the development of new appropriation mechanisms.

A first reason for the low patenting rate is the low maturity of ML technology in healthcare. So far, ML models with effective results in healthcare are still at the stage of tests and development and thus lead to scientific publications rather than technologies mature enough to be patented. One would then expect that the scientific discoveries being published today only need time to lead to the innovations of tomorrow as it is the case for all types of technologies (Fukuzawa and Ida, 2016; Mukherjee et al., 2017, Poege et al., 2019). Second, ML technologies are characterized by a very fast-changing environment with new, more efficient algorithms continuously displacing old ones (Sabater et al., 2019). This quickly changing environment is also favored by a culture of open science in the field that encourages researchers and computer scientists to publish their results to continuously increase the efficiency of the models used (Badawi et al., 2014). In this context, going through the process of paying for a patent granting a long-term monopoly for a technology that is quickly rendered obsolete by a new model makes little sense. Third, from a legislation perspective, ML technologies are based on the development of algorithms on a computer and are therefore subject to a very stringent patenting law requiring the connection of the innovation to a physical device which makes it more cumbersome to patent new ML solutions compared to other technologies (Guntersdorfer, 2003). Finally, and probably most interestingly, the low patenting rate can also be imputed to the emergence of new business models that are not based on owning the intellectual property of the invention but rather on the capacity of firms to establish an advantage at an early stage in terms of data accumulation and appropriation. The characteristics of the ML technology makes it only useful and superior if it is trained on massive amounts of high-quality data. Hence, holding the data is at least as important as mastering ML algorithms to ensure a competitive advantage on the technology. Similar low patenting strategies lead by new appropriation mechanisms have been recently identified in other high-tech industries such as mobile application development (Miric et al., 2019).

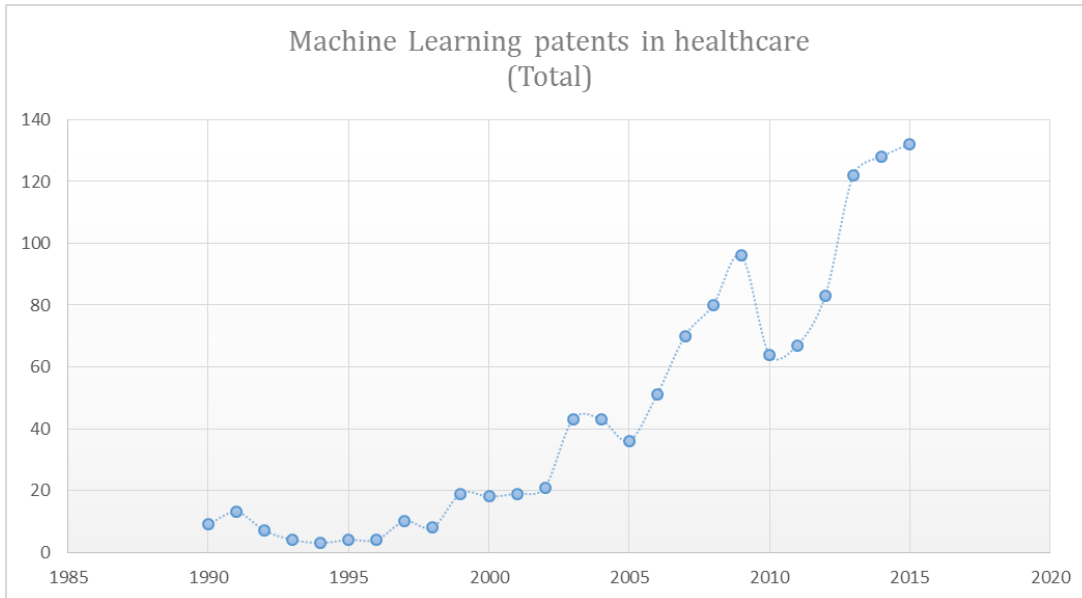
Another salient observation is the rather higher patenting rate of the United States compared to other countries. This observation can be imputed to the higher rate of patenting by the country in general as it also leads the ranking of patenting for all patents confounded. However, while the share of US-based patents in the international production represents an average of 13.2% over the 2006-2015 period, the share of US-based patents for ML in healthcare over the same period is three times higher (42.6%)<sup>69</sup>. This excessive domination of the United States in ML patenting is a sign of the leading technological role played by the United States in terms of ML applied to healthcare. However, it also reflects the less stringent law concerning software patents in the country. More specifically, the 1981 U.S. Supreme Court's decision (followed by the 1994 Federal Circuit's ruling) – allowing the patentability of software that produced “a useful,

---

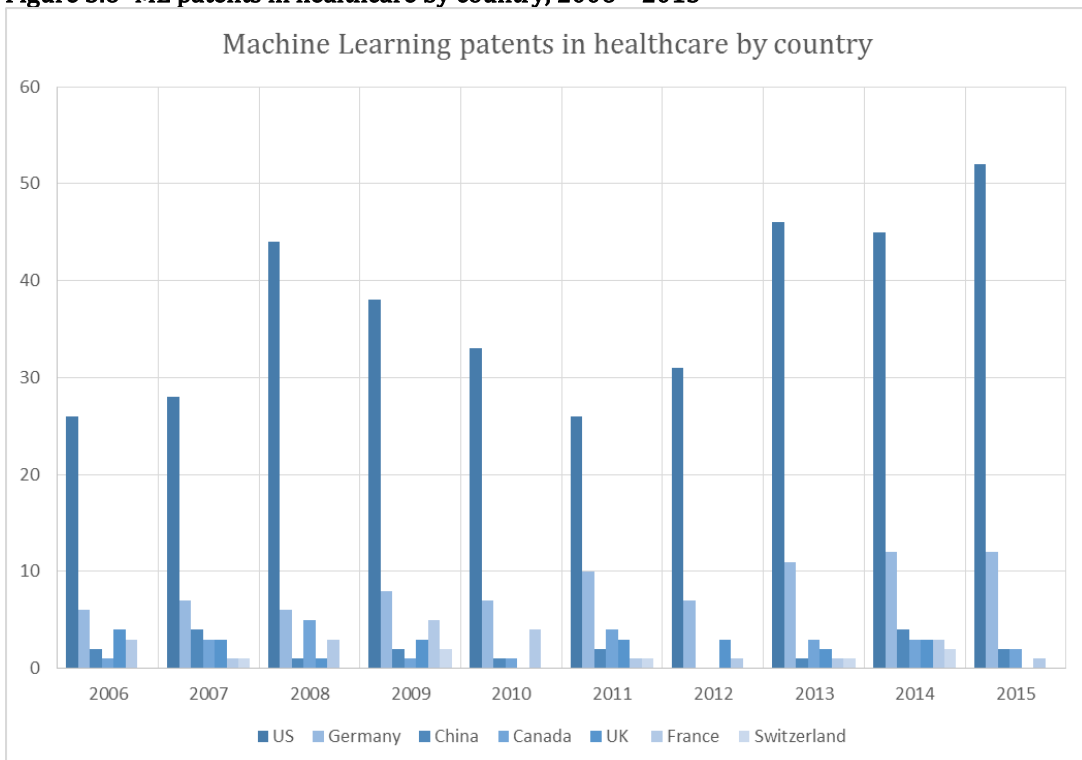
<sup>69</sup> Note that the share of US-based patents in ML in general (not only applied to healthcare) is similar around 40%.

concrete, and tangible result, even if that result was just on a computer screen” (Bessen and Hunt, 2007) – make it easier for US-based inventors to patent their inventions in ML than their European counterparts.

**Figure 5.5 - Evolution of ML patents in Healthcare 1990-2017 (World)**



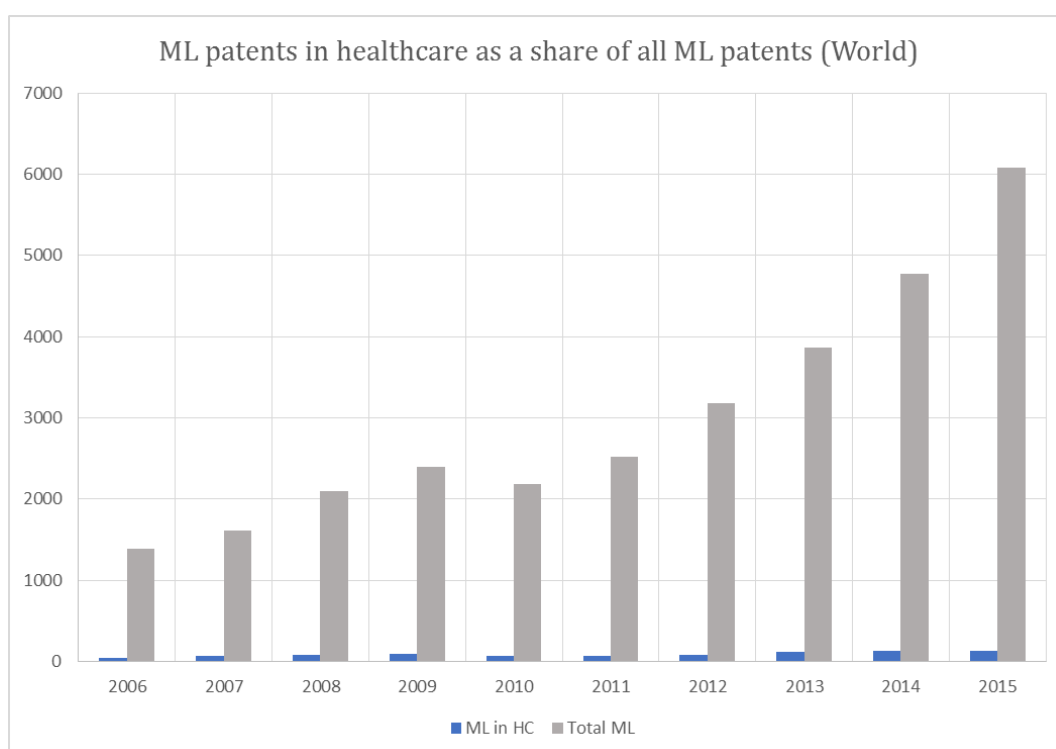
**Figure 5.6- ML patents in healthcare by country, 2006 – 2015**



Source – Patstat, authors’ calculation (see Appendix A5.1)

We also find that the share of patented healthcare applications relative to patents on all ML inventions is still rather tiny. Precisely, over the period 2006-2015, we observe that ML patents with applications in healthcare represented, on average, 3.2% of all ML patents granted. As a comparison, over the same time span, for all patents (not limited to ML), the share of healthcare patents represented around 4.6% (40% higher). These numbers suggest that the healthcare sector is not yet a leading field for ML applications (compared to computer sciences and other services and industries) (see figure 5.7).

**Figure 5.7- Evolution of the share of ML patents in healthcare (HC) relative to the total of ML patents- World, 2006 – 2015. (Total ML = 30'087, Total ML HC = 893, share of HC=3%)**



Source – Patstat, authors' calculation (see Appendix A5.1)

### 5.3.2 Main actors and entry of tech firms

Concerning the producers of ML related scientific knowledge in healthcare, both academic (universities) and industrial actors (firms) have a role to play. On the one hand, retrieving the publications by the firms with the heaviest R&D investment<sup>70</sup>, we observe that some traditional pure players of the software industry mastering ML techniques (such as IBM or Google) are producing scientific articles applied

---

<sup>70</sup> We consider the most active firms in terms of R&D investment in 2018 based on the 2018 EU Industrial R&D Investment Scoreboard (<https://iri.jrc.ec.europa.eu/scoreboard18.html>) published by the Joint Research Centre of the European Commission on a yearly basis since 2004.



to the healthcare sector (see table 5.1 for details). These scientific publications with health-related applications represent non-negligible shares of the global portfolio of articles in ML produced by these companies (around 9% on average), suggesting that they are performing diversification strategies towards healthcare. The application of their mastery of ML techniques in healthcare shows that these actors perceive a real potential for leveraging their know-how and data in healthcare. On the other hand, considering the publication activity by the top universities of ML articles applied to healthcare, we observe that they are larger providers than firms on average. This last observation is not surprising since the production of science is their primary role. Nevertheless, the share of the ML related articles that are applied to healthcare is not significantly higher than for companies with rates lower than 50% on average, suggesting that the focus of the research is mainly centered on core activities in computer science (see table 5.2 for details).

**Table 5.1- Publications in ML healthcare by IT companies with at least one scientific article in ML in healthcare.**

Software and IT									
Company	Country	Industry	Pubs total	Pubs ML	Pubs ML HC	Ratio ML/Pubs	Ratio MLHC/Pubs	Ratio MLHC/ML	
SAMSUNG ELECTRONICS	South Korea	Electronic & Electrical Equipment	37967	1602	476	4,22%	1,25%	29,7%	
IBM	US	Software & Computer Services	104307	6521	441	6,25%	0,42%	6,8%	
PHILIPS	Netherlands	General Industrials	38000	1148	416	3,02%	1,09%	36,2%	
SIEMENS	Germany	Electronic & Electrical Equipment	39634	1915	291	4,83%	0,73%	15,2%	
MICROSOFT	US	Software & Computer Services	30767	6233	168	20,26%	0,55%	2,7%	
ALPHABET	US	Software & Computer Services	9645	2810	117	29,13%	1,21%	4,2%	
INTEL	US	Technology Hardware & Equipment	26694	1501	42	5,62%	0,16%	2,8%	
TENCENT	China	Software & Computer Services	983	564	14	57,38%	1,42%	2,5%	
AMAZON.COM	US	General Retailers	656	209	9	31,86%	1,37%	4,3%	
HUAWEI	China	Technology Hardware & Equipment	6397	748	7	11,69%	0,11%	0,9%	
FACEBOOK	US	Software & Computer Services	1706	561	7	32,88%	0,41%	1,2%	
ORACLE	US	Software & Computer Services	1414	78	5	5,52%	0,35%	6,4%	
APPLE	US	Technology Hardware & Equipment	1062	108	3	10,17%	0,28%	2,8%	
<b>Total</b>			<b>261265</b>	<b>22396</b>	<b>1520</b>	<b>8,57%</b>	<b>0,58%</b>	<b>6,8%</b>	

Source: EU Scoreboard 2018 and Scopus

**Table 5.2- Publications in ML in Healthcare (HC) by Top 20 Universities in Computer Science**

QS rank in CS (2018)	Institution	Country	Pubs MLHC	Ratio ML/Total Pubs	Ratio MLHC/ML
1	Massachusetts Institute of Technology (MIT)	US	458	1.40%	13.41%
2	Stanford University	US	699	0.99%	23.50%
3	Carnegie Mellon University	US	298	4.37%	7.74%
4	University of California Berkeley (UCB)	US	223	0.91%	9.89%
5	University of Cambridge	UK	253	0.60%	16.54%
6	University of Oxford	UK	398	0.56%	26.91%
7	Harvard	US	1034	0.46%	44.24%
8	Ecole Polytechnique Federale de Lausanne (EPFL)	Switzerland	123	1.44%	11.63%
9	Swiss Federal Institute of Technology (ETH Zurich)	Switzerland	213	1.06%	13.99%
10	National University of Singapore (NUS)	Singapore	343	1.82%	12.18%
11	University of Toronto	Canada	497	0.58%	26.84%
12	Nanyang Technological University (NTU)	Singapore	393	3.24%	11.81%
13	Princeton University	US	110	0.96%	10.19%
14	University of California Los Angeles (UCLA)	US	499	0.54%	30.24%
15	Imperial College London	UK	423	0.66%	23.23%
16	Tsinghua University	China	246	2.32%	5.38%
17	University of Washington	US	447	0.75%	21.39%
18	Columbia University	US	435	0.70%	25.20%
19	Peking University	China	223	1.14%	11.82%
20	New York University (NYU)	US	311	0.67%	25.49%

Source: QS ranking in Computer Science 2018 and Scopus

Firms have been known to publish for several reasons, such as accessing external knowledge, retaining human capital, and signaling<sup>71</sup>. Moreover, since we observe a low level of patenting in the field of ML applied to healthcare, we believe that the rate and direction of scientific publications give the best indicator of firms' strategies. Therefore, to support the disruptive GPT hypothesis on horizontal propagation in the case of ML (see section 5.2.3), we proceed to a more systematic measure of the entry of big companies – such as IBM, Microsoft or Google – into healthcare as a sector of application for ML by monitoring their scientific publications. We observe that firms such as IBM or Microsoft are publishing works in ML with applications in healthcare at levels similar to some of the top universities (see tables 5.1 and 5.2). Furthermore, comparing tables 5.1 and 5.3 suggests that the bulk of firms' publications in ML for healthcare is produced by computer sciences and software firms (the producers of the ML GPT) rather than by firms specialized in healthcare (the application sector). This last observation suggests that ML is a disruptive GPT with inventors of the technology entering application fields rather than merely providing the new technology to traditional actors.

**Table 5.3- Publications in ML Healthcare by healthcare companies with at least one scientific article in ML in healthcare.**

Healthcare									
Company	Country	Industry	Pubs total	Pubs ML	Pubs ML HC	Ratio ML/Pub	Ratio MLHC/Pub	Ratio MLHC/M	
MEDTRONIC PUBLIC LIMITED	Ireland	Health Care Equipment & Services	3147	134	87	4,26%	2,76%	64,9%	
THERMO FISHER SCIENTIFIC	US	Health Care Equipment & Services	4948	155	65	3,13%	1,31%	41,9%	
FRESENIUS	Germany	Health Care Equipment & Services	1095	58	52	5,30%	4,75%	89,7%	
BAXTER INTERNATIONAL	US	Health Care Equipment & Services	2297	38	30	1,65%	1,31%	78,9%	
BECTON DICKINSON	US	Health Care Equipment & Services	1314	34	25	2,59%	1,90%	73,5%	
BOSTON SCIENTIFIC	US	Health Care Equipment & Services	1091	31	24	2,84%	2,20%	77,4%	
EDWARDS LIFESCIENCES	US	Health Care Equipment & Services	540	15	11	2,78%	2,04%	73,3%	
CARL ZEISS	Germany	Health Care Equipment & Services	2571	42	10	1,63%	0,39%	23,8%	
STRYKER	US	Health Care Equipment & Services	305	8	4	2,62%	1,31%	50,0%	
OLYMPUS	Japan	Health Care Equipment & Services	1074	22	4	2,05%	0,37%	18,2%	
MCKESSON	US	Health Care Equipment & Services	137	4	3	2,92%	2,19%	75,0%	
ZIMMER BIOMET	US	Health Care Equipment & Services	395	4	2	1,01%	0,51%	50,0%	
<b>Total</b>			<b>15372</b>	<b>407</b>	<b>228</b>	<b>2,65%</b>	<b>1,48%</b>	<b>56,0%</b>	

Source: EU Scoreboard 2018 and Scopus

As a further analysis, we run two regressions to evaluate the role of each type of actor in the integration of ML in healthcare. To do so, we collect ML publications in healthcare for the main firms publishing and for top universities, their citation count, and year of publication.

On the one hand, we observe that publishing more in the field of ML, in general, is associated with a higher publication rate in ML applied to healthcare (table 5.4a). This observation is in line with the disruptiveness of ML with inventors of the technology investigating the application field. On the other hand, we find that universities' articles in ML for healthcare are, on average, more impactful than firms'

<sup>71</sup> See Camerani et al., 2018 for a full review of the literature on the subject.

publications. Specifically, table 5.4b results suggest that universities' ML publications in healthcare receive, on average, around three citations per year more than publications by firms. However, considering the recent time trend (with the dummy post\_2009 equal to 1 if the publication is after 2009<sup>72</sup>), we observe that firms' ML publications in healthcare partly offset this difference to be on average only slightly less cited as university publications in the last ten years. This last observation indicates that firms' scientific production in ML for healthcare is more and more considered as impactful by actors of the field.

**Table 5.4a**

	(1) OLS Pubs in ML Healthcare
Pubs in ML	0.066*** (0.016)
Pubs Total	0.001 (0.000)
Pubs in healthcare	0.002*** (0.001)
University	-60.762 (37.603)
Constant	-82.721* (42.922)
Dummy Industry	Yes
Dummy Country	Yes
Observations	61
R2	0.966

In reporting the statistical significance of the coefficients, we apply the standard thresholds, i.e., \*\*\* p<0.01, \*\* p<0.05, \* p<0.1.

**Table 5.4b**

	(1) OLS Citations per year
Firm	-2.613*** (0.934)
Firm*Post_2009	1.759* (1.004)
Post_2009	-2.941 (4.269)
Constant	8.307** (4.228)
Dummy Year	Yes
Observations	8,901
R2	0.021

In reporting the statistical significance of the coefficients, we apply the standard thresholds, i.e., \*\*\* p<0.01, \*\* p<0.05, \* p<0.1.

---

<sup>72</sup> The choice of the year (2009), follows the observation of Cockburn et al. (2019) suggesting that the year 2009 represents a turning point in terms of publications in AI with learning functions.

### 5.3.3 New business models

#### **Towards a new appropriability regime**

While labor-saving solutions show high potential for cost reduction (see section 5.2), the most promising contribution for ML in healthcare is probably in opening the way for new business opportunities spurring innovation in the field. In fact, at the AI age, the technology structure of innovative ideas is partially shifting from technologies defined as processes and products to technologies defined as services (Varian, 2018). The ML innovation model brings new modes of doing business in terms of legal protection and commercialization. As suggested by the results of section 5.3.1, patents do not appear to be very well adapted, and the old model of patent and licensing out seems hardly applicable to ML solutions. The main reason for this is that the value lies in the data, i.e., the capacity of the startup or the inventor to secure ‘indefinite’ access to extensive medical databases through strong connections with data sources (whether the hospital or directly the patient). Based on this data access and use, the innovation then consists in selling a unique and exclusive service. In this context, building and preserving a data advantage becomes a central business strategy (Cockburn et al., 2019). By effectively generating and protecting clinical and medical data, firms can erect a “data-driven barrier to entry” that can ensure market dominance over at least the medium term. The need for high-quality data as a requirement for implementing ML innovations also creates complementary business opportunities for data production, collection, and analysis.

