# Comparing human and machine performances in transcribing 18th century handwritten Venetian script

Sofia Ares Oliveira

Digital Humanities Laboratory
Ecole Polytechnique Fédérale de Lausanne

DH2018, June 2018, Mexico City

Objective :

Make collections of digitized archival records accessible through textual search

Usually two options for transcription of large collections :

Crowdsourcing

Automatization

Can automatic transcription be as good as crowdsourcing approach ?

# Dataset

Subset of 18th century fiscal documents from the Venetian State Archives

Catastici delle parrocchie

# Dataset

Subset of 18th century fiscal documents from the Venetian State Archives

## Catastici delle parrocchie

Indici

# Dataset

Subset of 18th century fiscal documents from the Venetian State Archives

Catastici delle parrocchie

Indici

Quaderni dei Trasporti

# Dataset

Subset of 18th century fiscal documents from the Venetian State Archives

Catastici delle parrocchie

Indici

Quaderni dei Trasporti

Indici

Francesco ———————— Popiati figlio emancipato d'Angelo

Francesco e Scolastico e

del Cordon

ne in Calle de Fabri

Scola ——— degli ——— Ebrei Levantini di questa Città

Marco, e Agostin Bernar

Piero e Nipoti ————— Mecenati

Bernardo Valier

Nicolò é Fratello

Maria Mazzocco

Pupilli del q: Gio Antonio

Cervò ——————— Civi dal Camel

Girolamo ——— Bernardi di Giuseppe

Giacomo Bembo

Antonio dalla

g: g: di S. Francio di Paola

2ndo e fratti minelli

N. H. Felippo Tansedi

**23 000 units** (image segments) manually transcribed by trained archivists

**54 200 Venetian names** of persons and places

# Machine transcription

vs

Human transcription

# Neural network architecture : CRNN

A Convolutional Recurrent Neural Network combines the best of convolutional and recurrent neural networks



Sequence of descriptors

Per-frame label distribution

Highest probability sequence

B. Shi, X. Bai, and C. Yao, "An end-to-end trainable neural network for image-based sequence recognition and its application to scene text recognition," 2017

Height is fixed for all the image segments (but image ratio is kept)

Data augmentation (contrast, intensity, rotation)

20 712 image segments

48 628 words in total

8 848 vocabulary items

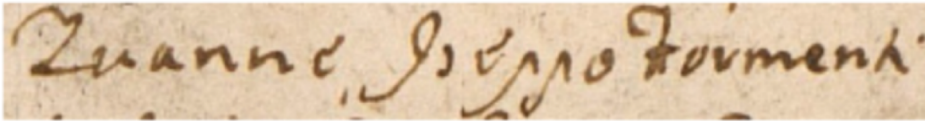'A-Za-z' characters + a few symbols for punctuation

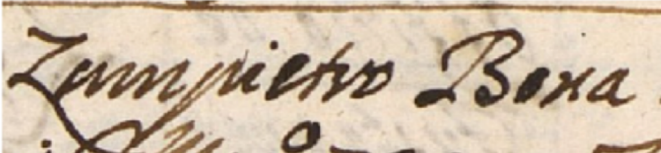# Evaluation



P: Zuanne Iseppo Formenti
GT: Zuanne, Iseppo Formenti

P: Zuanpietro Bona
GT: Zampietro Bona

P: Paulo Padre, e Do=
GT: Paulo Padre, Do=

P: Antonio Bazzerini da Villa Zappa
GT: Antonio Bazzerini da Villa Zoppa

**2 317** image segments
**5 559** words in total
**2 157** vocabulary items

Character Error Rate (CER) : **0.0804**

Machine transcription
vs

Human transcription

http://footage.framepool.com/shotimg/qf/775872743-archivist-historians-librarian-manuscript.jpg

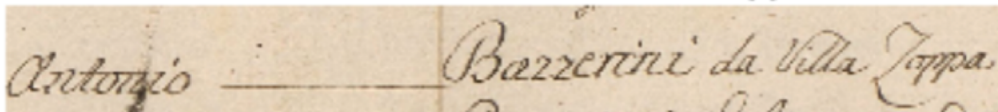