In terms of data generation, companies now have a higher incentive to develop devices that compile medical information about individuals both efficiently and reliably. For instance, in the United Kingdom, Butterfly, a company already valued at \$1.25 billion, aims at performing ultrasound scans on smartphones, to facilitate the generation of ultrasound images<sup>73</sup>. As for data collection and analysis, an interesting new business model is the one of SOPHiA Genetics, a Swiss-based start-up developing and operating a worldwide network of 850 leading hospitals around the world. All client hospitals of the company are connected through one data infrastructure. The hospitals provide clinical and medical data (genetic codes, medical images) and benefit in return from various types of health tech (ML-driven) applications. The core of the transactions between SOPHiA Genetics and the hospitals, as mediated by the platform, consists in selling unique services involving ML applications that analyze data and images and combine them with biological and clinical information to predict diseases’ evolution and support clinical decisions and strategies for individual patients. The transactional structure of the service is designed under the SaaS

---

<sup>73</sup> <https://www.forbes.com/sites/jamesmauroo/2019/09/27/butterfly-a-125-billion-healthtech-company-launches-new-ultrasound-technology-in-the-uk/#29f08359a3f0>

(Software as a Service) model based on a Pay per use schema. Because the success of the business model is strongly determined by network rules (direct and indirect network externalities), the critical strategic assets to sustain it involve: i) the capacity of the firm to continuously increase the size of the hospital network, ii) the algorithmic and data science capacity and iii) the marketing/selling capacity. In these models, patents are not central to capture a significant fraction of the social value of innovation (patents on new algorithms are primarily defensive). Following a business model akin to the one of SOPHiA Genetics, several startups are trying to leverage their technological advantage and convince health institutions to share their data to create predictive analytics innovations. In the United States, similar business models have been developed as the Federal Drug Administration (FDA) has granted its first acceptance to Arterys, a company using ML technologies for heart problems diagnosis.<sup>74</sup> In the EU also, the newly adopted General Data Protection Regulation (GDPR) does not seem to slow down the entry of new firms aiming to produce, collect or exploit healthcare data as VC investment in healthcare has been multiplied by six between 2015 and 2019<sup>75</sup>. More generally, entry of innovators (entrepreneurs, startups, and disruptive outsiders) is becoming a crucial determinant of innovation and transformation, and it was not the case in the past.

### **Conflicting with an open-science policy and academic culture**

One of the reasons for the steady increase in the number of publications in ML is the evolution of the academic culture in computer science towards data, code, and results sharing (Badawi et al., 2014). Accordingly, new business models based on building a competitive advantage through data accumulation and appropriation are likely to conflict with the increasingly dominant academic trend about open science and open data. A key challenge for firms to leverage their data advantage will thus reside in their ability to overcome this cultural clash in their cooperation with ML specialists.

For instance, we met an ETHZ (Zürich Institute of technology) Professor in bioinformatics who undertakes fundamental research in computer science, ML, and analytics. He develops methodological tools and works on a broad range of applications, including healthcare. He consistently works under the open-source regime and publishes everything – codes, prototypes, and the data used to do the research. Securing data access and the reproducibility of the research findings are core principles for him, and so is the case of most of his fellow academic computer scientists. However, this approach creates issues for collaboration with companies: Company X, developing a fertility tracker, contacted him to develop better tools to

---

<sup>74</sup> <https://www.forbes.com/sites/bernardmarr/2017/01/20/first-fda-approval-for-clinical-cloud-based-deep-learning-in-healthcare/#5b8cfd8c161c>

<sup>75</sup> <https://2019.stateofeuropeantech.com/chapter/investments/article/investment-industry/>

improve their prediction performance but refused any open source and data conditions – so the collaboration never started.

There is still a long way to go in terms of institutional and legal creativity to create a zone of compatibility between the emerging academic rules ensuring data accessibility and research findings reproducibility and the new business models – in short, to reconcile long-term welfare effect of open data and short-term welfare effect of delivering better products to the market.

## 5.4 Demand-side: Hospitals transformation and needs

The diffusion of ML technology in the healthcare sector does not depend only on the supply of new solutions offering efficiency gains and performance. A necessary component for the successful implementation of these innovations is the ability of the demand-side to adopt the new technology, help improve it, and even co-invent new, more adapted innovations. The main actor on the demand side for innovation applied to healthcare is hospitals. There are numerous challenges for hospitals to be capable of fully benefiting from the positive spillovers generated by the development of ML applications. Challenges involve upgrading skills, forming new management capabilities, investing in IT infrastructures, building cooperation with fundamental computer science, implementing radical organizational changes (including, for example, the increasing substitution of in-patient treatments by out-patient treatments and the adoption of a more patient-centric approach); as well as adapting healthcare processes to new business models.

In a 2017 McKinsey study to evaluate the predisposition of economic sectors for adopting AI technologies<sup>76</sup>, healthcare appeared as one of the laggards in terms of readiness for implementing AI technologies. Hence, to better understand the readiness of hospitals and their strategy for integrating and developing ML solutions, we ran a survey among the biggest hospitals in Switzerland in terms of daily operations. Our survey data provides numerous information about the state of the art and the future commitments of hospitals towards the digital revolution in healthcare (figures 5.8 to 5.15<sup>77</sup>).

We monitor that considerable investments have been made for the last years and will continue in the coming years. Most institutions expect to invest more in applications (data warehouses, AI, big data) than in computing capacities (figure 5.8). Most institutions are running projects on AI and big data (figure

---

<sup>76</sup><https://www.mckinsey.com/~media/mckinsey/industries/advanced%20electronics/our%20insights/how%20artificial%20intelligence%20can%20deliver%20real%20value%20to%20companies/mgi-artificial-intelligence-discussion-paper.ashx>

<sup>77</sup> For clarity reasons, figures 5.8 to 5.15 are reported in appendix A5.2.

5.9), and they start mostly in domains such as logistics and administrations, which are “less risky” than medical applications (figure 5.10). It also seems that doctors, other clinical personals, and patients are not highly involved in AI projects (figure 5.11).

The in-house development of data science capacities remains very low (figure 5.12), and cooperating with fundamental research institutions seems to be a less popular mechanism for knowledge sourcing than contracting with IT suppliers (figure 5.13). Generating inventive ideas through startups is a source that is not exploited (figure 5.14). A final element is that very few hospitals seem to have an innovation strategy for digital transformation (figure 5.15). We also noted this absence of a clear innovation strategy in the open comments section: “*All ingredients are there, but what is the vision, the strategy? How and why developing predictive analytics – the transformation of care production and coordination to minimize cost, increase productivity, forge new business models, provide new services?*” (a respondent of the survey).

Any new generic technology in its early stage – like ML – actively raises the problems of innovational complementarities (section 5.2.3). In our case, fully exploiting the ML potential to improve healthcare production, coordination and management will require significant changes in the organization, human capital composition as well as the diffusion of a new «epistemic culture» within hospitals. Most hospitals have started to fill adjustment, restructuring, and implementation gaps, but there is still a long way to go. A key takeaway from our qualitative investigation within hospitals is that even when the investment efforts and the strategic decisions are in favor of practical testing and use of ML solutions, having access to qualified human capital capable of implementing these solutions is very hard. As stated by Gofman and Jin (2019), human capital is essential for AI-driven innovation, which increases the scarcity of qualified experts and raises their price. The value of AI technologies being very high for large IT and software companies, they attract most of the engineering force with attractive financial conditions (Rock, 2019; Alekseeva et al., 2020) even drying out the academic market of qualified AI professors (Gofman and Jin, 2019). In healthcare, it seems that the value hospitals can draw from ML engineers is not yet perceived as high enough to compete for this valuable human capital. However, having access to qualified human capital is a necessary condition for leveraging ML and big data disruptive potential (Alekseeva et al., 2020). The difficult access to the necessary expertise for health institutions thus appears as a significant limitation to a proper integration of ML solutions. Our survey in hospitals suggests that the talent necessary for implementing ML algorithms and the investments in intangible assets are far from ready for creating innovative solutions. A key challenge is to form individuals capable of bridging the nexus between the technical knowledge of ML and the highly demanding and specific medical knowledge, which requires changing the way individuals are trained in higher education institutions (Trajtenberg, 2018).

Furthermore, the kind of problems ML can tackle requires both a deep understanding of medical notions and strong algorithmic knowledge. This complexity is visible in the most promising projects of ML in healthcare, usually bringing together ML engineers with medical experts. The work of Akselrod-Ballin et al. (2019), for instance, required a joint effort of IBM engineers working with health specialists from the department of imaging of the Assuta Medical Centers to produce link healthcare records with mammograms results and improve breast cancer prediction. More generally, the ability to leverage the potential of ML solutions in healthcare will heavily depend on the successful implementation of collaborative strategies between the providers of the technology and adopters (Trajtenberg, 2018, Varian, 2018). These collaborations will be necessary at the data collection part -with health experts and patients helping in the constitution of actionable databases. At the algorithm development stage, collaborations will be useful when ML specialists need the knowledge of medical experts to tackle the critical problems and properly train algorithms.

## 5.5 Discussion

ML-based innovation in healthcare corresponds to the story of highly disruptive innovations in a highly regulated sector. ML has the properties of a GPT progressively diffusing in several application sectors. In healthcare, it has the potential to offer promising solutions to help alleviate the healthcare cost burden and improve care delivery (see section 5.2). Nevertheless, healthcare is a peculiar field with strong organizational rigidities and institutional blockages that have limited the impact of previous GPTs on the sector. The innovation rate in the sector is far from optimal and suffers from a lack of organizational capacities that have curbed the impact of past GPTs. Our empirical observations suggest, however, that ML is enabling disruptiveness by tech companies considering applications in healthcare. This evolution offers an unforeseen opportunity to reduce externalities and spur innovation in the field, but it also raises new concerns and limitations.

### 5.5.1 The disruptiveness of ML as a solution to spur innovation in healthcare

The empirical observations of section 5.3 suggest that the supply side is experiencing a new phenomenon: the emergence of the tech giants such as IBM and Alphabet (parent company of Google) as a potentially dominant innovator in healthcare (see table 5.1 on scientific publication by companies). For a range of ML applications (requiring strong data science skills and engineering capacities and for which the business model is relatively clear and “easy”), Alphabet, Amazon, and Apple in the US and Tencent in China have the capacity to achieve the innovation in a faster and more cost-effective way than academic research laboratories. These companies thus need to acquire good quality healthcare data to train ML algorithms and develop useful innovations (Taddy, 2018), and they are already doing so. Alphabet has sealed an



agreement with Ascension<sup>78</sup>, the second-largest hospital system in the United States, to exploit the data for improving the predictive capacity of ML algorithms with healthcare applications. Similarly, Amazon has an agreement with a Boston based medical center<sup>79</sup> to train ML algorithms, Tencent is collaborating with the Chinese Academy of Medical Sciences to predict cervical cancers<sup>80</sup>, while Apple is leveraging its own produced Apple watch data to predict heart diseases<sup>81</sup>. Tech companies also acquire or develop in-house healthcare firms to increase their knowledge of the field and improve their algorithms.<sup>82</sup> All these examples give anecdotal evidence of the intention of the providers of ML to enter an application market rather than simply sell the service to adopters. This disruptiveness of ML as a GPT is rendered possible by the increasing returns to scale of the technology. In fact, for firms providing historical GPTs such as the steam engine, electricity, or computers, it was too expensive (in terms of human capital formation, transaction costs, and management costs) and too hazardous to expand to application sectors such as healthcare or transport. ML is different in this concern as the steeper learning curve, and the possibility to test new opportunities with data now outweigh the classical costs of expansion (Iansiti and Lakhani, 2020). This characteristic enables ML as a highly disruptive GPT, which entails substantial welfare implications. On the one side, the entry of tech firms in the field permits the internalization of externalities. It reduces transaction costs between inventors and co-inventors, which increases the effective adoption of ML technologies in healthcare. This internalization of externalities by tech firms offers an effective path to tackle one of the main historical drags on innovation in healthcare. However, the welfare advantages offered by the disruption of tech firms in healthcare raises new issues that we discuss in the following section.

### 5.5.2 Limitations and framework conditions

Although having a promising potential, the implementation of ML technologies in healthcare poses numerous challenges for private companies and policymakers alike.

#### **Competition policy**

We have seen that, at least so far, there seems to be a large number of firms in healthcare innovation attempting to take advantage of ML through the accumulation and preservation of data assets. This new business model arising from ML innovation poses additional threats and adds to the rising concerns

---

<sup>78</sup> <https://www.nytimes.com/2019/11/11/business/google-ascension-health-data.html>

<sup>79</sup> <https://aws.amazon.com/fr/blogs/machine-learning/improving-patient-care-with-machine-learning-at-beth-israel-deaconess-medical-center/>

<sup>80</sup> <https://www.bloomberg.com/press-releases/2019-05-22/tencent-miying-launches-ai-supported-auxiliary-diagnostic-system>

<sup>81</sup> <https://www.apple.com/newsroom/2017/11/apple-heart-study-launches-to-identify-irregular-heart-rhythms/>

<sup>82</sup> Alphabet, for instance, already owns three healthcare companies (Verily, Cityblock Health and DeepMind Health) and a biotech company (Calico). <https://www.economist.com/business/2018/02/03/apple-and-amazons-moves-in-health-signal-a-coming-transformation>

on competition policy (Gutiérrez and Philippon, 2017, 2018). Autor et al. (2017) show that “superstar firms” are continuously increasing their market shares and outperforming their rivals. The ability to collect, organize, and analyze massive amounts of data with ML algorithms is a central reason for their increased market power (Iansiti and Lakhani, 2020). The capability of these firms to enter new markets like healthcare, and their intention to do so, increases concerns about their excessive market domination and led some economists to call for new antitrust policies (Grullon et al., 2017; Shapiro, 2019). The knowledge of customers that tech companies develop with their mastery of ML algorithm and data allows them to spread over various sectors, but it also grants them a power that makes it difficult for any competition to emerge effectively. Unlike the music or the streaming industry where tech companies have built substantial market shares, the healthcare sector holds more stakes. Consequently, the reaction of legal authorities and their attitude towards the acquisition of healthcare firms, like the acquisition of Fitbit by Alphabet<sup>83</sup> for accessing health data, will shape the evolution of the innovation in the sector. In parallel, as suggested by (Cockburn et al. 2019), encouraging larger accessibility of data beyond the borders of these giant firms is crucial. Most developed countries have, accordingly, conducted public initiatives to help in the construction of large enough databases that could be accessible for all healthcare providers. Examples of these initiatives include the Cancer Genome Atlas, aiming to register large scale imaging data on cancers (Tomczak et al., 2015) in the US, the UK Biobank, an open-access database, registers detailed information on over 500,000 patients on “a wide range of health-related outcomes” (Sudlow et al., 2015) in the UK and, the *Health Data Hub*<sup>84</sup> centralizing health data from patients, insurance companies, and hospitals in France.

### **Social acceptance and technical threats**

The entry of tech giants also poses privacy and societal acceptance concerns (Acquisti et al., 2016). In a 2018 survey on 4,000 US adults on digital health adoption, Day and Zweig (2018) find that only 11% of respondents are willing to share their data with tech companies compared to 72% with their physician and around 50% with their health insurance company. New entrants in the healthcare market will have to overcome this strong defiance barrier by customers if they want to succeed in being impactful players in the sector, especially in a context of increased data privacy legislation in the EU and the US<sup>85</sup>. Beyond the societal acceptance problem, the implementation of ML technologies has to overcome increasing technical threats in terms of data protection and exploitation. On one side, the collection of ever-

---

<sup>83</sup> <https://www.reuters.com/article/us-usa-antitrust-congress/key-antitrust-lawmaker-frustrated-with-googles-fitbit-deal-idUSKBN1XN2WY>

<sup>84</sup> <https://www.health-data-hub.fr/>

<sup>85</sup> See for instance the National Conference of State Legislature for a summary of the evolution of the legislation on data policy: <https://www.ncsl.org/research/telecommunications-and-information-technology/consumer-data-privacy.aspx>

growing amounts of health data to feed ML algorithms increases security threats and risks of hacking, especially with the high sensitivity of health data (Abouelmehdi et al., 2017). On the other side, since most ML algorithms use image analysis to produce predictions, the system is also vulnerable to adversarial attacks, i.e., small perturbations to the data, not perceivable by the human eye that can completely fool ML algorithms (Kurakin et al., 2016; Huang et al., 2017). Increasing security measures and innovating in the tools to cope with these technical challenges is also a condition for the success of ML innovations, especially in a sector like healthcare where stakes are very high. In fact, most profitable innovations with ML technologies have so far been in consumer-oriented services, advertising, and marketing (Bresnahan, 2019). One of the reasons for this confinement to specific sectors is that stakes associated with errors in those sectors are low. While in marketing applications an error of the ML algorithm simply leads to a failed selling attempt, in healthcare (and in transport), it can have tragic consequences. Hence, to be more widely socially accepted, ML innovations in healthcare will have to prove higher reliability than in the classical sectors where it so far thrived.

Both social acceptance problems need certification and regulation measures to reduce vulnerability and improve the confidence of patients in the new processes. To date, there is still considerable uncertainty about future regulation and future processes of certification in the domain of medical device innovation involving data collection and analysis and predictive analytics. The current disclaimer- «do not use it for clinical decision.» -which is applied to several innovations of this type does not help marketing. However, the recent US FDA approval for clinical cloud-based deep learning in healthcare<sup>25</sup> shows that the institutional processes on this matter are progressively being established.

## 5.6 Conclusion

ML and big data can provide many potential solutions for fixing operational inefficiencies in the organization of healthcare, generate smarter processes of healthcare provision and coordination and therefore create new sources of productivity increase. That is why the healthcare system is today at a crossroads. It is benefitting from a steady supply of scientific knowledge and academic skills, from innovative institutions and institutional networks in the field of ML and big data healthcare applications as well as from a significant number of entrepreneurial initiatives, developing and testing new business models. Such strong supply dynamics should translate into the generation and diffusion of numerous innovative solutions to many healthcare coordination and delivery problems. However, the system is not yet fully ready on the demand side to harness the fruit of the emerging revolution. It is clear that, with notable exceptions, hospitals have not yet achieved the digital transition, particularly in the domain of using big data and predictive analytics to support medical and clinical processes as well as to make better operational decisions. Because ML is a GPT, innovational complementarities between the development of new

applications and their adoption in critical settings (hospitals) will be central to realize the potential of ML and big data in terms of productivity. The impact of these innovations will depend eventually on systemic changes in the sector, involving: i) professional development of various categories of healthcare workers who are now low-skilled and unprepared for the new technology; ii) strategic move towards digitalization in healthcare organization and this includes, in particular, building computing facilities and acquiring the adequate scientific skills in data analytics; iii) cultural change at the level of practitioners and patients. We find here a familiar “GPT story,” according to which leveraging GPT productivity potential (such as ML and big data) requires significant changes in organization and human capital composition which take time and represent high adjustment costs, implementation and restructuring lags but if well done could create productivity gains across several sectors. Nevertheless, the disruptiveness of the ML technology, the importance of technical skills and access to data, all give a more central role to pure players of ML in the story with both positive and negative welfare consequences. This new story is likely to create a longer time lag in healthcare than in most other service sectors because of the specificities of the sector that render all revolutions lengthy and costly.



# Chapter 6 Conclusion

## 6.1 Scope and findings

Our societies are currently facing many large-scale challenges that will undoubtedly shape the years to come (Foray et al., 2012). These “grand challenges” include issues as diverse as climate change, growing inequalities, ageing populations, and now pandemics. Although very different in nature and origin, all these challenges would greatly benefit from the creation and diffusion of knowledge, enabling societies to cope with these rising threats. This situation calls for the design of adequate policy schemes and institutions for incentivizing knowledge creation and dissemination. The empirical evidence and the theoretical discussions exhibited in this thesis aim to inform this policy design.

The necessary innovations in energy, health, and computer sciences alike are increasingly relying on scientific knowledge (Fleming et al. 2019, Bikard and Marx, 2019). The first two essays of this thesis evaluate the factors enhancing the production and diffusion of high-quality scientific research. The two essays (chapters 2 and 3) find that current research funding schemes are effectively promoting the production and dissemination of scientific knowledge in quality and quantity. However, several welfare gains could be attained by promoting the participation in funding schemes and by properly calibrating team compositions. By theoretically discussing the contribution to publicly available knowledge, the third essay (chapter 4) offers insights that could promote the diffusion of knowledge to a larger crowd, and thus reduce inequalities in access to knowledge. Finally, the last essay (chapter 5) evaluates the potential of machine learning as a solution to improve healthcare provision by assessing the conditions for an integration of the technology in healthcare services and the challenges it poses.

## 6.2 Future developments

The questions explored in this dissertation open the way for several future studies on the institutions governing knowledge production and diffusion and on the incentive mechanisms that can drive economic agents towards better generation and dissemination of new ideas. This last section describes four such projects which I am currently working on but are not formally part of this thesis.

The first two essays of this dissertation explore two different phases of the process of applying to a research grant: chapter 3 focuses on the phase preceding the selection when the application is crafted, while chapter 2 focuses on the outcomes after the selection is decided. However, in order to

properly evaluate funding systems, one has to consider the selectivity patterns during the evaluation process. In fact, there are growing concerns about a Matthew effect within the scientific community (Perc, 2014; Azoulay et al., 2014), with reputed researchers concentrating all the funds. In parallel, concerns exist about excessive use of bibliometrics in evaluations and the selection for funding (Stephan et al., 2017, Ayoubi et al. 2019a). It is thus vital to understand who has higher chances of being awarded with funds and what characteristics increase winning chances. In the current format of competitive research financing, two sets of factors affect the reception of funds by researchers: the profile of the applicant and the content of the research proposal. Two papers currently in progress aim at studying these two sets of factors. The first one - currently under review in the *Journal of Economic Behavior and Organization* - evaluates the effect of the past profile of a researcher on the probability of applying for funds and of receiving a grant. The second one uses detailed data on applicants to the Sloan Research Fellowships to evaluate the characteristics of proposals that increase funding success.

Applying for funds is a costly and time-consuming activity for researchers. Peer-reviewing also takes a considerable share of researchers' time with little apparent benefit. Therefore, building on the measure developed in chapter 3 of this thesis, I am currently working on evaluating the learning benefits of reviewing a research paper for the referee. The project is exploiting data from a randomized control trial (RCT) in which reviewers are assigned randomly to papers to review, which helps in solving the selection bias issue that would arise in a standard setting.

Finally, coping with climate change is a challenge that requires the contribution of all individuals to be properly effective (Schelling, 1996). Hence, considering the insights of chapter 4, the diffusion of knowledge about climate risks seems like an effective mechanism to raise awareness in the population about environmental issues. Higher awareness can then be decisive to promote environmentally friendly behavior. In a project currently underway, we are considering the implications of integrating moral preferences in the design of optimal environmental policy.

The goal of these projects is further to explore the determinants of knowledge production and diffusion and offer useful insights for policymaking in the field. My current work has found applications in healthcare and environmental policy, but it could be extended to other objectives such as improving working conditions and reducing inequalities.

# References

## References for chapter 1:

- Alger, I. and Weibull, J. W. (2013). Homo moralis—preference evolution under incomplete information and assortative matching. *Econometrica*, 81(6):2269–2302.
- Arrow, K. J. (1962). Economic Welfare and the Allocation of Resources for Invention. In *The Rate and Direction of Inventive Activity: Economic and Social Factors: A Conference of the Unive* (pp. 609-626).
- Ayoubi, C. and Thurm, B. (2018). Exploring the diversity of preferences: Is a heterogeneous population evolutionarily stable under assortative matching?
- Azoulay, P., Graff Zivin, J. S., Li, D., & Sampat, B. N. (2019). Public R&D investments and private-sector patenting: evidence from NIH funding rules. *The Review of economic studies*, 86(1), 117-152.
- Baumol, W. J. (2012). *The cost disease: Why computers get cheaper, and health care doesn't*. Yale university press.
- Boh, W. F., Ren, Y., Kiesler, S., & Bussjaeger, R. (2007). Expertise and collaboration in the geographically dispersed organization. *Organization Science*, 18(4), 595-612.
- Boudreau, K. J. and Lakhani, K. R. (2015). “open” disclosure of innovations, incentives and follow-on reuse: Theory on processes of cumulative innovation and a field experiment in computational biology. *Research Policy*, 44(1):4–19.
- Breschi, S., & Lissoni, F. (2001). Knowledge spillovers and local innovation systems: a critical survey. *Industrial and corporate change*, 10(4), 975-1005.
- Brynjolfsson, E., & Mitchell, T. (2017). What can machine learning do? Workforce implications. *Science*, 358(6370), 1530-1534.
- Cassiman, B., & Veugelers, R. (2006). In search of complementarity in innovation strategy: Internal R & D and external knowledge acquisition. *Management Science*, 52(1), 68–82.  
<https://doi.org/10.1287/mnsc.1050.0470>
- Cohen, W. M., & Levinthal, D. A. (1990). Absorptive capacity: A new perspective on learning and innovation. *Administrative science quarterly*, 128-152.
- Cowan, R., David, P. A., & Foray, D. (2000). The explicit economics of knowledge codification and tacitness. *Industrial and corporate change*, 9(2), 211-253.
- Cutler, D. M. (2011). Where are the health care entrepreneurs? The failure of organizational innovation in health care. *Innovation Policy and the Economy*, 11(1), 1-28.
- Dasgupta, P., & David, P. A. (1994). Toward a new economics of science. *Research Policy*, 23(5), 487-521.



- David, P. A. (1993). Intellectual property institutions and the panda's thumb: patents, copyrights, and trade secrets in economic theory and history. *Global dimensions of intellectual property rights in science and technology*, 19, 29.
- Debackere, K., & Veugelers, R. (2005). The role of academic technology transfer organizations in improving industry science links. *Research Policy*, 34(3), 321-342.
- Feldman, M. P., & Kogler, D. F. (2010). Stylized facts in the geography of innovation. In *Handbook of the Economics of Innovation* (Vol. 1, pp. 381-410). North-Holland.
- Foray, D. (2004). *Economics of knowledge*. MIT press.
- Gallus, J. (2017). Fostering public good contributions with symbolic awards: A large-scale natural field experiment at Wikipedia. *Management Science*, 63(12):3999-4015.
- Gertler, M. S. (2003). Tacit knowledge and the economic geography of context, or the undefinable tacitness of being (there). *Journal of economic geography*, 3(1), 75-99.
- Grossman, G. M., & Helpman, E. (1994). Endogenous innovation in the theory of growth. *Journal of Economic Perspectives*, 8(1), 23-44.
- Gush, J., Jaffe, A., Larsen, V., & Laws, A. (2018). The effect of public funding on research output: the New Zealand Marsden Fund. *New Zealand Economic Papers*, 52(2), 227-248.
- Ioannidis, J. P. (2011). Fund people not projects. *Nature*, 477(7366), 529-531.
- Jacob, B. A., & Lefgren, L. (2011). The impact of research grant funding on scientific productivity. *Journal of public economics*, 95(9-10), 1168-1177.
- Jaffe, A. B. (1986). *Technological opportunity and spillovers of R&D: evidence from firms' patents, profits and market value* (No. w1815). National bureau of economic research.
- Jaffe, A. B. (1989). Real effects of academic research. *The American economic review*, 957-970.
- Jaffe, A. B., Trajtenberg, M., & Henderson, R. (1993). Geographic localization of knowledge spillovers as evidenced by patent citations. *The Quarterly Journal of Economics*, 108(3), 577-598.
- Jones, B. F., Wuchty, S., & Uzzi, B. (2008). Multi-university research teams: Shifting impact, geography, and stratification in science. *Science*, 322(5905), 1259-1262.
- Kocher, R., & Sahni, N. R. (2011). Hospitals' race to employ physicians—the logic behind a money-losing proposition. *N Engl J Med*, 364(19), 1790-1793.
- Krugman, P. R. (1991). *Geography and trade*. MIT press.
- Myerson, R. B. (1999). Nash equilibrium and the history of economic theory. *Journal of Economic Literature*, 37(3), 1067-1082.

- Nelson, R. R. (1959). The simple economics of basic scientific research. *Journal of Political Economy*, 67(3), 297-306.
- Romer, P. M. (1986). Increasing returns and long-run growth. *Journal of Political Economy*, 94(5), 1002-1037.
- Romer, P. M. (1990). Endogenous technological change. *Journal of Political Economy*, 98(5, Part 2), S71-S102.
- Stephan, P. E. (1996). The economics of science. *Journal of Economic Literature*, 34(3):1199–1235.
- Stephan, P. E. (2010). The economics of science-funding for research. *International Centre for Economic Research Working Paper*, (12).
- Stephan, P. E. (2012). *How economics shapes science* (Vol. 1). Cambridge, MA: Harvard University Press.
- Stiglitz, J. E. (1999). Knowledge as a global public good. *Global public goods*, 1(9), 308-326.
- Tirole, J. (2017). *Economics for the common good*. Princeton University Press.
- Williamson, P. J., & De Meyer, A. (2012). Ecosystem advantage: How to successfully harness the power of partners. *California management review*, 55(1), 24-46.
- World Intellectual Property Organization. (2019). *WIPO Technology Trends 2019: Artificial Intelligence*. World Intellectual Property Organization.
- Wuchty, S., Jones, B. F., & Uzzi, B. (2007). The increasing dominance of teams in production of knowledge. *Science*, 316(5827), 1036-1039.
- Xu, L., Nian, T., and Cabral, L. (2020). What makes geeks tick? A study of stack overflow careers. *Management Science*, 66(2):587–604.

## References for chapter 2:

- Aghion, P., Bloom, N., Blundell, R., Griffith, R., & Howitt, P. (2005). Competition and innovation: An inverted-U relationship. *The quarterly journal of economics*, *120*(2), 701-728.
- Angrist, J. D., & Pischke, J. S. (2008). *Mostly harmless econometrics: An empiricist's companion*. Princeton university press.
- Arora, Ashish, and Alfonso Gambardella, "The impact of NSF support for basic research in economics," *Annales d'Economie et de Statistique* 79/80 (2005), 91-117.
- Arrow, K. J. (1962). Economic welfare and the allocation of resources for invention. In *The Rate and Direction of Inventive Activity: Economic and Social Factors: A Conference of the Unive*, pages 609-626.
- Ayoubi, C., Pezzoni, M., & Visentin, F. (2017). At the origins of learning: Absorbing knowledge flows from within the team. *Journal of Economic Behavior & Organization*, *134*, 374-387.
- Azoulay, P., Graff Zivin, J. S., & Manso, G. (2011). Incentives and creativity: evidence from the academic life sciences. *The RAND Journal of Economics*, *42*(3), 527-554.
- Azoulay, P., Stuart, T., & Wang, Y. (2014). Matthew: Effect or fable?. *Management Science*, *60*(1), 92-109.
- Azoulay, P., Zivin, J. S. G., Li, D., & Sampat, B. N. (2015). *Public R&D Investments and Private-sector Patenting: Evidence from NIH Funding Rules* (No. w20889). National Bureau of Economic Research.
- Carayol, N., & Lanoë, M. (2017). *The Impact of Project-Based Funding in Science: Lessons from the ANR Experience* (No. 2017-04). Groupe de Recherche en Economie Théorique et Appliquée.
- Chavarro, D., Tang, P., & Rafols, I. (2014). Interdisciplinarity and research on local issues: evidence from a developing country. *Research Evaluation*, *23*(3), 195-209.
- Chubin, D. E., Hackett, E. J., & Hackett, E. J. (1990). *Peerless science: Peer review and US science policy*. Suny Press.
- Cole, J. R., & Zuckerman, H. (1984). The productivity puzzle. *Advances in Motivation and Achievement. Women in Science*. JAI Press, Greenwich, CT.
- Cumming, G. (2013). *Understanding the new statistics: Effect sizes, confidence intervals, and meta-analysis*. Routledge.
- Cumming, G. (2014). The new statistics: Why and how. *Psychological science*, *25*(1), 7-29.
- Defazio, D., Lockett, A., & Wright, M. (2009). Funding incentives, collaborative dynamics and scientific productivity: Evidence from the EU framework program. *Research Policy*, *38*(2), 293-305.
- Etzkowitz, H. (2003). Research groups as 'quasi-firms': the invention of the entrepreneurial university. *Research Policy*, *32*(1), 109-121.

- Fox, M. F. (1983). Publication productivity among scientists: A critical review. *Social studies of science*, 13(2), 285-305.
- Graves, N., Barnett, A. G., & Clarke, P. (2011). Funding grant proposals for scientific research: retrospective analysis of scores by members of grant review panel. *Bmj*, 343, d4797.
- Gush, J., Jaffe, A. B., Larsen, V., & Laws, A. (2015). *The Effect of Public Funding on Research Output: the New Zealand Marsden Fund* (No. w21652). National Bureau of Economic Research.
- Huber, G. P. (1991). Organizational learning: The contributing processes and the literatures. *Organization science*, 2(1), 88-115.
- Ioannidis, J. P. (2011). Fund people not projects. *Nature*, 477(7366), 529-531.
- Jacob, B. A., & Lefgren, L. (2011). The impact of research grant funding on scientific productivity. *Journal of public economics*, 95(9-10), 1168-1177.
- Jaffe, A. B. (2002). Building programme evaluation into the design of public research-support programmes. *Oxford Review of Economic Policy*, 18(1), 22-34.
- Lane, J., & Bertuzzi, S. (2011). Measuring the results of science investments. *Science*, 331(6018), 678-680.
- Lee, S., & Bozeman, B. (2005). The impact of research collaboration on scientific productivity. *Social studies of science*, 35(5), 673-702.
- Long, J. S. (1992). Measures of sex differences in scientific productivity. *Social Forces*, 71(1), 159-178.
- Mairesse, J., & Pezzoni, M. (2015). Does gender affect scientific productivity?. *Revue économique*, 66(1), 65-113.
- Mora, R., & Reggio, I. (2019). Alternative diff-in-diffs estimators with several pretreatment periods. *Econometric Reviews*, 38(5), 465-486.
- Rafols, I., Leydesdorff, L., O'Hare, A., Nightingale, P., & Stirling, A. (2012). How journal rankings can suppress interdisciplinary research: A comparison between innovation studies and business & management. *Research Policy*, 41(7), 1262-1282.
- Ruben, A. (2017). Another tenure-track scientist bites the dust. *Science*, 361(6409), 801.
- Stephan, P. E. (1996). The economics of science. *Journal of Economic Literature*, 34(3):1199-1235.
- Stephan, P. E. (2010). The economics of science-funding for research. *International Centre for Economic Research Working Paper*, (12).
- Uzzi, B., Mukherjee, S., Stringer, M., & Jones, B. (2013). Atypical combinations and scientific impact. *Science*, 342(6157), 468-472.
- Wooldridge, J. M. (2012). *Introductory econometrics: A modern approach*. Cengage Learning.

### References for chapter 3:

- Adams, J. D., Black, G. C., Clemmons, J. R., & Stephan, P. E. (2005). Scientific teams and institutional collaborations: Evidence from US universities, 1981–1999. *Research Policy*, *34*(3), 259-285.
- Agrawal, A., Cockburn, I., & McHale, J. (2003). *Gone but not forgotten: Labor flows, knowledge spillovers, and enduring social capital* (No. w9950). National Bureau of Economic Research.
- Agrawal, A., Kapur, D., & McHale, J. (2008). How do spatial and social proximity influence knowledge flows? Evidence from patent data. *Journal of urban economics*, *64*(2), 258-269.
- Apestequia, J., Azmat, G., & Iriberry, N. (2012). The impact of gender composition on team performance and decision making: Evidence from the field. *Management Science*, *58*(1), 78-93.
- Argote, L., & Miron-Spektor, E. (2011). Organizational learning: From experience to knowledge. *Organization science*, *22*(5), 1123-1137.
- Bercovitz, J., & Feldman, M. (2011). The mechanisms of collaboration in inventive teams: Composition, social networks, and geography. *Research Policy*, *40*(1), 81-93.
- Bikard, M., Murray, F., & Gans, J. S. (2015). Exploring trade-offs in the organization of scientific work: Collaboration and scientific reward. *Management science*, *61*(7), 1473-1495.
- Börner, K., Contractor, N., Falk-Krzesinski, H. J., Fiore, S. M., Hall, K. L., Keyton, J., ... & Uzzi, B. (2010). A multi-level systems perspective for the science of team science. *Science Translational Medicine*, *2*(49), 49cm24-49cm24.
- Burt, R. S. (2004). Structural holes and good ideas. *American journal of sociology*, *110*(2), 349-399.
- Campbell, T. A., & Campbell, D. E. (1997). Faculty/student mentor program: Effects on academic performance and retention. *Research in higher education*, *38*(6), 727-742.
- Cummings, J. N., & Kiesler, S. (2008, November). Who collaborates successfully? Prior experience reduces collaboration barriers in distributed interdisciplinary research. In *Proceedings of the 2008 ACM conference on Computer supported cooperative work* (pp. 437-446).
- Defazio, D., Lockett, A., & Wright, M. (2009). Funding incentives, collaborative dynamics and scientific productivity: Evidence from the EU framework program. *Research Policy*, *38*(2), 293-305.
- Ding, W. W., Levin, S. G., Stephan, P. E., & Winkler, A. E. (2010). The impact of information technology on academic scientists' productivity and collaboration patterns. *Management Science*, *56*(9), 1439-1461.
- Fiore, S. M. (2008). Interdisciplinarity as teamwork: How the science of teams can inform team science. *Small Group Research*, *39*(3), 251-277.
- Freeman, R. B., & Huang, W. (2015). Collaborating with people like me: Ethnic coauthorship within the United States. *Journal of Labor Economics*, *33*(S1), S289-S318.

- Granovetter, M. S. (1973). The Strength of Weak Ties'. *American Journal of Sociology*, 78(6), 1360-1380.
- Huber, G. P. (1991). Organizational learning: The contributing processes and the literatures. *Organization science*, 2(1), 88-115.
- Jaffe, A. B., Trajtenberg, M., & Henderson, R. (1993). Geographic localization of knowledge spillovers as evidenced by patent citations. *the Quarterly journal of Economics*, 108(3), 577-598.
- Jones, B. F., Wuchty, S., & Uzzi, B. (2008). Multi-university research teams: Shifting impact, geography, and stratification in science. *science*, 322(5905), 1259-1262.
- Katz, J. S., & Martin, B. R. (1997). What is research collaboration?. *Research Policy*, 26(1), 1-18.
- Knoben, J., & Oerlemans, L. A. (2006). Proximity and inter-organizational collaboration: A literature review. *international Journal of management reviews*, 8(2), 71-89.
- Lee, S., & Bozeman, B. (2005). The impact of research collaboration on scientific productivity. *Social studies of science*, 35(5), 673-702.
- Levitt, B., & March, J. G. (1988). Organizational learning. *Annual review of sociology*, 14(1), 319-338.
- Lin, M., Lucas Jr, H. C., & Shmueli, G. (2013). Too big to fail: large samples and the p-value problem. *Information Systems Research*, 24(4), 906-917.
- Mairesse, J., & Pezzoni, M. (2015). Does gender affect scientific productivity?. *Revue économique*, 66(1), 65-113.
- Mairesse, J., & Turner, L. (2005). *Measurement and explanation of the intensity of co-publication in scientific research: An analysis at the laboratory level* (No. w11172). National Bureau of Economic Research.
- McPherson, M., Smith-Lovin, L., & Cook, J. M. (2001). Birds of a feather: Homophily in social networks. *Annual review of sociology*, 27(1), 415-444.
- Misra, S., Stokols, D., & Cheng, L. (2015). The transdisciplinary orientation scale: Factor structure and relation to the integrative quality and scope of scientific publications. *Journal of Translational Medicine and Epidemiology*, 3(2), 1042.
- Mowatt, G., Shirran, L., Grimshaw, J. M., Rennie, D., Flanagan, A., Yank, V., ... & Bero, L. A. (2002). Prevalence of honorary and ghost authorship in Cochrane reviews. *Jama*, 287(21), 2769-2771.
- Nooteboom, B., Van Haverbeke, W., Duysters, G., Gilsing, V., & Van den Oord, A. (2007). Optimal cognitive distance and absorptive capacity. *Research Policy*, 36(7), 1016-1034.
- Pezzoni, M., Mairesse, J., Stephan, P., & Lane, J. (2016). Gender and the publication output of graduate students: A case study. *PLoS One*, 11(1), e0145146.

- Porac, J. F., Wade, J. B., Fischer, H. M., Brown, J., Kanfer, A., & Bowker, G. (2004). Human capital heterogeneity, collaborative relationships, and publication patterns in a multidisciplinary scientific alliance: a comparative case study of two scientific teams. *Research Policy*, *33*(4), 661-678.
- Salomon, G., & Perkins, D. N. (1998). Individual and social aspects of learning. *Review of research in education*, *23*(1), 1-24.
- SNSF, 2011. Regulation on Sinergia Grants. National Research Council. [http://www.snf.ch/SiteCollectionDocuments/sinergia\\_reglement\\_e.pdf](http://www.snf.ch/SiteCollectionDocuments/sinergia_reglement_e.pdf)
- Stokols, D., Hall, K. L., Taylor, B. K., & Moser, R. P. (2008). The science of team science: overview of the field and introduction to the supplement. *American journal of preventive medicine*, *35*(2), S77-S89.
- Uzzi, B., Mukherjee, S., Stringer, M., & Jones, B. (2013). Atypical combinations and scientific impact. *Science*, *342*(6157), 468-472.
- Whitfield, J. (2008). Group theory; What makes a successful team? John Whitfield looks at research that uses massive online databases and network analysis to come up with some rules of thumb for productive collaborations. *Nature*, *455*(7214), 720-724.
- Woolley, A. W., Chabris, C. F., Pentland, A., Hashmi, N., & Malone, T. W. (2010). Evidence for a collective intelligence factor in the performance of human groups. *science*, *330*(6004), 686-688.
- Wuchty, S., Jones, B. F., & Uzzi, B. (2007). The increasing dominance of teams in production of knowledge. *Science*, *316*(5827), 1036-1039.
- Zenger, T. R., & Lawrence, B. S. (1989). Organizational demography: The differential effects of age and tenure distributions on technical communication. *Academy of Management journal*, *32*(2), 353-376.

## References for chapter 4:

- Ahmadpoor, M. and Jones, B. F. (2017). The dual frontier: Patented inventions and prior scientific advance. *Science*, 357(6351):583–587.
- Alger, I., Van Leeuwen, B., and Weibull, J. W. (2019). Estimating social preferences and kantian morality in strategic interactions.
- Alger, I. and Weibull, J. W. (2013). Homo moralis—preference evolution under incomplete information and assortative matching. *Econometrica*, 81(6):2269–2302.
- Alger, I. and Weibull, J. W. (2016). Evolution and Kantian morality. *Games and Economic Behavior*, 98:56–67.
- Andreoli-Versbach, P. and Mueller-Langer, F. (2014). Open access to data: An ideal professed but not practised. *Research Policy*, 43(9):1621–1633.
- Andreoni, J. (1988). Privately provided public goods in a large economy: the limits of altruism. *Journal of public Economics*, 35(1):57–73.
- Andreoni, J. (1990). Impure altruism and donations to public goods: A theory of warm-glow giving. *The Economic Journal*, 100(401):464–477.
- Andreoni, J. (1995). Cooperation in public-goods experiments: kindness or confusion? *The American Economic Review*, pages 891–904.
- Arrow, K. J. (1962). Economic welfare and the allocation of resources for invention. In *The Rate and Direction of Inventive Activity: Economic and Social Factors: A Conference of the Unive*, pages 609–626.
- Auriol, E. and Renault, R. (2008). Status and incentives. *The RAND Journal of Economics*, 39(1):305–326.
- Awad, E., Dsouza, S., Shariff, A., Rahwan, I., and Bonnefon, J.-F. (2020). Universals and variations in moral decisions made in 42 countries by 70,000 participants. *Proceedings of the National Academy of Sciences*, 117(5):2332–2337.
- Ayoubi, C., Pezzoni, M., Visentin, F., et al. (2019). Does it pay to do novel science? the selectivity patterns in science funding. *UNU-MERIT Working Paper Series*, 2019(037).
- Ayoubi, C. and Thurm, B. (2018). Exploring the diversity of preferences: Is a heterogeneous population evolutionarily stable under assortative matching?
- Badia, A. (2014). Data, information, knowledge: An information science analysis. *Journal of the Association for Information Science and Technology*, 65(6):1279–1287.
- Bell, A., Chetty, R., Jaravel, X., Petkova, N., and Van Reenen, J. (2019). Who becomes an inventor in america? the importance of exposure to innovation. *The Quarterly Journal of Economics*, 134(2):647–713.
- Besley, T. and Ghatak, M. (2008). Status incentives. *American Economic Review*, 98(2):206–11.



- Bodle, R. (2011). Regimes of sharing: Open apis, interoperability, and facebook. *Information, Communication & Society*, 14(3):320–337.
- Boudreau, K. J. and Lakhani, K. R. (2015). “open” disclosure of innovations, incentives and follow-on reuse: Theory on processes of cumulative innovation and a field experiment in computational biology. *Research Policy*, 44(1):4–19.
- Bruvoll, A., Halvorsen, B., and Nyborg, K. (2002). Households’ recycling efforts. *Resources, Conservation and recycling*, 36(4):337–354.
- Burton-Chellew, M. N., El Mouden, C., and West, S. A. (2016). Conditional cooperation and confusion in public-goods experiments. *Proceedings of the National Academy of Sciences*, 113(5):1291– 1296.
- Capraro, V. and Rand, D. G. (2018). Do the right thing: Experimental evidence that preferences for moral behavior, rather than equity or efficiency per se, drive human prosociality. *Forthcoming in Judgment and Decision Making*.
- Carpenter, J. and Matthews, P. H. (2017). Using raffles to fund public goods: Lessons from a field experiment. *Journal of Public Economics*, 150:30–38.
- Chen, Y., Farzan, R., Kraut, R., YeckehZaare, I., and Zhang, A. F. (2018). Motivating contributions to public information goods: A field experiment on wikipedia. Technical report, Working Paper.
- Cockburn, I. M., Henderson, R., & Stern, S. (2019). The Impact of Artificial Intelligence on Innovation. *The Economics of Artificial Intelligence: An Agenda*, 115.
- Cowan, R., David, P. A., and Foray, D. (2000). The explicit economics of knowledge codification and tacitness. *Industrial and corporate change*, 9(2):211–253.
- Cowan, R. and Foray, D. (1997). The economics of codification and the diffusion of knowledge. *Industrial and corporate change*, 6(3):595–622.
- Dasgupta, P. and David, P. A. (2002). Toward a new economics of science. *Science Bought and Sold: Essays in the Economics of Science*, page 219.
- David, P. A. (2004). Understanding the emergence of ‘open science’ institutions: functionalist economics in historical context. *Industrial and corporate change*, 13(4):571–589.
- David, P. A. and Foray, D. (1996). Information distribution and the growth of economically valuable knowledge: a rationale for technological infrastructure policies. In *Technological Infrastructure Policy*, pages 87–116. Springer.
- David, P. A., Waterman, A., and Arora, S. (2003). Floss-us the free/libre/open source software survey for 2003. *Stanford Institute for Economic Policy Research*, pages 1–39.
- De la Croix, D., Doepke, M., and Mokyr, J. (2018). Clans, guilds, and markets: Apprenticeship institutions and growth in the preindustrial economy. *The Quarterly Journal of Economics*, 133(1):1–70.

- De Rassenfosse, G., Palangkaraya, A., and Webster, E. (2016). Why do patents facilitate trade in technology? testing the disclosure and appropriation effects. *Research Policy*, 45(7):1326–1336.
- Debackere, K. and Veugelers, R. (2005). The role of academic technology transfer organizations in improving industry science links. *Research Policy*, 34(3):321–342.
- Falk, A., Becker, A., Dohmen, T., Enke, B., Huffman, D., and Sunde, U. (2018). Global evidence on economic preferences. *The Quarterly Journal of Economics*, 133(4):1645–1692.
- Falk, A. and Ichino, A. (2003). Clean evidence on peer pressure.
- Farnham, A., Kurz, C., Öztürk, M. A., Solbiati, M., Myllyntaus, O., Meekes, J., Pham, T. M., Paz, C., Langiewicz, M., Andrews, S., et al. (2017). Early career researchers want open science. *Genome biology*, 18(1):221.
- Fehr, E. and Gächter, S. (1998). Reciprocity and economics: The economic implications of *Homo Reciprocans*. *European Economic Review*, 42(3-5):845–859.
- Fershtman, C. and Gandal, N. (2007). Open source software: Motivation and restrictive licensing. *International Economics and Economic Policy*, 4(2):209–225.
- Fischbacher, U., Gächter, S., and Fehr, E. (2001). Are people conditionally cooperative? evidence from a public goods experiment. *Economics letters*, 71(3):397–404.
- Fleming, L., Greene, H., Li, G., Marx, M., and Yao, D. (2019). Government-funded research increasingly fuels innovation. *Science*, 364(6446):1139–1141.
- Furman, J. L. and Stern, S. (2011). Climbing atop the shoulders of giants: The impact of institutions on cumulative research. *American Economic Review*, 101(5):1933–63.
- Gächter, S., von Krogh, G., and Haefliger, S. (2010). Initiating private-collective innovation: The fragility of knowledge sharing. *Research Policy*, 39(7):893–906.
- Gallus, J. (2017). Fostering public good contributions with symbolic awards: A large-scale natural field experiment at wikipedia. *Management Science*, 63(12):3999–4015.
- Gangopadhyay, K. and Mondal, D. (2012). Does stronger protection of intellectual property stimulate innovation? *Economics Letters*, 116(1):80–82.
- Garry, M. (2009). Big retailers share data: Gma. *Supermarket News*.
- Gewin, V. (2016). Data sharing: An open mind on open data. *Nature*, 529(7584):117–119.
- Ghoshal, A., Kumar, S., and Mookerjee, V. (2018). Dilemma of data sharing alliance: when do competing personalizing and non-personalizing firms share data. *Production and Operations Management*.
- Giles, J. (2005). Internet encyclopaedias go head to head.
- Gneezy, U., Meier, S., and Rey-Biel, P. (2011). When and why incentives (don't) work to modify behavior. *Journal of Economic Perspectives*, 25(4):191–210.

- Gneezy, U. and Rustichini, A. (2000). A fine is a price. *The Journal of Legal Studies*, 29(1):1–17.
- Gould, E. D. and Kaplan, T. R. (2011). Learning unethical practices from a co-worker: the peer effect of Jose Canseco. *Labour Economics*, 18(3):338–348.
- Greenstein, S. and Zhu, F. (2012). Is Wikipedia biased? *American Economic Review*, 102(3):343–48.
- Guedj, D. and Ramjoué, C. (2015). European Commission policy on open access to scientific publications and research data in horizon 2020. *Biomed Data J*, 1(1).
- Gupta, A. K. and Nadarajah, S. (2004). *Handbook of beta distribution and its applications*. CRC press.
- Haeussler, C. (2011). Information-sharing in academia and the industry: A comparative study. *Research Policy*, 40(1):105–122.
- Hergueux, J., Algan, Y., Benkler, Y., and Morell, M. F. (2015). Cooperation in a peer production economy experimental evidence from Wikipedia.
- Hippel, E. v. and Krogh, G. v. (2003). Open-source software and the “private-collective” innovation model: Issues for organization science. *Organization Science*, 14(2):209–223.
- Hopenhayn, H. A. and Squintani, F. (2016). Patent rights and innovation disclosure. *The Review of Economic Studies*, 83(1):199–230.
- Johansson, M. A., Reich, N. G., Meyers, L. A., and Lipsitch, M. (2018). Preprints: An underutilized mechanism to accelerate outbreak science. *PLoS Medicine*, 15(4).
- Kant, I. (1870). *Grundlegung zur metaphysik der sitten*, volume 28. L. Heimann.
- Kim, Y. and Stanton, J. M. (2016). Institutional and individual factors affecting scientists’ data-sharing behaviors: A multilevel analysis. *Journal of the Association for Information Science and Technology*, 67(4):776–799.
- Kollock, P. (1998). Social dilemmas: The anatomy of cooperation. *Annual review of sociology*, 24(1):183–214.
- Kopel, M., Lamantia, F., and Szidarovszky, F. (2014). Evolutionary competition in a mixed market with socially concerned firms. *Journal of Economic Dynamics and Control*, 48:394–409.
- Kurataa, H. and Van Longb, N. (2019). Corporate environmentalism at the kantian equilibrium.
- Lakhani, K. R. and Von Hippel, E. (2004). How open-source software works: “free” user-to-user assistance. In *Produktentwicklung mit virtuellen Communities*, pages 303–339. Springer.
- Lakhani, K. R. and Wolf, R. G. (2003). Why hackers do what they do: Understanding motivation and effort in free/open-source software projects.
- Lee, G. K. and Cole, R. E. (2000). The Linux kernel development as a model of open source knowledge creation. *unpub. MS, Haas School of Business, UC Berkeley*.

- Lee, G. K. and Cole, R. E. (2003). From a firm-based to a community-based model of knowledge creation: The case of the Linux kernel development. *Organization science*, 14(6):633–649.
- Lerner, J. (2009). The empirical impact of intellectual property rights on innovation: Puzzles and clues. *American Economic Review*, 99(2):343–48.
- Lerner, J. and Tirole, J. (2002). Some simple economics of open source. *The journal of industrial economics*, 50(2):197–234.
- Lerner, J. and Tirole, J. (2005). The economics of technology sharing: Open source and beyond. *Journal of Economic Perspectives*, 19(2):99–120.
- Lichbach, M. I. (1996). *The cooperator's dilemma*. University of Michigan Press.
- Liu, J. and Ram, S. (2018). Using big data and network analysis to understand Wikipedia article quality. *Data & Knowledge Engineering*, 115:80–93.
- Longo, D. L. and Drazen, J. M. (2016). Data sharing.
- Luna-Reyes, L. F. and Najafabadi, M. M. (2019). The US open data initiative: The road ahead. *Information Polity*, 24(2):163–182.
- McKiernan, E. C., Bourne, P. E., Brown, C. T., Buck, S., Kenall, A., Lin, J., McDougall, D., Nosek, B. A., Ram, K., Soderberg, C. K., et al. (2016). Point of view: How open science helps researchers succeed. *Elife*, 5:e16800.
- Merton, R. K. (1973). *The sociology of science: Theoretical and empirical investigations*. University of Chicago Press.
- Merton, R. K. (1988). The Matthew effect in science, ii: Cumulative advantage and the symbolism of intellectual property. *isis*, 79(4):606–623.
- Miettinen, T., Kosfeld, M., Fehr, E., and Weibull, J. (2020). Revealed preferences in a sequential prisoners' dilemma: A horse-race between six utility functions. *Journal of Economic Behavior & Organization*, 173:1–25.
- Moser, P. (2013). Patents and innovation: evidence from economic history. *Journal of Economic Perspectives*, 27(1):23–44.
- Munafò, M. R., Nosek, B. A., Bishop, D. V., Button, K. S., Chambers, C. D., Du Sert, N. P., Simonsohn, U., Wagenmakers, E.-J., Ware, J. J., and Ioannidis, J. P. (2017). A manifesto for reproducible science. *Nature human behaviour*, 1(1):1–9.
- Ockenfels, P. (1993). Cooperation in prisoners' dilemma: An evolutionary approach. *European Journal of Political Economy*, 9(4):567–579.
- Olson, M. (1965). *Logic of collective action: Public goods and the theory of groups (Harvard economic studies. v. 124)*. Harvard University Press.

- Owens, B. (2016). Montreal institute going 'open' to accelerate science.
- Poegel, F., Harhoff, D., Gaessler, F., and Baruffaldi, S. (2019). Science quality and the value of inventions. *Science Advances*, 5(12):eaay7323.
- Polanyi, M. (1958). Personal knowledge: Towards a post-critical philosophy Routledge. *Paul, London*.
- Polanyi, M. (1967). The tacit dimension. Anchor. *Garden City, NY*.
- Rigney, D. (2010). *The Matthew effect: How advantage begets further advantage*. Columbia University Press.
- Roemer, J. E. (2015). Kantian optimization: A microfoundation for cooperation. *Journal of Public Economics*, 127:45–57.
- Samuelson, P. A. (1954). The pure theory of public expenditure. *The Review of Economics and Statistics*, pages 387–389.
- Sandel, M. J. (2013). Market reasoning as moral reasoning: why economists should re-engage with political philosophy. *Journal of Economic Perspectives*, 27(4):121–40.
- Singh, J. (2005). Collaborative networks as determinants of knowledge diffusion patterns. *Management Science*, 51(5):756–770.
- Smith, E., Haustein, S., Mongeon, P., Shu, F., Ridde, V., and Larivière, V. (2017). Knowledge sharing in global health research—the impact, uptake and cost of open access to scholarly literature. *Health research policy and systems*, 15(1):73.
- Spichiger, D. (2018). Open science. *Chimia*, 72(5):342–344.
- Squazzoni, F., Bravo, G., and Takacs, K. (2013). Does incentive provision increase the quality of peer review? An experimental study. *Research Policy*, 42(1):287–294.
- Stephan, P. E. (1996). The economics of science. *Journal of Economic Literature*, 34(3):1199–1235.
- Stephan, P. E. (2012). *How economics shapes science*, volume 1. Harvard University Press, Cambridge, MA.
- Stern, S. (2004). Do scientists pay to be scientists? *Management Science*, 50(6):835–853.
- Stiglitz, J. E. (1989). *The economic role of the state*. Wiley-Blackwell.
- Stiglitz, J. E. (1999). Knowledge as a global public good. *Global public goods*, 1(9):308–326.
- Stiglitz, J. E. (2007). Economic foundations of intellectual property rights. *Duke LJ*, 57:1693.
- Stiglitz, J. E. (2014). Intellectual property rights, the pool of knowledge, and innovation. Technical report, National Bureau of Economic Research.
- Sundt, S. and Rehdanz, K. (2015). Consumers' willingness to pay for green electricity: A metaanalysis of the literature. *Energy Economics*, 51:1–8.

- Teplitskiy, M., Lu, G., and Duede, E. (2017). Amplifying the impact of open access: Wikipedia and the diffusion of science. *Journal of the Association for Information Science and Technology*, 68(9):2116–2127.
- Thaler, R. H. and Sunstein, C. R. (2008). *Nudge: Improving Decisions about Health, Wealth, and Happiness*. Yale University Press.
- Thompson, N. and Hanley, D. (2018). Science is shaped by Wikipedia: Evidence from a randomized control trial.
- Thursby, J. G., Haeussler, C., Thursby, M. C., and Jiang, L. (2018). Prepublication disclosure of scientific results: Norms, competition, and commercial orientation. *Science advances*, 4(5):eaar2133.
- Von Krogh, G., Spaeth, S., and Lakhani, K. R. (2003). Community, joining, and specialization in open source software innovation: a case study. *Research Policy*, 32(7):1217–1241.
- Wasko, M. M. and Faraj, S. (2000). “it is what one does”: why people participate and help others in electronic communities of practice. *The journal of strategic information systems*, 9(2-3):155–173.
- Xu, L., Nian, T., and Cabral, L. (2020). What makes geeks tick? A study of stack overflow careers. *Management Science*, 66(2):587–604.
- Zhang, X. M. and Zhu, F. (2011). Group size and incentives to contribute: A natural experiment at Chinese Wikipedia. *American Economic Review*, 101(4):1601–15.
- Zins, C. (2007). Conceptual approaches for defining data, information, and knowledge. *Journal of the American Society for information science and technology*, 58(4):479–493.

## References for chapter 5:

- Acquisti, A., Taylor, C., & Wagman, L. (2016). The economics of privacy. *Journal of economic Literature*, 54(2), 442-92.
- Agrawal, A., McHale, J., & Oettl, A. (2018). *Finding needles in haystacks: Artificial intelligence and recombinant growth* (No. w24541). National Bureau of Economic Research.
- Agrawal, A., Gans, J. S., & Goldfarb, A. (2019). Artificial intelligence: the ambiguous labor market impact of automating prediction. *Journal of Economic Perspectives*, 33(2), 31-50.
- Aizcorbe, A., & Nestoriak, N. (2011). Changing mix of medical care services: stylized facts and implications for price indexes. *Journal of Health Economics*, 30(3), 568-574.
- Akselrod-Ballin, A., Chorev, M., Shoshan, Y., Spiro, A., Hazan, A., Melamed, R., Barkan, E., Herzal, E., Naor, S., Karavani, E. and Koren, G., (2019). Predicting Breast Cancer by Applying Deep Learning to Linked Health Records and Mammograms. *Radiology*, p.182622.
- Alekseeva, L., Azar, J., Gine, M., Samila, S., & Taska, B. (2020). The Demand for AI Skills in the Labor Market.
- Atun, R., 2015. Transitioning health systems for multimorbidity. *The Lancet*, 386(9995), pp.721-722.
- Autor, D., Dorn, D., Katz, L., & Patterson, C. J. Van Reenen (2017). "The Fall of the Labor Share and the Rise of Superstar Firms." *NBER Working Paper*, (23396).
- Autor, D., & Salomons, A. (2018). *Is Automation Labor-Displacing? Productivity Growth, Employment, and the Labor Share* (No. 24871). National Bureau of Economic Research, Inc.
- Badawi, O., Brennan, T., Celi, L. A., Feng, M., Ghassemi, M., Ippolito, A., ... & Moses, C. (2014). Making big data useful for health care: a summary of the inaugural mit critical data conference. *JMIR medical informatics*, 2(2), e22.
- Baumol, W. J. (1967). Macroeconomics of unbalanced growth: the anatomy of urban crisis. *The American economic review*, 57(3), 415-426.
- Baumol, W. J. (1993). Health care, education and the cost disease: A looming crisis for public choice. *Public choice*, 77(1), 17-28.
- Baumol, W. J. (2012). *The cost disease: Why computers get cheaper and health care doesn't*. Yale university press.
- Bessen, J., & Hunt, R. M. (2007). An empirical look at software patents. *Journal of Economics & Management Strategy*, 16(1), 157-189.
- Bessen, J. E., Impink, S. M., Reichensperger, L., & Seamans, R. (2018). The Business of AI Startups. *Boston Univ. School of Law, Law and Economics Research Paper*, (18-28).

- Bibault, J. E., Giraud, P., Durdux, C., Taieb, J., Berger, A., Coriat, R.,... & Burgun, A. (2018). Deep Learning and Radiomics predict complete response after neo-adjuvant chemoradiation for locally advanced rectal cancer. *Scientific reports*, *8*(1), 12611.
- Bojke, C., Castelli, A., Grašič, K., & Street, A. (2017). Productivity growth in the English National Health Service from 1998/1999 to 2013/2014. *Health economics*, *26*(5), 547-565.
- Bresnahan, T. F., & Trajtenberg, M. (1995). General purpose technologies 'Engines of growth?'. *Journal of econometrics*, *65*(1), 83-108.
- Bresnahan, T. F., & Greenstein, S. M. (1997). *Technical progress and co-invention in computing and in the uses of computers* (Vol. 477). Institute of Government and Public Affairs, University of Illinois.
- Bresnahan, T. (2010). General purpose technologies. In *Handbook of the Economics of Innovation* (Vol. 2, pp. 761-791). North-Holland.
- Bresnahan, T. (2019). Artificial Intelligence technologies and Aggregate Growth Prospects. Available at: [https://web.stanford.edu/~tbres/AI\\_Technologies\\_in\\_use.pdf](https://web.stanford.edu/~tbres/AI_Technologies_in_use.pdf)
- Brynjolfsson, E., & Mitchell, T. (2017). What can machine learning do? Workforce implications. *Science*, *358*(6370), 1530-1534.
- Brynjolfsson, E., Mitchell, T., & Rock, D. (2018). What Can Machines Learn, and What Does It Mean for Occupations and the Economy?. In *AEA Papers and Proceedings* (Vol. 108, pp. 43-47).
- Brynjolfsson, E., Rock, D., & Syverson, C. (2019). Artificial Intelligence and the Modern Productivity Paradox. *The Economics of Artificial Intelligence: An Agenda*, 23.
- Camerani, R., Rotolo, D., & Grassano, N. (2018). Do firms publish? A multi-sectoral analysis. *A Multi-Sectoral Analysis (October 2018)*. SWPS, 21.
- Chandra, A., & Skinner, J. (2012). Technology growth and expenditure growth in health care. *Journal of Economic Literature*, *50*(3), 645-80.
- Chesbrough, H., & Rosenbloom, R. S. (2002). The role of the business model in capturing value from innovation: evidence from Xerox Corporation's technology spin-off companies. *Industrial and corporate change*, *11*(3), 529-555.
- Christensen, C. M. (1997). The Innovator's Dilemma: When New Technologies Cause Great Firms to Fail.
- Cockburn, I. M., Henderson, R., & Stern, S. (2019). The Impact of Artificial Intelligence on Innovation. *The Economics of Artificial Intelligence: An Agenda*, 115.
- Cutler, D. M. (2011). Where are the health care entrepreneurs? The failure of organizational innovation in health care. *Innovation Policy and the Economy*, *11*(1), 1-28.
- Cutler, D. M. (2017). Rising medical costs mean more rough times ahead. *Jama*, *318*(6), 508-509.



- Cutler, D. M., & Richardson, E. (1998). The value of health: 1970-1990. *The American economic review*, 88(2), 97-100.
- Day, S., & Zweig, M. (2018). Beyond Wellness For the Healthy: Digital Health Consumer Adoption, 2018.
- De Fauw, J., Ledsam, J. R., Romera-Paredes, B., Nikolov, S., Tomasev, N., Blackwell, S.,... & van den Driessche, G. (2018). Clinically applicable deep learning for diagnosis and referral in retinal disease. *Nature medicine*, 24(9), 1342.
- DeMonaco, H. J., Ali, A., & Von Hippel, E. (2006). The major role of clinicians in the discovery of off-label drug therapies. *Pharmacotherapy: The Journal of Human Pharmacology and Drug Therapy*, 26(3), 323-332.
- Feldstein, M. (2017). Underestimating the real growth of GDP, personal income, and productivity. *Journal of Economic Perspectives*, 31(2), 145-64.
- Fukuzawa, N., & Ida, T. (2016). Science linkages between scientific articles and patents for leading scientists in the life and medical sciences field: The case of Japan. *Scientometrics*, 106(2), 629-644.
- Gardner, R. L., Cooper, E., Haskell, J., Harris, D. A., Poplau, S., Kroth, P. J., & Linzer, M. (2019). Physician stress and burnout: the impact of health information technology. *Journal of the American Medical Informatics Association*, 26(2), 106-114.
- Gofman, M., & Jin, Z. (2019). Artificial Intelligence, Human Capital, and Innovation. *Human Capital, and Innovation (August 20, 2019)*.
- Goldfarb, B. (2005). Diffusion of general-purpose technologies: understanding patterns in the electrification of US Manufacturing 1880-1930. *Industrial and Corporate Change*, 14(5), 745-773.
- Goolsbee, A., & Klenow, P. J. (2002). Evidence on learning and network externalities in the diffusion of home computers. *The Journal of Law and Economics*, 45(2), 317-343.
- Guntersdorfer, M. (2003). Software patent law: United States and Europe compared. *Duke Law & Technology Review*, 2(1), 1-12.
- Gutiérrez, G., & Philippon, T. (2017). *Declining Competition and Investment in the US* (No. w23583). National Bureau of Economic Research.
- Gutiérrez, G., & Philippon, T. (2018). *How EU markets became more competitive than US markets: A study of institutional drift* (No. w24700). National Bureau of Economic Research.
- Haenssle, H. A., Fink, C., Schneiderbauer, R., Toberer, F., Buhl, T., Blum, A.,... & Uhlmann, L. (2018). Man against machine: diagnostic performance of a deep learning convolutional neural network for dermoscopic melanoma recognition in comparison to 58 dermatologists. *Annals of Oncology*, 29(8), 1836-1842.
- Helpman, E., & Trajtenberg, M. (1996). *Diffusion of general purpose technologies* (No. w5773). National bureau of economic research.

- Hendel, I., & Spiegel, Y. (2014). Small steps for workers, a giant leap for productivity. *American Economic Journal: Applied Economics*, 6(1), 73-90.
- Hendrich, A., Chow, M. P., Skierczynski, B. A., & Lu, Z. (2008). A 36-hospital time and motion study: how do medical-surgical nurses spend their time?. *The Permanente Journal*, 12(3), 25.
- Hildebrandt, M. (2018). Law as computation in the era of artificial legal intelligence: Speaking law to the power of statistics. *University of Toronto Law Journal*, 68(supplement 1), 12-35.
- Iansiti, M., & Lakhani, K. R. (2020) *Competing in the Age of AI: Strategy and Leadership When Algorithms and Networks Run the World*. Boston: Harvard Business Review Press, in press.
- James, S. D., Leiblein, M. J., & Lu, S. (2013). How firms capture value from their innovations. *Journal of management*, 39(5), 1123-1155.
- Kocher, R., & Sahni, N. R. (2011). Hospitals' race to employ physicians—the logic behind a money-losing proposition. *N Engl J Med*, 364(19), 1790-1793.
- Kurakin, A., Goodfellow, I., & Bengio, S. (2016). Adversarial examples in the physical world. *arXiv preprint arXiv:1607.02533*.
- Lee, J., McCullough, J. S., & Town, R. J. (2013). The impact of health information technology on hospital productivity. *The RAND Journal of Economics*, 44(3), 545-568.
- Mihet, R., & Philippon, T. (2019). The Economics of Big Data and Artificial Intelligence. *Disruptive Innovation in Business and Finance in the Digital World (International Finance Review, Vol. 20)*, Emerald Publishing Limited, 29-43.
- Miric, M., Boudreau, K. J., & Jeppesen, L. B. (2019). Protecting their digital assets: The use of formal & informal appropriability strategies by App developers. *Research Policy*, 48(8), 103738.
- Mukherjee, S., Romero, D. M., Jones, B., & Uzzi, B. (2017). The nearly universal link between the age of past knowledge and tomorrow's breakthroughs in science and technology: The hotspot. *Science advances*, 3(4), e1601315.
- Oliveira, P., Zejnilovic, L., Canhão, H., & von Hippel, E. (2015). Innovation by patients with rare diseases and chronic needs. *Orphanet Journal of Rare Diseases*, 10(1), 41.
- Rajkomar, A., Dean, J., & Kohane, I. (2019). Machine learning in medicine. *New England Journal of Medicine*, 380(14), 1347-1358.
- Rao, S. K., Kimball, A. B., Lehrhoff, S. R., Hidrue, M. K., Colton, D. G., Ferris, T. G., & Torchiana, D. F. (2017). The impact of administrative burden on academic physicians: results of a hospital-wide physician survey. *Academic Medicine*, 92(2), 237-243.
- Roberts, H., Cows, J., Morley, J., Taddeo, M., Wang, V., & Floridi, L. (2019). The Chinese Approach to Artificial Intelligence: an Analysis of Policy and Regulation. *Available at SSRN 3469784*.

- Rock, D. (2019). *Engineering Value: The Returns to Technological Talent and Investments in Artificial Intelligence*.
- Sabater-Mir, J., Torra, V., & Aguiló, I. (Eds.). (2019). *Artificial Intelligence Research and Development: Proceedings of the 22nd International Conference of the Catalan Association for Artificial Intelligence* (Vol. 319). IOS Press.
- Sahni, N. R., Huckman, R. S., Chigurupati, A., & Cutler, D. M. (2017). The IT Transformation Health Care Needs. *HARVARD BUSINESS REVIEW*, 95(6), 129-136.
- Scott, R. D. (2009). The direct medical costs of healthcare-associated infections in US hospitals and the benefits of prevention.
- Shapiro, C. (2019). Protecting competition in the American economy: Merger control, tech titans, labor markets. *Journal of Economic Perspectives*, 33(3), 69-93.
- Sudlow, C., Gallacher, J., Allen, N., Beral, V., Burton, P., Danesh, J., ... & Liu, B. (2015). UK biobank: an open access resource for identifying the causes of a wide range of complex diseases of middle and old age. *PLoS medicine*, 12(3).
- Sun, R., Limkin, E.J., Vakalopoulou, M., Dercle, L., Champiat, S., Han, S.R., Verlingue, L., Brandao, D., Lancia, A., Ammari, S. and Hollebecque, A., (2018). A radiomics approach to assess tumour-infiltrating CD8 cells and response to anti-PD-1 or anti-PD-L1 immunotherapy: an imaging biomarker, retrospective multi-cohort study. *The Lancet Oncology*, 19(9), pp.1180-1191.
- Taddy, M. (2018). *The technological elements of artificial intelligence* (No. w24301). National Bureau of Economic Research.
- Tomczak, K., Czerwińska, P., & Wiznerowicz, M. (2015). The Cancer Genome Atlas (TCGA): an immeasurable source of knowledge. *Contemporary oncology*, 19(1A), A68.
- Trajtenberg, 2018, *AI as the next GPT: a political economy perspective*, NBER wp 24245.
- Varian, H. (2018). *Artificial intelligence, economics, and industrial organization* (No. w24839). National Bureau of Economic Research.
- Webb, M. (2019). The Impact of Artificial Intelligence on the Labor Market. *Available at SSRN 3482150*.
- World Intellectual Property Organization. (2019). *WIPO Technology Trends 2019: Artificial Intelligence*. World Intellectual Property Organization.
- Yu, K. H., Beam, A. L., & Kohane, I. S. (2018). Artificial intelligence in healthcare. *Nature biomedical engineering*, 2(10), 719-731.

## References for chapter 6:

- Ayoubi, C., Pezzoni, M., & Visentin, F. (2019). The important thing is not to win, it is to take part: What if scientists benefit from participating in research grant competitions?. *Research Policy*, *48*(1), 84-97. a
- Ayoubi, C., Pezzoni, M., & Visentin, F. (2019). Does it pay to do novel science? The selectivity patterns in science funding. *UNU-MERIT Working Paper Series*, *2019* (037). b
- Ayoubi, C., Barbosu, S., Pezzoni, M., & Visentin, F. (2020). What matters in funding: The value of research coherence and alignment in evaluators' decisions. *UNU-MERIT Working Papers*, *2020* (010).
- Ayoubi, C. and Thurm, B. (2020). Environmental care and morality: An economic model with heterogeneous agents. *25th Annual Conference of the European Association of Environmental and Resource Economists (EAERE)*.
- Azoulay, P., Stuart, T., & Wang, Y. (2014). Matthew: Effect or fable?. *Management Science*, *60*(1), 92-109.
- Bikard, M., & Marx, M. (2019). Bridging Academia and Industry: How Geographic Hubs Connect University Science and Corporate Technology. *Management Science*.
- Fleming, L., Greene, H., Li, G., Marx, M., and Yao, D. (2019). Government-funded research increasingly fuels innovation. *Science*, *364*(6446):1139–1141.
- Foray, D., Mowery, D. C., & Nelson, R. R. (2012). Public R&D; and social challenges: What lessons from mission R&D; programs?. *Research Policy*, *41*(ARTICLE), 1697-1702.
- Perc, M. (2014). The Matthew effect in empirical data. *Journal of The Royal Society Interface*, *11*(98), 20140378.
- Schelling, T. C. (1996). The economic diplomacy of geoengineering. *Climatic Change*, *33*(3), 303-307.
- Stephan, P., Veugelers, R., & Wang, J. (2017). Reviewers are blinkered by bibliometrics. *Nature*, *544*(7651), 411-412.



# Appendix

## Appendix A2 for chapter 2:

### Appendix A2.1 – Discipline definition

We assign each scientist to one or more disciplines according to the journal, where she published during the five years preceding the application year. A researcher is assigned to a discipline if she published at least one article in one of the journals classified in that discipline. We retrieve the classification of journals using the journal classification of the Scopus Source List<sup>86</sup>. When excluding Humanities and Social Sciences, the journal classification table reports 20 disciplines. We exclude generalist journals and journals belonging to more than one discipline to avoid ambiguous classification of researchers. Table A2.1 reports the non-exclusive attribution of disciplines to the sample of applicants.

**Table A2.1: List of disciplines of applicants (1,060 observations).**

SCOPUS disciplines	Mean	Std. Dev.	Min	Max
Chemistry	0.14	0.34	0	1
Engineering	0.09	0.29	0	1
Materials Science	0.05	0.22	0	1
Physics and Astronomy	0.20	0.40	0	1
Agricultural and Biological Sciences	0.10	0.30	0	1
Biochemistry, Genetics and Molecular Biology	0.42	0.49	0	1
Environmental Science	0.07	0.25	0	1
Computer Science	0.08	0.28	0	1
Mathematics	0.02	0.15	0	1
Medicine	0.32	0.47	0	1
Health Professions	0.00	0.03	0	1
Nursing	0.00	0.04	0	1
Earth and Planetary Sciences	0.08	0.27	0	1
Energy	0.01	0.09	0	1
Pharmacology, Toxicology and Pharmaceutics	0.06	0.25	0	1
Chemical Engineering	0.02	0.13	0	1
Neuroscience	0.17	0.38	0	1
Veterinary	0.01	0.11	0	1
Immunology and Microbiology	0.11	0.32	0	1
Dentistry	0.01	0.08	0	1

The table reports mean, standard deviation, minimum and maximum for the 20 dummy variables used to identify the scientists' discipline(s). The attribution of a discipline to a scientist relies on the discipline classification of the journals where she publishes. The classification of the journals is obtained from the bibliometric database Scopus Source List (Elsevier). The attribution of a scientist to a discipline is not exclusive; a scientist might be classified in more than one discipline.

<sup>86</sup> <https://www.scopus.com/sources>

## Appendix A2.2 – Dealing with endogeneity issues in assessing the effect of applying for a grant

In this section, we consider the possible endogeneity concerns that could be raised in estimating the effects of applying for a grant. We thus propose an alternative estimation strategy using a two-stage least square (2SLS) approach. We find that the core results presented in the main text are robust to the choice of estimation strategy.

Better researchers might self-select themselves into applying since they believe they have higher chances to succeed in being funded. Having high scientific profiles, applicants might perform better than other scientists regardless of the involvement in the application process. To address this selection bias concern, we follow an instrumental variable approach in addition to the propensity score matching strategy described in paragraph 2.3.2 of the main text. Specifically, we identify two excluded instruments for the *Applicant* variable. As first instrument, we consider the attractiveness of the researcher's discipline for funding agencies (*Funds in the US*). As second instrument, we consider the presence in the researcher co-authorship network of a scientist who applied for SINERGIA in the previous years (*Network Applicant*).

The first excluded instrument, *Funds in the US*, is a continuous variable that measures the availability of funds in the scientist's discipline during the application year<sup>87,88</sup>. The researcher who belongs to a discipline characterized by an extensive availability of funds is expected to perceive higher chances of success for her application, and she is then incentivized to apply. Given this correlation between *Funds in the US* and the dummy *Applicant*, the instrument is strong. It is also a valid instrument since the perception of higher chances to be awarded is not expected to influence the researcher's outcomes. The funds considered are those available in the researcher's discipline in the US (and not in Switzerland) to avoid any possible correlation with the unobserved factors that determine the researchers' outcomes in our sample. Given that scientific trends are not constrained by a country's boundaries, we expect disciplines highly funded in the US to be also highly funded in Switzerland. We attribute a value of funding to each researcher-application year according to the discipline(s) where the research is active (see Appendix A2.1 for discipline attributions)<sup>89</sup>.

---

<sup>87</sup> Data on yearly US funds by discipline are available at <https://www.aaas.org/page/historical-rd-data>

<sup>88</sup> The available funds are adjusted for the size of each discipline as represented by the number journals classified as belonging to the discipline in the Scopus Source List 2017 journal classification table. (See Appendix A2.1)

<sup>89</sup> If the researcher is classified in more than one discipline, we consider the sum of the funds available for each discipline (i.e. the propensity to apply of a researcher conducting her research in medicine and engineering will be affected by the funds made available in both disciplines).

The second excluded instrument, *Network Applicant*, is a dummy variable that equals one if at least one of the co-authors of the researcher or, one of the co-authors' co-authors, has applied for a SINERGIA grant in the five years preceding the application year<sup>90</sup>. We expect *Network Applicant* to be negatively correlated with the probability of applying for SINERGIA. The reason for this expected negative correlation is that the SINERGIA call requires an applicant to find partners for her application. An applicant usually looks for co-applicants in her professional network. If the researchers in her network have already applied, they are less likely to be part of another application regardless of the result of their application. If they were awarded, there is a weak incentive for them to apply again. If they were not awarded, they might be discouraged from applying again. In both cases, the pool of potential co-applicants of the focal researcher is reduced with a detrimental effect on her probability to apply because she has greater difficulties in finding valuable partners for the application. Hence, we expect *Network Applicant* to be a strong instrument. Concerning the validity of the instrument, as suggested by the homophily mechanism of Cummings and Kiesler (2008)<sup>91</sup>, the choice of collaboration partners is mainly based on sociological characteristics or geographical considerations. The excluded instrument *Network Applicant* is, therefore, not directly correlated with the unobserved ability of researchers driving her expected scientific outcomes. Table A2.2 reports the results of our estimation exercise.

We test for the presence of endogeneity in the estimations using a *Durbin-Wu-Hausman* test. The test does not reject the null hypothesis that the variables are exogenous. Given the results of the *Durbin-Wu-Hausman*, OLS appears as the most efficient estimation, and its estimations are consistent with our main results (Table A2.3).

---

<sup>90</sup> For a potential applicant, we consider as application year the year when her matched applicant submitted her application.

<sup>91</sup> Cummings, Jonathon N., Kiesler, Sara, "Who collaborates successfully? prior experience reduces collaboration barriers in distributed interdisciplinary research" *Proceedings of the 2008 ACM Conference on Computer Supported Cooperative Work ACM*, pp. 437–446 (2008).



**Table A2.2: Regression results for the scientific outcomes comparing applicants to potential applicants. IV estimations (2SLS).**

	(1) IV	(2) IV	(3) IV	(4) IV	(5) IV	(6) First step Applicant
	Log(Publication count)	Log(Average IF)	Log(Average citations)	Co-applicant collaboration	Log(Learn- ing)	
Applicant	0.69*** (0.47 ; 0.91)	0.12 (-0.037 ; 0.29)	-0.31*** (-0.52 ; -0.11)	0.46*** (0.32 ; 0.60)	0.66*** (0.39 ; 0.93)	
Funds in the US						0.082*** (0.025 ; 0.14)
Network applicants						-0.35*** (-0.40 ; -0.30)
Controls	Yes	Yes	Yes	Yes	Yes	Yes
Observations	2,120	2,120	2,120	2,120	2,120	2,120
R2	0.642	0.496	0.498	0.584	0.486	0.403
Durbin-Wu-Hausman	0.16	0.60	0.37	0.84	0.24	

The table reports an instrumental variable estimation of the effect of applying for a SINERGIA grant (*Applicant*). The controls include the scientist's demographic and bibliometric characteristics. The five scientific outcomes considered are *Publication count*, *Average IF*, *Average citations*, *Co-applicant collaboration*, and *Learning*. Column 6 reports the first stage equation of the 2SLS models. The dependent variable in column 6 is the dummy variable *Applicant*. The two excluded instruments are the variables *Funds in the US* and *Network applicants*. The sample includes 2,120 observations, i.e. 1,060 Applicant-application pairs, and 1,060 Potential applicants. Robust standard errors are clustered around the application. In reporting the statistical significance of the coefficients, we apply the standard thresholds, namely \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ . 95% confidence intervals are reported in parenthesis below each coefficient.

**Table A2.3: Regression results for the scientific outcomes comparing applicants to potential applicants. OLS estimations.**

	(1) OLS	(2) OLS	(3) OLS	(4) OLS	(5) OLS
	Log(Publication count)	Log(Average IF)	Log(Average citations)	Co-applicant col- laboration	Log(Learning)
Applicant	0.53*** (0.46 ; 0.60)	0.081*** (0.021 ; 0.14)	-0.41*** (-0.47 ; -0.34)	0.48*** (0.42 ; 0.53)	0.50*** (0.41 ; 0.59)
Controls	Yes	Yes	Yes	Yes	Yes
Observations	2,120	2,120	2,120	2,120	2,120
R2	0.645	0.496	0.499	0.584	0.489

The table reports an OLS estimation of the effect of applying to the SINERGIA grant (*Applicant*). The controls include the scientist's demographic and bibliometric characteristics. The five scientific outcomes considered are *Publication count*, *Average IF*, *Average citations*, *Co-applicant collaboration*, and *Learning*. The sample includes 2,120 observations, 1,060 Applicant-application pairs, and 1,060 Potential applicants. Robust standard errors are clustered around the application. In reporting the statistical significance of the coefficients, we apply the standard thresholds, namely \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ . 95% confidence intervals are reported in parenthesis below each coefficient.

### Appendix A2.3 – Entering new research domains

As discussed in Section 2.6 of the main text, we measure the entry in new research domains relying on the knowledge mobilized by a researcher in her scientific publications. Precisely, we measure the cognitive distance between the researcher’s learning -i.e., the new journals cited after the application date by a scientist- and her initial stock of knowledge -i.e., the journals cited in the articles published before the application year. To do so, we build on Ayoubi et al. (2017) to construct a journal distance matrix measuring the distance between all existing pairs of journals. This section reports the details of the construction of the matrix and the steps leading to the creation of the variable *Journal distance* used as dependent variable in Table 2.8 in the main text.

#### *The journal distance matrix*

We build the journal distance matrix based on the assumption that the closer two journals are in terms of scientific content, the more likely they are to be co-cited in scientific publications. Exploiting this assumption, we calculate the closeness of two journals as the ratio of the number of publications in which they are co-cited to the number of publications in which the least cited journal of the two is cited. Then, the distance reported in our journal distance matrix is the inverse of the computed closeness. Specifically, we derive the distance separating two journals  $i$  and  $j$  as the inverted ratio between the number of publications in which  $i$  and  $j$  are cited together and the minimum number of publications where  $i$  or  $j$  are cited depending on which one is the least cited. This computed distance measure ranges from 1 to infinity. For the infinite values, namely when  $i$  and  $j$  are never cited together in a publication, for computational purposes, we consider as distance the maximum non-infinite distance found in the matrix. Equation A2.1 shows the calculation of the journal distance for two journals  $i$  and  $j$ .

$$D(i, j) = \frac{1}{\frac{\text{\#pubs where } i \text{ and } j \text{ are co-cited}}{\min(\text{\#pubs where } i \text{ is cited}, \text{\#pubs where } j \text{ is cited})}}$$

(Equation A2.1)

We calculate the journal distance matrix based on the publications of Swiss scientists active in the top 12 universities in Switzerland (excluding humanities and social sciences) in the period 2003-2015, for a total of more than 200,000 articles. The journal distance matrix includes 3,013 journals. The minimum distance value is 1, while the maximum distance is 5,628. In the computed matrix, 65% of the journal pairs never appear together in the reference list of an article. For those pairs, we fix the distance value at the maximum. For the sake of illustration, we consider an extraction from the journal distance matrix, including two journals in physics, i.e., *Physical Review Letters* and *Annals of Physics*, and a

journal in molecular biology, i.e., *The EMBO journal*. Table A2.4 shows the distance between the three journals in the journal distance matrix extraction.

**Table A2.4: Journal distance matrix extraction.**

	<i>Annals of Physics</i>	<i>The EMBO journal</i>	<i>Physical Review Letters</i>
<i>Annals of Physics</i>	-	119	2
<i>The EMBO journal</i>	119	-	66
<i>Physical Review Letters</i>	2	66	-

The table shows an extraction of the journal distance matrix resulting from the calculation of the distance between each pair of journals. The distance calculated between two journals is based on the assumption that the closer two journals are in terms of scientific content, the more likely they are to be co-cited in scientific publications.

As expected, the two journals in physics are close, meaning that they are frequently co-cited in the references of published articles ( $D(\textit{Physical Review Letters}, \textit{Annals of Physics})=2$ ). These two journals in physics are far from the journal in biology (*The EMBO journal*) with a distance value of respectively  $D(\textit{The EMBO journal}, \textit{Annals of Physics})=119$  and  $D(\textit{The EMBO journal}, \textit{Physical Review Letters})=66$ .

We use the journal distance matrix to calculate the average cognitive distance between the scientist's stock of knowledge before the application year and the knowledge she acquires after the application year. To calculate this cognitive distance, we proceed in four steps. First, we extract both the list of distinct journals cited by the scientist before the application and the list of the new journals cited after the application. Second, we construct the pairs resulting from all the possible combinations of distinct journals in the two lists of the first step. Third, we attribute to each journal pair its corresponding distance value in the journal distance matrix. Finally, we create the variable *Journal distance* as the average of the distances retrieved in step three. The variable *Journal distance*, in logarithmic form, is the dependent variable of the regression shown in Table 2.8 in the main text.

## **Appendix A2.4 Robustness checks**

### **Study sample restricted to the applicant-potential applicant pairs characterized by a highly similar propensity score**

The sample of potential applicants used in the main text has been constructed applying the propensity score matching technique. However, some matched applicant-potential applicant pairs, show a higher similarity in their propensity scores than others. We construct an alternative study sample restricted only to the applicant-potential applicant pairs, where the potential applicants have a propensity score highly similar to the one of applicants. Precisely, we restrict our analysis to the applicant-potential applicant pairs for which the propensity score difference is below 2.09%. We choose as threshold the median value of the differences of propensity scores of the sample used in the main analysis (2.09%). The threshold choice restricts the sample to 2,120 observations, i.e., half of the total number of the observations of the main analysis. Using this alternative study sample, we find results comparable with the ones reported in Table 2.6 in the main text (See Table A2.5 below).

**Table A2.5: Regression results for the estimation of Equation 2.1 comparing applicants to potential applicants (Study sample constructed with the constraint of having a propensity score difference smaller than the median difference (2.09%)).**

VARIABLES	(1)	(2)	(3)	(4)	(5)
	OLS Log(Publication count)	OLS Log(Average IF)	OLS Log(Average citations)	Probit Co-applicant Collaboration	OLS Log(Learning)
Applicant*Post-Application	0.47*** (0.38 ; 0.56)	0.12*** (0.030 ; 0.20)	-0.29*** (-0.38 ; -0.20)	0.28*** (0.14 ; 0.41)	0.44*** (0.33 ; 0.55)
Applicant	0.24*** (0.15 ; 0.32)	-0.034 (-0.12 ; 0.053)	-0.11* (-0.22 ; 0.0042)	0.36*** (0.29 ; 0.43)	0.013 (-0.064 ; 0.090)
Post-Application	-0.24*** (-0.31 ; -0.16)	-0.11*** (-0.19 ; -0.036)	-0.21*** (-0.29 ; -0.12)	-0.068 (-0.16 ; 0.025)	-0.60*** (-0.68 ; -0.53)
Seniority	0.030*** (0.025 ; 0.036)	-0.0013 (-0.0056 ; 0.0030)	-0.0015 (-0.0068 ; 0.0037)	-0.0015* (-0.0032 ; 0.00023)	0.012*** (0.0073 ; 0.016)
Other active funding	0.29*** (0.19 ; 0.39)	0.20*** (0.13 ; 0.28)	0.16*** (0.062 ; 0.27)	-0.017 (-0.051 ; 0.018)	0.23*** (0.15 ; 0.32)
Previous expired funding	0.25*** (0.15 ; 0.36)	-0.040 (-0.13 ; 0.045)	-0.025 (-0.14 ; 0.091)	0.012 (-0.023 ; 0.047)	0.092** (0.0022 ; 0.18)
Constant	1.63*** (1.39 ; 1.87)	1.01*** (0.82 ; 1.21)	0.87*** (0.66 ; 1.09)		3.17*** (2.95 ; 3.40)
Dummy Application year	Yes	Yes	Yes	Yes	Yes
Dummy Discipline	Yes	Yes	Yes	Yes	Yes
Dummy Affiliation	Yes	Yes	Yes	Yes	Yes
R2 / Pseudo R2	0.418	0.256	0.237	0.40	0.435
N. of Researchers	1,060	1,060	1,060	1,060	1,060
Observations	2,120	2,120	2,120	2,120	2,120

The table shows a difference-in-differences estimation in the equivalent regression formulation, where the effect of the treatment, i.e., applying for a SINERGIA grant, can be read in the coefficient of the interaction variable *Applicant\*Post-Application*. The five scientific outcomes considered are *Publication count*, *Average IF*, *Average citations*, *Co-applicant collaboration*, and *Learning*. The sample includes 2,120 observations. The 530 Applicant-application pairs and the 530 Potential applicants (keeping only the pairs having a propensity score difference smaller than the median difference (2.09%)), i.e., 1,060 Applicant-application pairs/Potential applicants-matched application pairs, are observed in two periods, namely before and after the application year. Columns 1, 2, 3, and 5 report OLS estimates, whereas Column 4 reports the marginal effect of a Probit estimation that considers the binary nature of the dependent variable *Co-applicant collaboration*. Robust standard errors are clustered around the application. In reporting the statistical significance of the coefficients, we apply the standard thresholds, namely \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ . 95% confidence intervals are reported in parenthesis below each coefficient.

**Control sample where applicants and potential applicants have no statistical difference on the variable *N. of journals pre-application***

The sample of potential applicants used in the analysis reported in the main text shows a statistically significant difference in the average value of the number of journals referenced before applying to SINERGIA (*N. of journals pre-application*) when we compare applicants and potential applicants. We constructed an alternative control sample where, in searching for potential applicants, we impose the constraint of exhibiting no significant statistical difference on the variable *N. of journals pre-application*, ignoring the other characteristics of the potential applicants. Table A2.6 reports the results of the regression exercise that replicates the analysis of the main text (Section 2.4) with the use of this alternative control sample. The results are consistent with our main findings with the only exception of the non-significance of the estimated coefficient of the variable *Applicant\*Post-Application* in the regression having *Average IF* as dependent variable (column 2 in Table A2.6). However, the *Applicant\*Post-Application* coefficient is likely to be positive since most of its range lies in the strictly positive part of the estimated confidence interval, and the coefficient of interest (0.037), although rather small, is positive and similar to the one estimated in Table 2.6 in the main text (0.070).

Taking apart the *N. of journals pre-application*, this alternative control sample shows a low-quality matching. Applicants and potential applicants differ significantly according to several important characteristics (not included in the matching criteria): *Seniority*, *Publication count pre-application*, *Average citations pre-application*, *Co-applicant collaboration pre-application*, *Average authors pre-application*, *Other active funding*, *Previous expired funding*, and *Productivity break*. This low matching quality refrained us to apply this alternative control sample in the analysis reported in the main text.

**Table A2.6: Regression results for the estimation of Equation 2.1 comparing applicants to potential applicants (Control sample where applicants and potential applicants have no statistical difference on the variable *N. of journals pre-application*).**

	(1) OLS Log(Publication count)	(2) OLS Log(Average IF)	(3) OLS Log(Average citations)	(4) Probit Co-applicant Collaboration	(5) OLS Log(Learning)
Applicant*Post-Application	0.37*** (0.30 ; 0.43)	0.037 (-0.017 ; 0.091)	-0.18*** (-0.26 ; -0.11)	0.18*** (0.11 ; 0.24)	0.45*** (0.38 ; 0.53)
Applicant	0.41*** (0.34 ; 0.48)	-0.036 (-0.10 ; 0.032)	-0.32*** (-0.40 ; -0.24)	0.43*** (0.35 ; 0.50)	-0.062** (-0.11 ; -0.014)
Post-Application	-0.19*** (-0.25 ; -0.13)	-0.034 (-0.084 ; 0.016)	-0.30*** (-0.37 ; -0.24)	-0.0025 (-0.052 ; 0.047)	-0.72*** (-0.78 ; -0.65)
Seniority	0.034*** (0.031 ; 0.038)	-0.0036*** (-0.0062 ; -0.0011)	-0.0074*** (-0.011 ; -0.0038)	-0.0012 (-0.0028 ; 0.00049)	0.011*** (0.0087 ; 0.014)
Other active funding	0.28*** (0.22 ; 0.35)	0.16*** (0.11 ; 0.21)	0.14*** (0.073 ; 0.21)	-0.012 (-0.041 ; 0.017)	0.19*** (0.12 ; 0.25)
Previous expired funding	0.22*** (0.15 ; 0.30)	0.0049 (-0.044 ; 0.054)	-0.044 (-0.11 ; 0.021)	0.041** (0.0097 ; 0.073)	0.094*** (0.031 ; 0.16)
Constant	1.41*** (1.24 ; 1.58)	1.26*** (1.11 ; 1.40)	1.42*** (1.23 ; 1.62)		3.50*** (3.30 ; 3.70)
Dummy Application year	Yes	Yes	Yes	Yes	Yes
Dummy Discipline	Yes	Yes	Yes	Yes	Yes
Dummy Affiliation	Yes	Yes	Yes	Yes	Yes
Appl./Potential appl.	2,120	2,120	2,120	2,120	2,120
Observations	4,240	4,240	4,240	4,240	4,240
R2 / Pseudo R2	0.495	0.240	0.250	0.39	0.444

The table shows a difference-in-differences estimation in the equivalent regression formulation, where the effect of the treatment, i.e., applying for a SINERGIA grant, can be read in the coefficient of the interaction variable *Applicant\*Post-Application*. The five scientific outcomes considered are *Publication count*, *Average IF*, *Average citations*, *Co-applicant collaboration*, and *Learning*. The sample includes 4,240 observations. The 1,060 Applicant-application pairs and the 1,060 Potential applicants selected based on the criterion of exhibiting no statistical difference on the variable *N. of journals pre-application*, i.e., 2,120 Applicant-application pairs/Potential applicants-matched application pairs, are observed in two periods, namely before and after the application year. Columns 1, 2, 3, and 5 report OLS estimates, whereas Column 4 reports the marginal effect of a Probit estimation that considers the binary nature of the dependent variable *Co-applicant collaboration*. Robust standard errors are clustered around the application. In reporting the statistical significance of the coefficients, we apply the standard thresholds, namely \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ . 95% confidence intervals are reported in parenthesis below each coefficient.

### **Joint estimation of the effects of applying for a grant and being awarded a grant**

In the main text, we evaluate the effect of applying for SINERGIA and the effect of being awarded a grant in two separate sets of regressions. Doing so, our estimations measure the average effect of applying mixing two different types of applicants: applicants awarded and non-awarded. Table A2.7 reports an estimated model where we include the two interaction terms *Awarded\*Post-Application* and *Applicant\*Post-Application* in the same model. We find that applicants, regardless of the funding decision, are more productive in terms of number of papers published and of average impact factor, receive fewer citations, learn more, and have a greater probability of establishing co-authorship with their co-applicants. Awarded scientists show no significant differences with non-awarded scientists except for the higher probability of co-authoring with their co-applicants. The results in Table A2.7 are consistent with the results reported in Table 2.6 and Table 2.7 in the main text.



**Table A2.7: Regression results of the joint estimation of the effects of applying for a grant and being awarded a grant**

VARIABLES	(1) OLS Log(Publication count)	(3) OLS Log(Average IF)	(2) OLS Log(Average citations)	(4) Probit Co-applicant Collaboration	(5) OLS Log(Learning)
Applicant*Post-Application	0.41*** (0.34 ; 0.49)	0.058* (-0.0015 ; 0.12)	-0.37*** (-0.45 ; -0.30)	0.14*** (0.063 ; 0.22)	0.37*** (0.28 ; 0.46)
Awarded*Post-Application	0.034 (-0.049 ; 0.12)	0.027 (-0.024 ; 0.079)	0.086 (-0.019 ; 0.19)	0.11*** (0.056 ; 0.17)	-0.018 (-0.14 ; 0.10)
Applicant	0.25*** (0.17 ; 0.32)	-0.083** (-0.16 ; -0.0053)	-0.15*** (-0.24 ; -0.057)	0.44*** (0.37 ; 0.51)	0.058 (-0.025 ; 0.14)
Awarded	-0.032 (-0.12 ; 0.055)	0.12** (0.014 ; 0.22)	0.087 (-0.040 ; 0.21)	-0.022 (-0.067 ; 0.023)	-0.032 (-0.14 ; 0.075)
Post-Application	-0.25*** (-0.30 ; -0.20)	-0.067*** (-0.11 ; -0.019)	-0.15*** (-0.21 ; -0.096)	-0.014 (-0.073 ; 0.046)	-0.62*** (-0.68 ; -0.57)
Seniority	0.031*** (0.027 ; 0.034)	-0.0019 (-0.0046 ; 0.00085)	-0.00049 (-0.0035 ; 0.0026)	-0.0014 (-0.0031 ; 0.00027)	0.012*** (0.0092 ; 0.015)
Other active funding	0.33*** (0.27 ; 0.40)	0.18*** (0.13 ; 0.23)	0.16*** (0.10 ; 0.22)	-0.011 (-0.042 ; 0.020)	0.25*** (0.19 ; 0.30)
Previous expired funding	0.22*** (0.15 ; 0.29)	-0.024 (-0.075 ; 0.026)	-0.067** (-0.13 ; -0.0012)	0.031** (0.0016 ; 0.060)	0.098*** (0.040 ; 0.16)
Constant	1.76*** (1.60 ; 1.91)	1.18*** (1.05 ; 1.31)	0.99*** (0.84 ; 1.14)		3.36*** (3.22 ; 3.51)
Dummy Application year	Yes	Yes	Yes	Yes	Yes
Dummy Discipline	Yes	Yes	Yes	Yes	Yes
Dummy Affiliation	Yes	Yes	Yes	Yes	Yes
Appl./Potential appl.	2,120	2,120	2,120	2,120	2,120
Observations	4,240	4,240	4,240	4,240	4,240
R2 / Pseudo R2	0.435	0.265	0.246	0.39	0.466

The table shows a difference-in-differences estimation, in the equivalent regression formulation, including two treatments: applying for (*Applicant*) and being awarded (*Awarded*) a SINERGIA grant. The effect of the two treatments can be read, respectively, on the coefficient of the interaction variable *Applicant\*Post-Application* and on the coefficient of the interaction variable *Awarded\*Post-Application*. The controls include fixed effects for *Application year*, *Affiliation*, and *Discipline*. The five scientific outcomes considered are *Publication count*, *Average IF*, *Average citations*, *Co-applicant collaboration*, and *Learning*. The sample includes 4,240 observations. The 1,060 Applicant-application pairs and the 1,060 Potential applicants, i.e., 2,120 Applicant-application pairs/Potential applicants, are observed in two periods, before and after the application year. Columns 1, 2, 3, and 5 report OLS estimates, whereas Column 4 reports the marginal effects of a Probit estimation that considers the binary nature of the dependent variable *Co-applicant collaboration*. Robust standard errors are clustered around the application. In reporting the statistical significance of the coefficients, we apply the standard thresholds, namely \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ . 95% confidence intervals are reported in parenthesis below each coefficient.

**Estimation of the effect of being awarded a grant restricting the sample to the applications graded 3 and 4**

The Regression Discontinuity Design (RDD) approach is not suitable in our case due to the evaluation framework of the program under analysis, SINERGIA. The main reason is that a fine-grain evaluation scale is not available: applications receive a score from 1 to 6, and it is not possible to rank the quality of applications receiving the same grade. To test the robustness of our results with an approach similar to RDD, we restrict our sample to applications of comparable quality around the pay line, i.e., those graded 3 or 4. Table A2.8 shows the results of our exercise. The results are in line with the ones reported in Table 2.7 in the main text.

**Table A2.8: Regression results of the effect of being awarded a grant restricting the sample to applications graded 3 and 4.**

VARIABLES	(1) OLS Log(Publication count)	(2) OLS Log(Average IF)	(3) OLS Log(Average citations)	(4) Probit Co-applicant Collaboration	(5) OLS Log(Learning)
Awarded*Post-Application	0.049 (-0.097 ; 0.20)	0.065 (-0.030 ; 0.16)	0.026 (-0.14 ; 0.20)	0.15** (0.0095 ; 0.30)	-0.10 (-0.34 ; 0.14)
Awarded	-0.21** (-0.36 ; -0.047)	0.0078 (-0.15 ; 0.17)	0.021 (-0.20 ; 0.24)	-0.14** (-0.28 ; -0.0051)	-0.15* (-0.31 ; 0.017)
Post-Application	0.052 (-0.042 ; 0.15)	-0.035 (-0.095 ; 0.024)	-0.60*** (-0.70 ; -0.49)	0.19*** (0.082 ; 0.29)	-0.36*** (-0.50 ; -0.23)
Seniority	0.023*** (0.015 ; 0.031)	0.0047 (-0.0015 ; 0.011)	0.0083* (-0.00077 ; 0.017)	-0.0043 (-0.0098 ; 0.0011)	0.011*** (0.0042 ; 0.018)
Female	-0.13 (-0.30 ; 0.030)	-0.054 (-0.21 ; 0.10)	0.097 (-0.080 ; 0.28)	-0.0091 (-0.17 ; 0.15)	-0.11 (-0.25 ; 0.042)
Other active funding	0.27*** (0.14 ; 0.40)	0.10* (-0.00013 ; 0.20)	0.12 (-0.038 ; 0.27)	0.063 (-0.035 ; 0.16)	0.17*** (0.078 ; 0.27)
Previous expired funding	0.22*** (0.068 ; 0.37)	0.024 (-0.082 ; 0.13)	-0.077 (-0.21 ; 0.062)	0.051 (-0.057 ; 0.16)	0.11 (-0.022 ; 0.23)
Swiss team	-0.048 (-0.31 ; 0.22)	0.086 (-0.13 ; 0.30)	0.028 (-0.19 ; 0.24)	0.079 (-0.11 ; 0.26)	0.18** (0.019 ; 0.34)
At least one female researcher	-0.018 (-0.13 ; 0.095)	-0.081 (-0.22 ; 0.058)	-0.040 (-0.23 ; 0.15)	0.093 (-0.027 ; 0.21)	-0.088 (-0.25 ; 0.074)
Log(Amount Requested)	0.070 (-0.095 ; 0.23)	0.27*** (0.070 ; 0.47)	0.37*** (0.10 ; 0.63)	-0.30*** (-0.45 ; -0.15)	-0.046 (-0.27 ; 0.18)
Log(N. of co-applicants)	0.074 (-0.11 ; 0.25)	0.0035 (-0.20 ; 0.20)	0.12 (-0.18 ; 0.42)	0.24** (0.042 ; 0.43)	0.043 (-0.18 ; 0.26)
Log(N. of disciplines)	0.058 (-0.046 ; 0.16)	0.095 (-0.037 ; 0.23)	-0.017 (-0.19 ; 0.16)	-0.055 (-0.13 ; 0.021)	0.21*** (0.093 ; 0.33)
Science & Medicine	-0.12 (-0.31 ; 0.071)	0.54*** (0.31 ; 0.76)	0.50*** (0.20 ; 0.79)	-0.12 (-0.32 ; 0.078)	0.76*** (0.51 ; 1.02)
Log(1+Distance hours)	-0.0079 (-0.092 ; 0.076)	0.14** (0.030 ; 0.25)	0.068 (-0.062 ; 0.20)	-0.10** (-0.19 ; -0.012)	0.13*** (0.034 ; 0.22)
Previous SINERGIA application	0.13* (-0.014 ; 0.27)	-0.082 (-0.24 ; 0.073)	-0.24** (-0.48 ; -0.0020)	-0.0036 (-0.16 ; 0.15)	0.19*** (0.049 ; 0.33)
Previous SINERGIA awarded	-0.011 (-0.19 ; 0.17)	-0.030 (-0.26 ; 0.20)	0.10 (-0.22 ; 0.43)	0.17* (-0.026 ; 0.36)	-0.058 (-0.29 ; 0.17)
Constant	1.20 (-1.19 ; 3.59)	-3.42** (-6.38 ; -0.45)	-5.09** (-8.96 ; -1.22)		3.45** (0.39 ; 6.51)
Dummy Application year	Yes	Yes	Yes	Yes	Yes
Dummy Discipline	Yes	Yes	Yes	Yes	Yes
Dummy Affiliation	Yes	Yes	Yes	Yes	Yes
Applicant-application pairs	356	356	356	356	356
Observations	712	712	712	712	712
R2 / Pseudo R2	0.475	0.514	0.469	0.23	0.661

The table shows a difference-in-differences estimation in the equivalent regression formulation, where the effect of the treatment, i.e., being awarded a SINERGIA grant, can be read in the coefficient of the interaction variable *Awarded\*Post-Application*. Robust standard errors are clustered around the application. In reporting the statistical significance of the coefficients, we apply the standard thresholds, namely \*\*\* p<0.01, \*\* p<0.05, \* p<0.1. 95% confidence intervals are reported in parenthesis below each coefficient.

### Poisson difference-in-differences estimation for *Publication count* and *Learning* in applicants vs. potential applicants regression exercise

The count variable nature of *Publication count* and *Learning* allows us to choose an alternative estimation method based on a Poisson difference-in-differences. We report this alternative estimation in Table A2.9 to test the robustness of the results reported in the main text. We find that the results in Table A2.9 are consistent with the ones reported in the main text in Table 2.6.

**Table A2.9: Poisson difference-in-differences estimation for the variables *Publication count* and *Learning*.**

	(1) Poisson Publication count	(2) Poisson Learning
Applicant*Post-Application	0.15*** (0.095 ; 0.21)	0.21*** (0.15 ; 0.27)
Applicant	0.096*** (0.033 ; 0.16)	0.087*** (0.040 ; 0.13)
Post-Application	-0.0066 (-0.048 ; 0.035)	-0.49*** (-0.53 ; -0.45)
Seniority	0.024*** (0.021 ; 0.027)	0.010*** (0.0081 ; 0.012)
Other active funding	0.27*** (0.20 ; 0.34)	0.16*** (0.12 ; 0.21)
Previous expired funding	0.18*** (0.11 ; 0.25)	0.045* (-0.00065 ; 0.091)
Constant	2.44*** (2.30 ; 2.59)	3.88*** (3.76 ; 3.99)
Dummy Application year	Yes	Yes
Dummy Discipline	Yes	Yes
Dummy Affiliation	Yes	Yes
Appl./Potential appl.	2,120	2,120
Observations	4,240	4,240
R2 / Pseudo R2	0.35	0.43

The table shows a difference-in-differences estimation in the equivalent regression formulation, where the effect of the treatment, i.e., applying for a SINERGIA grant, can be read in the coefficient of the interaction variable *Applicant\*Post-Application*. Two scientific outcomes are considered: *Publication count*, and *Learning*. The sample includes 4,240 observations. The 1,060 Applicant-application pairs and the 1,060 Potential applicants, i.e., 2,120 Applicant-application pairs/Potential applicants-matched application pairs, are observed in two periods, namely before and after the application year. Columns 1 and 2 report Poisson estimates for the two considered outcomes. Robust standard errors are clustered around the application. In reporting the statistical significance of the coefficients, we apply the standard thresholds, namely \*\*\*  $p < 0.01$ , \*\*  $p < 0.05$ , \*  $p < 0.1$ . 95% confidence intervals are reported in parenthesis below each coefficient.

## Appendix A3 for chapter 3:

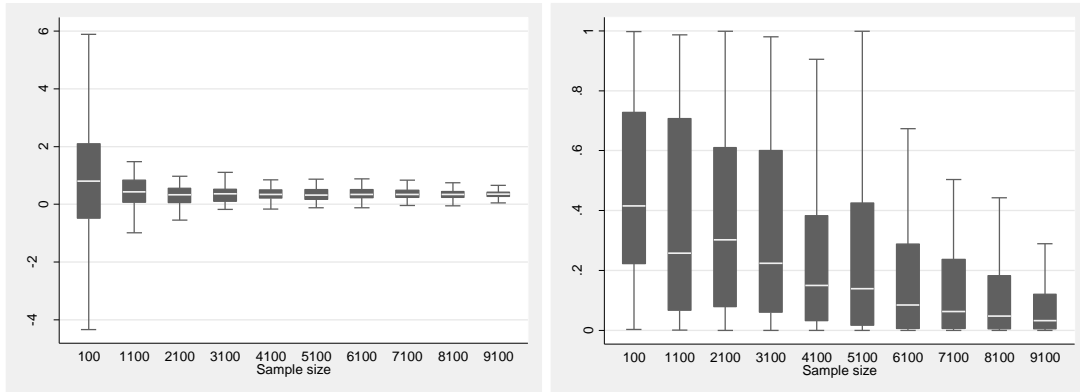
### **The statistical and economic significance of our results on cognitive distance.**

When econometric analyses are conducted on large samples, the standard econometric levels of significance of the estimated coefficients have to be treated carefully. With large samples, even regression coefficients with a negligible economic impact (i.e., a small size of the coefficient) might result in being statistically significant. In this paper, we have the case of a large sample of 106,898 observations (Table 3.5, column 4). In this appendix, we comment on the significance of the coefficients estimated for one of our main variables of interest. Figure 3.3 shows that the prediction of the probability of learning from within the team varies to an economically significant extent, from a probability of approximately 56% for low cognitive distance, up to 63% for medium cognitive distance and 40% for high cognitive distance<sup>92</sup>. We go beyond the statistical significance of the coefficient by considering a Monte-Carlo CPS chart to show that the impact of the linear and quadratic components of the cognitive distance is already significant for smaller samples randomly drawn (Lin et al., 2013). The Monte-Carlo CPS approach draws random observations for different sample sizes, ranging from 100 to 9,100 observations. It extracts 100 random samples for each sample size. Our econometric model is then estimated for the extracted samples. In Figure A3.1, we report the boxplots of the P-values of the linear term of the cognitive distance for each sample size. In Figure A3.2, we report the boxplots of the P-values of the quadratic term of the cognitive distance for each sample size. We find that the P-values converge very quickly to the values observed in the complete sample, well before reaching the 106,898 observations used in the regression in Table 3.5, column 4. Both the economically significant extent of the impact of cognitive distance and the Monte-Carlo CPS chart confirm that the impact of the cognitive distance is not a purely statistical artifact.

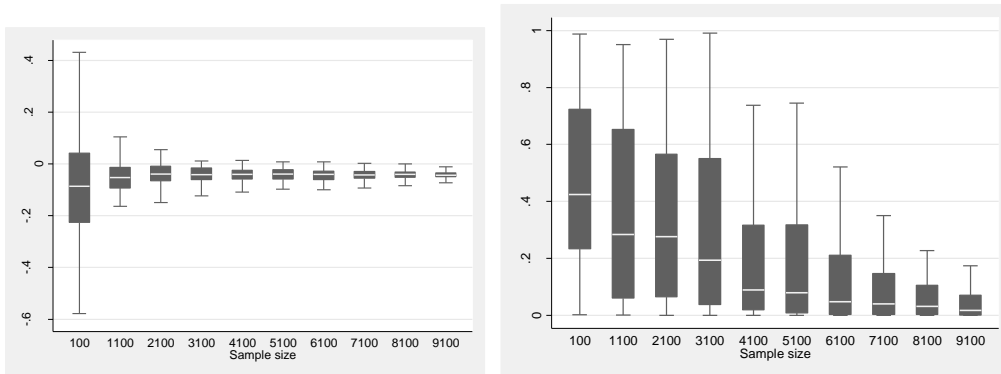
---

<sup>92</sup> All these figures are calculated for the representative individual described in Figure 3.3.

**Figure A3.1: Boxplots of the extent of the estimated coefficient (left side) and of the P-value (right side) of the linear term of the cognitive distance for each sample size.**



**Figure A3.2: Boxplots of the extent of the estimated coefficient (left side) and of the P-value (right side) of the quadratic term of the cognitive distance for each sample size.**



## Appendix A4 for chapter 4:

### Appendix A4.1: Proofs of theorems and propositions.

#### Proof of Theorem 1 (Sharing strategy)

*A homo moralis individual with a degree of morality  $\kappa_i$  involved in a sharing social dilemma plays a pure strategy ( $x_i \in \{0,1\}$ ) if and only if:  $\kappa_i \leq \kappa^L$  ( $x_i = 0$ ) or  $\kappa_i \geq \kappa^H$  ( $x_i = 1$ ).*

Where,  $\kappa^L = \inf_{(0,1]}[\psi_0]$ ,  $\kappa^H = \sup_{(0,1)}[\psi_1]$ , and  $\forall u \in (0,1), \psi_u(x) = \frac{C(x) - C(u)}{\xi(x) - \xi(u)}$

*Proof.* Following the maximization program (4.5), we know that  $x_i$  maximize the function defined for  $x$  in  $[0,1]$  as:  $u(x) = \kappa_i \cdot \xi(x) - C(x)$ . Therefore, noting that  $u(0) = 0$  by construction of  $C$  and  $\xi$ , individual  $i \in I$ , plays the pure strategy  $x_i = 0$  if and only if:  $\forall x \in (0,1], u(x) \leq 0$ . Hence, replacing  $u(\cdot)$  with its value, and noting that  $\xi(\cdot)$  is strictly positive on  $(0,1]$ , individual  $i \in I$ , plays the pure strategy  $x_i = 0$  if and only if:  $\forall x \in (0,1], \kappa_i \leq \frac{C(x)}{\xi(x)}$ . Considering the definition of  $\psi_0(\cdot)$ , this last condition is equivalent to  $\kappa_i \leq \inf_{(0,1]} \psi_0$ .

Similarly, individual  $i \in I$  plays the pure strategy  $x_i = 1$ , if and only if  $u(\cdot)$  is maximized in 1. In other words,  $x_i = 1$  if and only if:  $\forall x \in [0,1), u(x) \leq u(1)$ . Hence, knowing the formulation of function  $u$ , noting that  $(\xi(\cdot) - \xi(1))$  is strictly negative on  $[0,1)$  (because  $\xi(\cdot)$  is strictly increasing on  $[0,1]$ ), and considering the definition of  $\psi_1(\cdot)$ , we have that individual  $i \in I$ , plays the pure strategy  $x_i = 1$  if and only if:  $\kappa_i \geq \sup_{[0,1)} \psi_1$ .  $\square$

#### Proof of Corollary 1 (Population with pure strategies only)

*If the functions  $\psi_0$  and  $\psi_1$  are decreasing on  $[0,1]$ , then all the individuals in the population play pure strategies, and we have, for all  $i \in I$ :  $x_i = 0$  if  $\kappa_i \leq \psi_0(1)$  and  $x_i = 1$  otherwise.*

*Proof.* Since the functions  $C$  and  $\xi$  are continuous on  $[0,1]$ , the functions  $\psi_0$  and  $\psi_1$  are also continuous on  $[0,1]$ . Moreover, the interval  $[0,1]$  is bounded, therefore, according to the extreme value theorem, the continuous functions  $\psi_0$  and  $\psi_1$  attain their minimum and maximum on this interval. Furthermore, if the functions  $\psi_0$  and  $\psi_1$  are decreasing on  $[0,1]$ , then they attain their minimum in 1 and their maximum in 0. We thus have that:  $\kappa^L = \inf_{(0,1]} \psi_0 = \psi_0(1)$  and  $\kappa^H = \sup_{[0,1)} \psi_1(x) = \psi_1(0)$ .

Hence, noting that  $\psi_0(1) = \psi_1(0) = \frac{C(1)}{\xi(1)}$ , we have that:  $\kappa^L = \kappa^H = \frac{C(1)}{\xi(1)} = \psi_0(1)$ . Therefore, with theorem 1, we have that an individual  $i$  in  $I$  plays the pure strategy  $x_i = 0$  if  $\kappa_i \leq \psi_0(1)$  and plays the pure strategy  $x_i = 1$  if  $\kappa_i \geq \psi_0(1)$ . And no individual plays an interior strategy.  $\square$

**Proof of Theorem 2**

A homo moralis individual with a degree of morality  $\kappa_i$  involved in a sharing social dilemma has an interior strategy if and only if:  $\kappa^L < \kappa_i < \kappa^H$ . And the degree of sharing  $x_i \in (0,1)$  is then solution to:

$$\frac{\partial C(x_i)}{\partial x_i} = \kappa_i \frac{\partial \xi(x_i)}{\partial x_i}$$

*Proof.* Following Theorem 1, we know that the maximization problem (4.5) has a corner solution if and only if  $\kappa_i \leq \kappa^L$  or  $\kappa_i \geq \kappa^H$ . Therefore, an individual  $i$  in  $I$  has an interior strategy if and only if  $\kappa^L < \kappa_i < \kappa^H$  and the solution to the program in (4.5) is an interior solution  $x_i \in (0,1)$ . When we have an interior solution to a maximization program, it respects the first-order condition. In other words, noting for  $x \in [0,1]$   $u(x) = \kappa_i \cdot \xi(x) - C(x)$ , we have: When  $x_i \in (0,1)$  then it verifies  $\frac{\partial u(x_i)}{\partial x_i} = 0$ . Hence, replacing  $u$  with its value and noting that  $\frac{\partial \xi(x_i)}{\partial x_i}$  is strictly positive (because  $\xi$  is differentiable and strictly increasing on  $[0,1]$ ), we have that: When  $i$  plays an interior strategy,  $x_i \in (0,1)$  is solution to:  $\frac{\partial C(x_i)}{\partial x_i} = \kappa_i \frac{\partial \xi(x_i)}{\partial x_i}$ .  $\square$

**Appendix A4.2. Applications**

Application 2: Quadratic individual cost

Individuals' strategy:

- Detailed calculation of  $\kappa^H = \frac{2\gamma}{\beta_n}$ :

*Proof.* According to theorem 1, we know that  $\kappa^H = \sup_{[0,1]}[\psi_1(x)]$ .

Thus, in the case of Application 2, with the definition of  $C$  and  $\xi$  we have:

$$\forall x \in [0,1], \psi_1(x) = \frac{\gamma(x^2-1)}{\beta_n(x-1)} = \frac{\gamma(x+1)}{\beta_n}$$

$\psi_1$  is, therefore, a linear function attaining its maximum in  $x = 1$ , and we have  $\kappa^H = \psi_1(1) = \frac{2\gamma}{\beta_n}$ .  $\square$

- Comparative statics

Independently of the distribution of morality in the population, the average level of sharing in the population is decreasing in the cost factor  $\gamma$  and increasing in the size of the population  $n$ :

$$\frac{d\bar{x}}{d\gamma} = -\frac{\beta_n}{2\gamma^2} E[\kappa \mid \kappa < 2\gamma/\beta_n] < 0$$

$$\frac{d\bar{x}}{dn} = \frac{d\bar{x}}{d\beta_n} \frac{d\beta_n}{dn} = \frac{d\beta_n}{dn} \frac{1}{2\gamma} E[\kappa \mid \kappa < 2\gamma/\beta_n] > 0$$

*Proof.* Recall that the average degree of cooperation in the population satisfies:

$$\bar{x} = \frac{\beta_n}{2\gamma} E[\kappa \mid \kappa < 2\gamma/\beta_n] + 1 - F_\kappa(2\gamma/\beta_n)$$



When  $2\gamma \geq \beta_n$ , the average degree of sharing in the population simplifies to:  $\bar{x} = \frac{\beta_n}{2\gamma} E(\kappa)$ , and the above result is straightforward. When  $2\gamma < \beta_n$ , we use the Leibniz's rule for differentiation under the integral sign, and we have:

$$\begin{aligned} \frac{d\bar{x}}{d\gamma} &= -\frac{\beta_n}{2\gamma^2} E[\kappa \mid \kappa < 2\gamma/\beta_n] + \frac{\beta_n}{2\gamma} \frac{2\gamma}{\beta_n} f_\kappa(2\gamma/\beta_n) \frac{2}{\beta_n} - f_\kappa(2\gamma/\beta_n) \frac{2}{\beta_n} \\ &= -\frac{\beta_n}{2\gamma^2} E[\kappa \mid \kappa < 2\gamma/\beta_n] < 0 \end{aligned}$$

The calculation to derive  $\frac{\partial \bar{x}}{\partial n}$  follows the same logic. □

Application 3: Quadratic individual cost and peer pressure

- Individuals' strategy

Taking the functions of cost and externality defined in equations (4.6) and (4.7), and setting, for  $x \in [0,1]$ ,  $P(x, \bar{x}) = \delta(x - \bar{x})^2$ , the program defined in (4.9) becomes:

$$x_i \in \operatorname{argmax}_{x \in [0,1]} [x(a - bx)]$$

$$\text{where } a = \beta_n \kappa_i - \gamma + 2\bar{x}\delta(1 - \kappa_i)$$

$$\text{and } b = \delta(1 - \kappa_i)$$

A *homo moralis* individual, therefore, plays the strategy:

- $x_i = 0$  if and only if  $\kappa_i = 0 = \kappa^L$  and  $\bar{x} = 0$
- $x_i = 1$  if and only if  $\kappa_i \geq \kappa^H = \frac{2\gamma + 2(1 - \bar{x})\delta}{\beta_n + 2(1 - \bar{x})\delta}$
- Otherwise  $x_i = \frac{2\bar{x}\delta(1 - \kappa_i) + \beta_n \kappa_i}{2\bar{x}\delta(1 - \kappa_i) + 2\gamma}$ .

*Proof.* First, note that  $d > 0$  because  $\gamma$  and  $\delta$  are strictly positives. Thus, the  $u : x \rightarrow x(c - dx)$  is an inverted U-curve function, passing by the origin, with a global maximum that is attained in  $x = c/2d$ . Hence, since  $c/2d \geq 0$  the function  $u$  is maximized in 0 on  $[0,1]$  if and only if  $c = 0$ , which is only the case when  $\kappa_i = 0$  and  $\bar{x} = 0$ . Then, when  $c/2d > 0$ ,  $x_i = 1$  if and only if the global maximum of  $u$  is attained after 1, i.e., if  $c/2d > 1$ . In other words, we have  $x_i = 1$  if and only if  $c > 2d$  which can be re-ordered as  $\kappa_i \geq \frac{2\gamma + 2(1 - \bar{x})\delta}{\beta_n + 2(1 - \bar{x})\delta}$ . Otherwise,  $u$  reaches its global maximum on  $[0,1]$ , and we have:

$$x_i = \frac{c}{2d} = \frac{2\bar{x}\delta(1 - \kappa_i) + \beta_n \kappa_i}{2\bar{x}\delta(1 - \kappa_i) + 2\gamma}. \quad \square$$

- Comparative statics

Consequently, noting that  $\beta_n > \gamma$ , we observe the following comparative statics:

$$\begin{aligned}\frac{d\kappa^H}{d\gamma} &= \frac{2}{\beta_n + 2(1 - \bar{x})\delta} > 0 \\ \frac{d\kappa^H}{d\beta_n} &= -\frac{2(\gamma + (1 - \bar{x})\delta)}{(\beta_n + 2(1 - \bar{x})\delta)^2} < 0 \\ \frac{d\kappa^H}{d\delta|_i} &= \frac{2(1 - \bar{x})(\beta_n - 2\gamma)}{(\beta_n + 2(1 - \bar{x})\delta)^2} < 0 \Leftrightarrow \beta_n < 2\gamma\end{aligned}$$

Moreover, if  $i \in I$  plays an interior strategy, we have:

$$\begin{aligned}\frac{dx_i}{d\gamma} &= -\frac{2\bar{x}\delta(1 - \kappa_i) + \beta_n\kappa_i}{2(\delta(1 - \kappa_i) + \gamma)^2} < 0 \\ \frac{dx_i}{d\beta_n} &= \frac{\kappa_i}{2(\delta(1 - \kappa_i) + \gamma)} > 0 \\ \frac{dx_i}{d\delta} &= \frac{(1 - \kappa_i)(2\bar{x}\gamma - \beta_n\kappa_i)}{2(\delta(1 - \kappa_i) + \gamma)^2} > 0 \Leftrightarrow \kappa_i < \frac{2\bar{x}\gamma}{\beta_n}\end{aligned}$$

#### Application 5: Linear individual cost and peer pressure

Taking the linear functions of cost and externality, the program defined in (4.9) becomes:

$$x_i \in \operatorname{argmax}_{x \in [0,1]} [x(a - bx)]$$

$$\text{where } a = \beta_n\kappa_i - \gamma + 2\bar{x}\delta(1 - \kappa_i)$$

$$\text{and } b = \delta(1 - \kappa_i)$$

The optimal strategy for an individual  $i \in I$  is the following:

- $x_i = 0$  if and only if  $(\beta_n - 2\bar{x}\delta)\kappa_i < \gamma - 2\bar{x}\delta$
- $x_i = 1$  if and only if  $\kappa_i \geq \kappa^H = \frac{\gamma + 2(1 - \bar{x})\delta}{\beta_n + 2(1 - \bar{x})\delta}$
- Otherwise  $x_i = \bar{x} + \frac{\beta_n - \gamma}{2\delta(1 - \kappa_i)} - \frac{\beta_n}{2\delta}$ .

*Proof.* First, note that  $b = 0$  if and only if  $\kappa_i = 1$ , and then the optimal strategy is  $x_i = 1$  like in section 4.3.2 because  $\beta_n > \gamma$  by assumption of the sharing social dilemma. Otherwise,  $b > 0$  and the function  $v : x \rightarrow x(a - bx)$  is an inverted U-curve, passing by the origin, with a global maximum that is attained in  $x = a/2b$ . Hence, if  $a/2b < 0$  (i.e.

$(\beta_n - 2\bar{x}\delta)\kappa_i < \gamma - 2\bar{x}\delta$ ), then the function  $v$  is decreasing on  $[0,1]$  and  $x_i = 0$ . Similarly, if  $a/2b > 1$  (i.e.  $\kappa_i \geq \frac{\gamma + 2(1 - \bar{x})\delta}{\beta_n + 2(1 - \bar{x})\delta}$ ), then the function  $v$  is increasing on  $[0,1]$  and  $x_i = 1$ . Otherwise,  $a/2b \in (0,1)$

and the function  $v$  reaches its global maximum on  $[0,1]$  and  $x_i = \frac{a}{2b} = \bar{x} + \frac{\beta_n - \gamma}{2\delta(1 - \kappa_i)} - \frac{\beta_n}{2\delta}$   $\square$

An interesting result in this case is related to the strategy of *homo-oeconomicus*. Following the first bullet point result, we know that if  $2\bar{x}\delta > \gamma$  then  $x_i > 0$  for all  $i \in I$ . In other words, if the average level of sharing in the population is high enough and the peer pressure is strong enough relatively to the individual cost, then nobody in the population refrains from sharing, not even *homo-oeconomicus*.

For the rest, noting that  $\beta_n > \gamma$ , we can evaluate the impact of each variable with the following comparative statics:

$$\begin{aligned}\frac{d\kappa^H}{d\gamma} &= \frac{1}{\beta_n + 2(1 - \bar{x})\delta} > 0 \\ \frac{d\kappa^H}{d\beta_n} &= -\frac{\gamma + 2\delta(1 - \bar{x})}{(\beta_n + 2(1 - \bar{x})\delta)^2} < 0 \\ \frac{d\kappa^H}{d\delta} &= \frac{2(1 - \bar{x})(\beta_n - \gamma)}{(\beta_n + 2(1 - \bar{x})\delta)^2} > 0\end{aligned}$$

Moreover, if  $i \in I$  plays an interior strategy, we have:

$$\begin{aligned}\frac{dx_i}{d\gamma} &= -\frac{1}{\delta(1 - \kappa_i)} < 0 \\ \frac{dx_i}{d\beta_n} &= \frac{\kappa_i}{2\delta(1 - \kappa_i)} > 0 \\ \frac{dx_i}{d\delta} &= \frac{\gamma - \beta_n\kappa_i}{2(1 - \kappa_i)\delta^2} > 0 \Leftrightarrow \kappa_i < \gamma/\beta_n \\ \frac{dx_i}{d\bar{x}} &= 1 > 0\end{aligned}$$

## Appendix A5 for chapter 5:

### Appendix A5.1 Methodological appendix

#### *Publication data*

We collected publication data related to the subject of ML using the Scopus database. To collect the publications on ML in Healthcare, we used a query searching for typical combinations of words associated with ML in the titles, keywords, and abstracts of scientific articles in the Scopus database. The choice of combination of words for the query is based on the work of the World Intellectual Property Organization (WIPO) for performing a landscape of innovation in Artificial intelligence. After collecting all publications related to ML, we limit the dataset to the medical fields using the discipline categorization of Scopus. We then limited our query to *Health Sciences* as defined on Scopus. Note that all publication data has been collected on the online platform of Scopus in February 2020.

#### *Patent data*

We collected patent data using the Patstat database. We used the query used by WIPO for the identification of patents in ML. The query selects patents based on the IPC codes associated with ML. After constructing a database of patents on ML, we limited our sample to patents with medical applications by selecting the ones with IPC codes in Medical technology (IPC codes A61B, C, D, F, G, H, J, L, M, N; based on the classification of WIPO<sup>93</sup>). Patents from the same docdb family were counted as one observation. The country data is based on the nationality of the inventors listed on the patent application. The Patstat data availability limits the observations to 2015.

#### *Hospitals' survey*

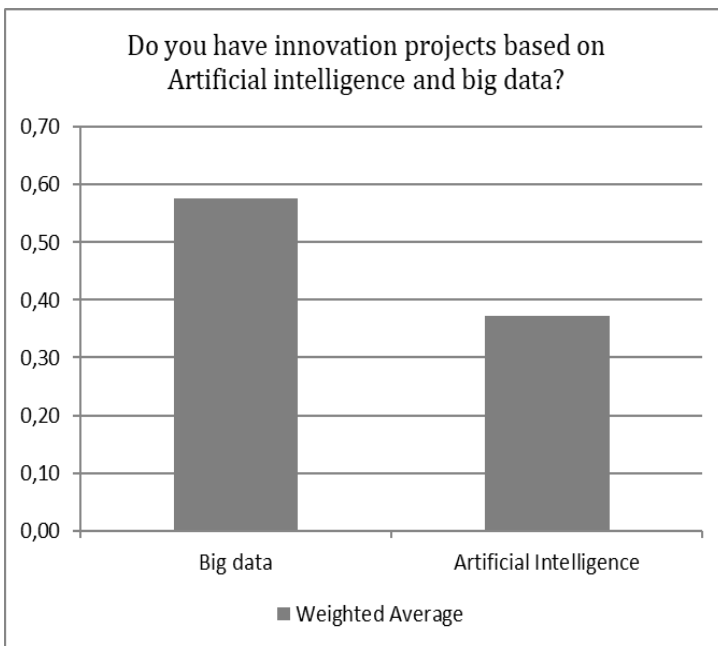
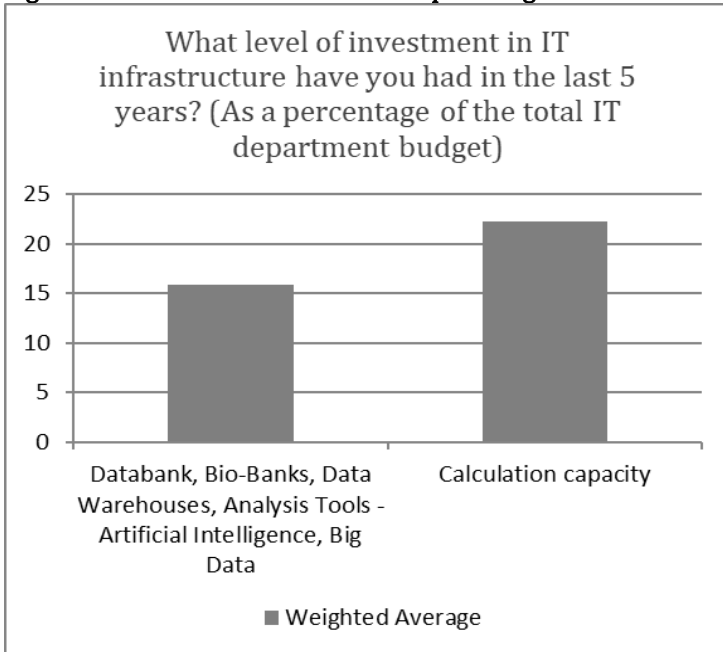
The list of Swiss hospitals was constructed by selecting Swiss medical institutions performing at least 10'000 medical interventions – such as consultations and operations- per year on average. We then retrieved their contact details on their official website one by one. The final list of hospitals was made of 229 institutions all over the country. The survey consisted of 23 questions on the integration of Artificial intelligence and Big data techniques in hospitals. We got 62 replies to the survey with 34 fully completed responses. N.Rosat (DSI-CHUV) provided inputs on a preliminary version of the survey.

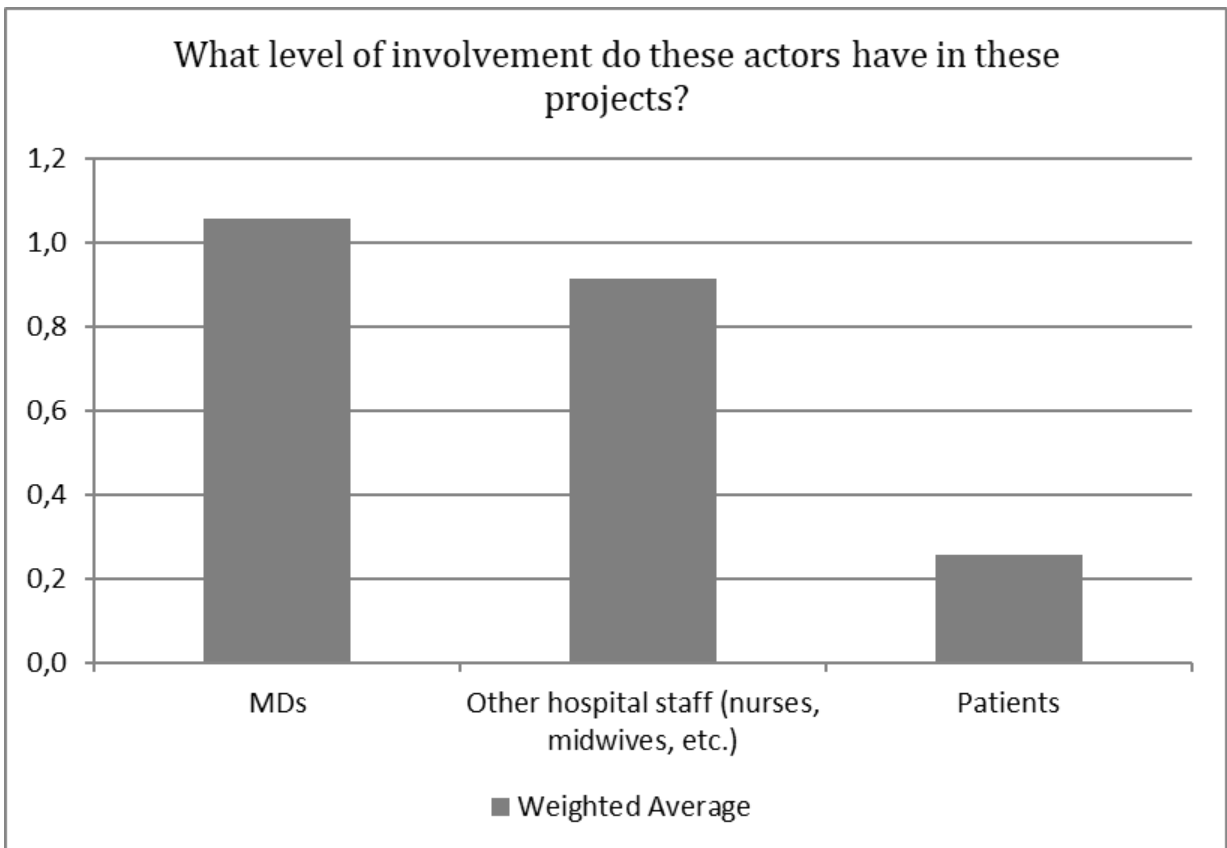
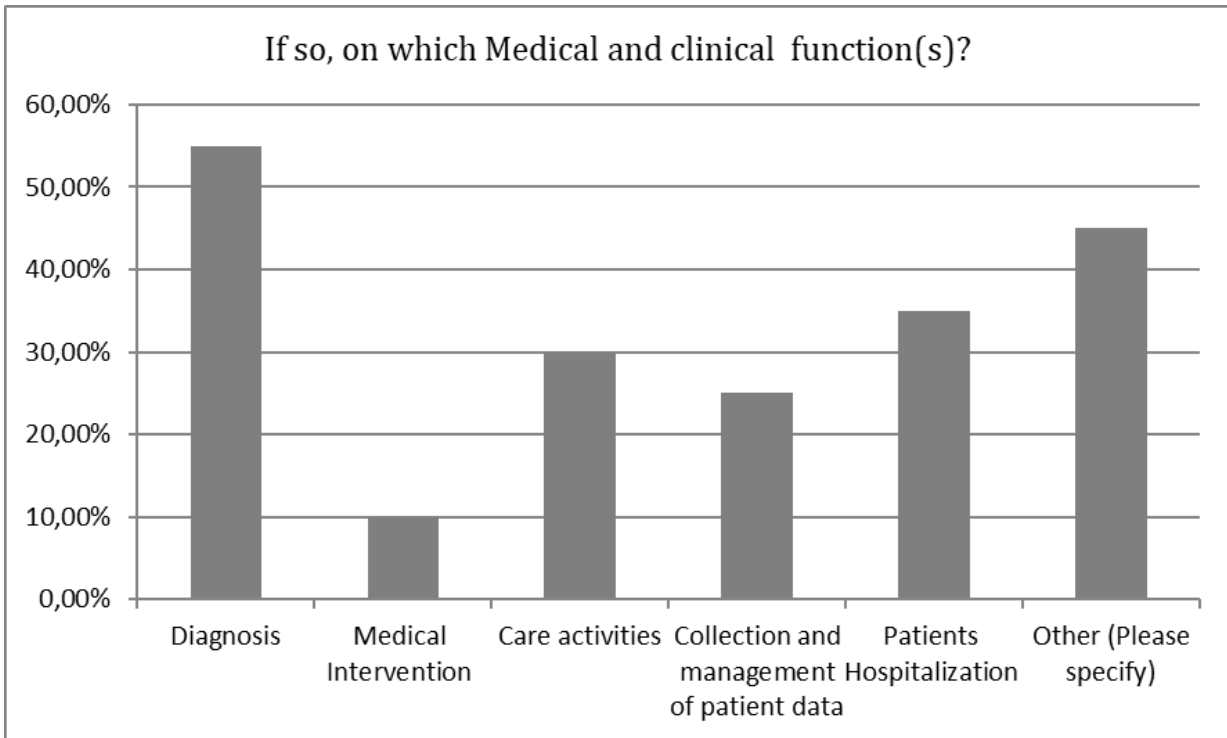
---

<sup>93</sup> See [http://www.wipo.int/export/sites/www/ipstats/en/statistics/patents/pdf/wipo\\_ipc\\_technology.pdf](http://www.wipo.int/export/sites/www/ipstats/en/statistics/patents/pdf/wipo_ipc_technology.pdf)

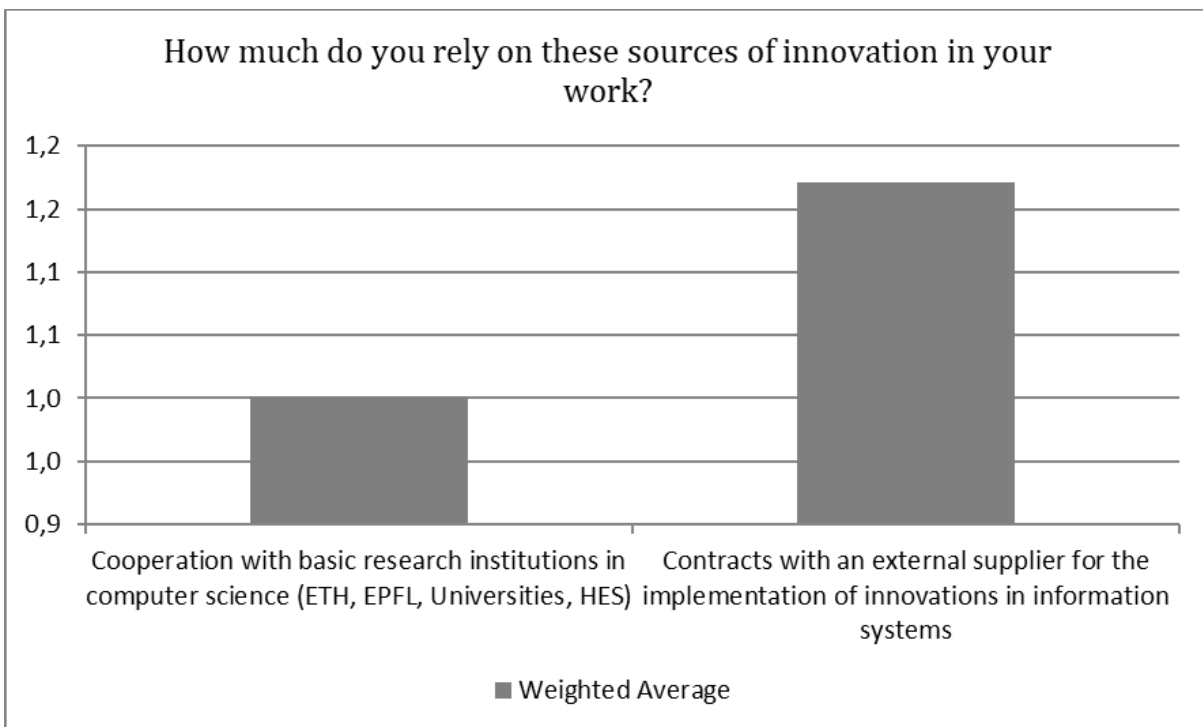
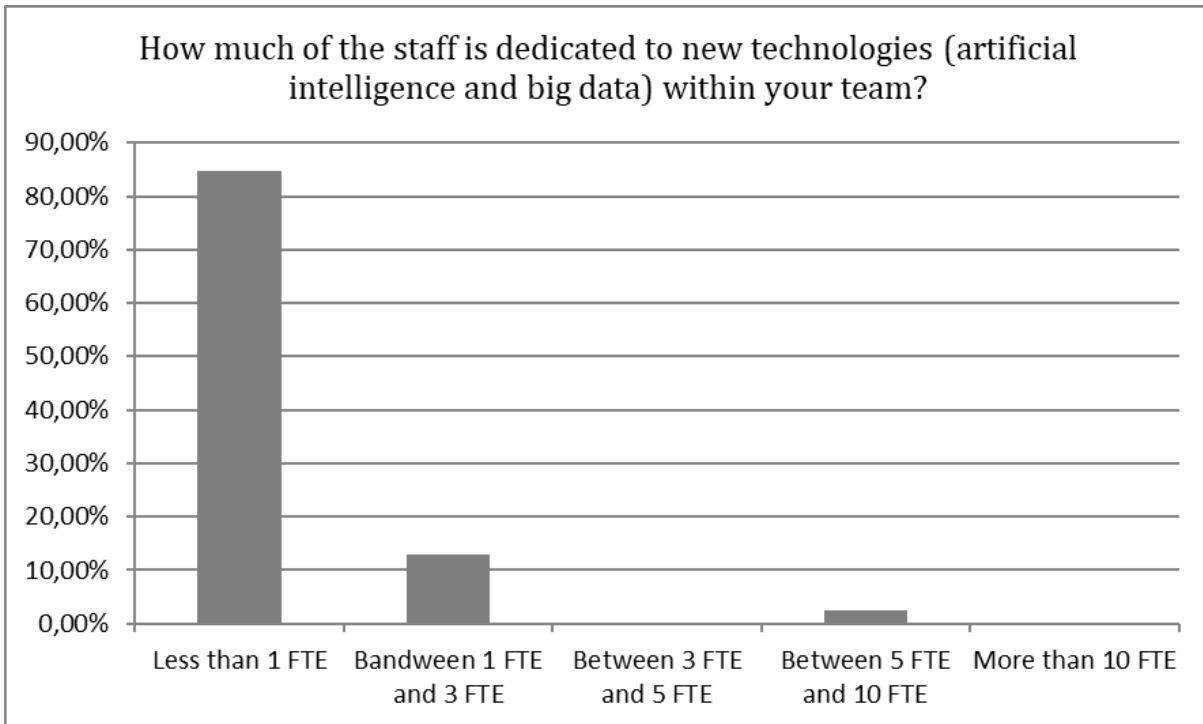
## Appendix A5.2 Survey results

Figures 5.8 to 5.15 – On the Swiss hospitals digital transition

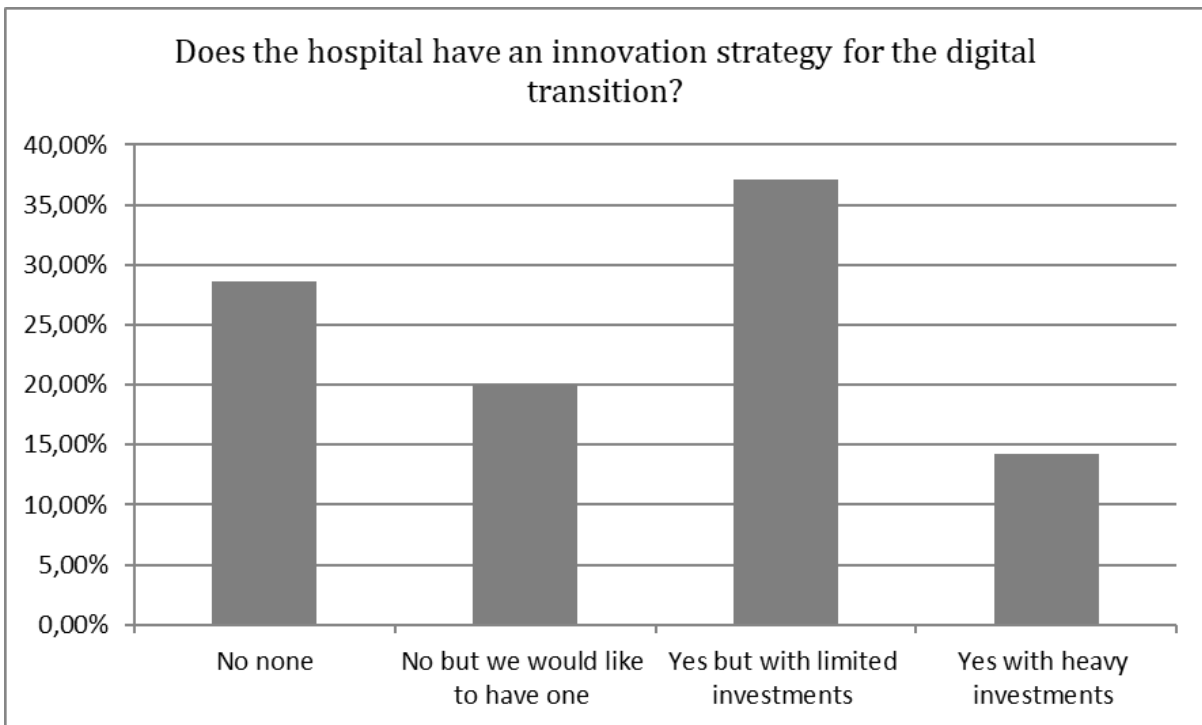
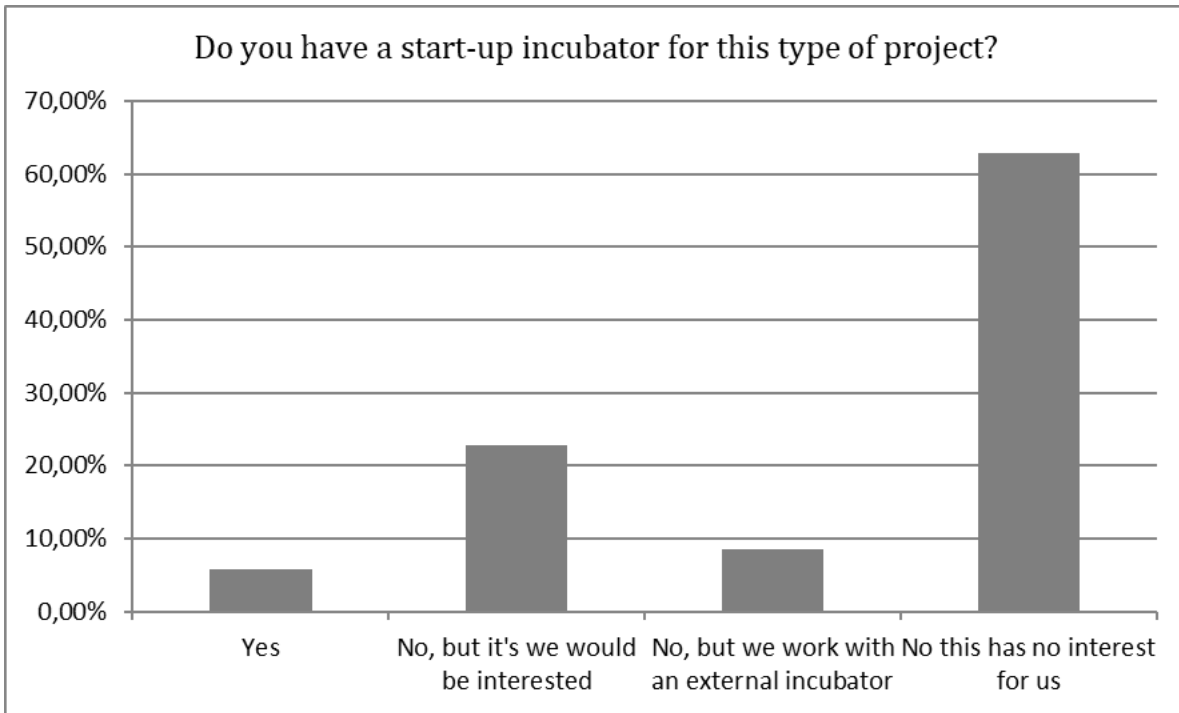




Values between 0 (no implication) and 3 (strong implication)



Values between 0 (not important at all) and 3 (highly important)



Source –authors' survey (see Appendix A5.1)



# Curriculum Vitae

## Charles Ayoubi

EPFL CDM ITPP CEMI  
 ODY 4 16 (Odyssea) Station 5  
 CH-1015 Lausanne, Switzerland  
 +41 78 835 72 15  
[charles.ayoubi@epfl.ch](mailto:charles.ayoubi@epfl.ch)  
 Orcid:0000-0002-0752-3328  
 Born on 13/10/1989 in Toulouse (France)  
 Citizenship: French

## Education

2015- Present	École Polytechnique Fédérale de Lausanne (EPFL)	Ph.D. Candidate in Economics <i>Supervisor:</i> Pr. Dominique Foray
2016	Gerzensee Program, 2016: Microeconomics, Macroeconomics and Econometrics	Doctoral program in Economics <i>Teachers:</i> Bo Honoré, Mark Watson, Jörgen Weibull, Klaus Schmidt
2014- 2015	Université Paris Sud	Master in Economics
2009- 2014	ESSEC Business School	Diplôme de Grande École de Commerce Master of Science in Management
2009- 2013	École Centrale Paris (ECP)	Diplôme de Grande École d'Ingénieur Master in Energy Engineering
2007- 2009	Lycée Louis le Grand	Classes Préparatoires Scientifiques aux Grandes Écoles

## Publications

At the origins of learning: Absorbing knowledge flows from within the team, (with Michele Pezzoni and Fabiana Visentin) *Journal of Economic Behavior and Organization*, 2017 134, 374-387.

The important thing is not to win, it is to take part: What if scientists benefit from participating in competitive grant races?, (with Michele Pezzoni and Fabiana Visentin) *Research Policy*, 2019, vol. 48, no 1, p. 84-97.

COVID-19\_Insights from Innovation Economists (with a collective of scholars primarily associated to the College of Management of Technology at EPFL, under the supervision of Profs Dominique Foray and Gaétan de Rassenfosse). *Science and Public Policy*, 2020, Forthcoming.

## Working papers

Does it Pay to Do Novel Science? The Selectivity Patterns in Science funding (with Michele Pezzoni and Fabiana Visentin). *Under Review in the Journal of Economic Behavior and Organization*

Machine Learning in Healthcare: Mirage or Miracle for Breaking the Costs deadlock? (with Dominique Foray).

Knowledge Diffusion and Morality: Why do we Freely Share Valuable Information with Strangers? (with Boris Thurm)

What matters in funding: The value of research coherence and alignment in evaluators' decisions (with Sandra Barbosu, Michele Pezzoni and Fabiana Visentin).

Exploring the diversity of social preferences: Is a heterogeneous population evolutionarily stable under assortative matching? (with Boris Thurm)

## Selected Work in Progress

Why do individuals care for Nature? (with Boris Thurm).

Learning by Reviewing: The Determinants of Knowledge Acquisition from Scientific Assessments (with Jacqueline Ng Lane).

## Writing for the General Public

It's not the winning but the taking part that counts: how the process of applying for competitive grants is of benefit to researchers ([LSE Impact Blog](#))

Taking part, not winning, counts most in grant applications ([Research Europe](#))

[Simply Applying for a Competitive Grant is a Win \(Social Science Space\)](#)

The SINERGIA Program: A policy evaluation. ([Policy Brief](#)) (2016/01). Dominique Foray, Fabiana Visentin, Michele Pezzoni, Charles Ayoubi & Jacques Mairesse, 2016, EPFL Lausanne.

## Main Conference Presentations

2020: \*(*Forthcoming*) 25th Annual Conference of the European Association of Environmental and Resource Economists (EAERE), Berlin. (Virtual conference due to COVID-19 outbreak).

\*Seminar at Exeter Business School, Science, Innovation, Technology, and Entrepreneurship (SITE), Exeter

2019: \*79th Annual Meeting of the Academy of Management (AoM), Boston, Massachusetts

\*EMAEE Economics, Governance and Management of AI, Robots and Digital Transformations at SPRU, Sussex

\*Organizing in the Era of Digital Technology, ETHZ, Monte Verità, Switzerland

\*Seminar at Boston University's Technology & Policy Research Initiative (TPRI), Boston, Massachusetts

2018: \*MPI Junior Researcher Workshop "From Science to Innovation," Munich, Germany

\*"Advancing the Science of Science Funding" Workshop at the NBER Summer Institute, Boston, Massachusetts

\*BRICK, 12th Workshop on The Organisation, Economics and Policy of Scientific Research, Bath, UK

2017: \*Gerzensee Alumni Conference, Gerzensee, Switzerland

\*Barcelona GSE (Universitat Pompeu Fabra), Barcelona, Spain

2016: \*International Symposium on Science of Science (Library of Congress), Washington DC.

\*DRUID 20<sup>th</sup> Anniversary Conference, Copenhagen, Denmark

## Awards and Grants

12/2018	NBER award on Advancing the Science of Science Funding (\$10,000)
10/2016-10/2018	Best Teaching Assistant for the Master program 2016-2018 Courses in Microeconomics and Economics of Innovation
10/2015-10/2017	Best Teaching Assistant for the Master program 2015-2017 Courses in Microeconomics and Economics of IP
10/2017	Collegio Carlo Alberto award for the paper: At the origins of learning: Absorbing knowledge flows from within the team

## Teaching Experience

Fall, 2018	Master thesis supervision: “Empirical evidence on the role of environmental awareness in the willingness to pay for green electricity.” Mohamed Detsouli (EPFL)
Fall, 2018, 2019	Technology Policy and Energy Transition, Teaching assistant EPFL
Spring, 2018	Master thesis supervision: “Research mapping of the emerging field of geoengineering.” Miguel Gómez Quintanilla (EPFL)
Spring, 2018	Master thesis supervision: “Student migration and environmental awareness.” Louis Delannoy (EPFL)
Spring, 2017	Economics of Innovation, Teaching assistant EPFL
Spring, 2016	Economics of Intellectual Property, Teaching assistant EPFL
Fall, 2015, 2016, 2017	Microeconomics, Teaching assistant EPFL

## Academic Services

Reviewer: *Research Policy*, *Industrial and Corporate Change*, *Journal of Institutional Economics*, *Academy of Management conference (AoM)*, *DRUID Conference*.

Ph.D. Students representative: EPFL, 2017-19

Organization committee: Gerzensee Alumni Conference, Engelberg, Switzerland (April 2018 and 2019)

## Other Skills

Languages: Native in French – Fluent in English, Arabic and Spanish - Intermediate in German

Computer programs/languages: Stata, R, Pascal, Latex, MySQL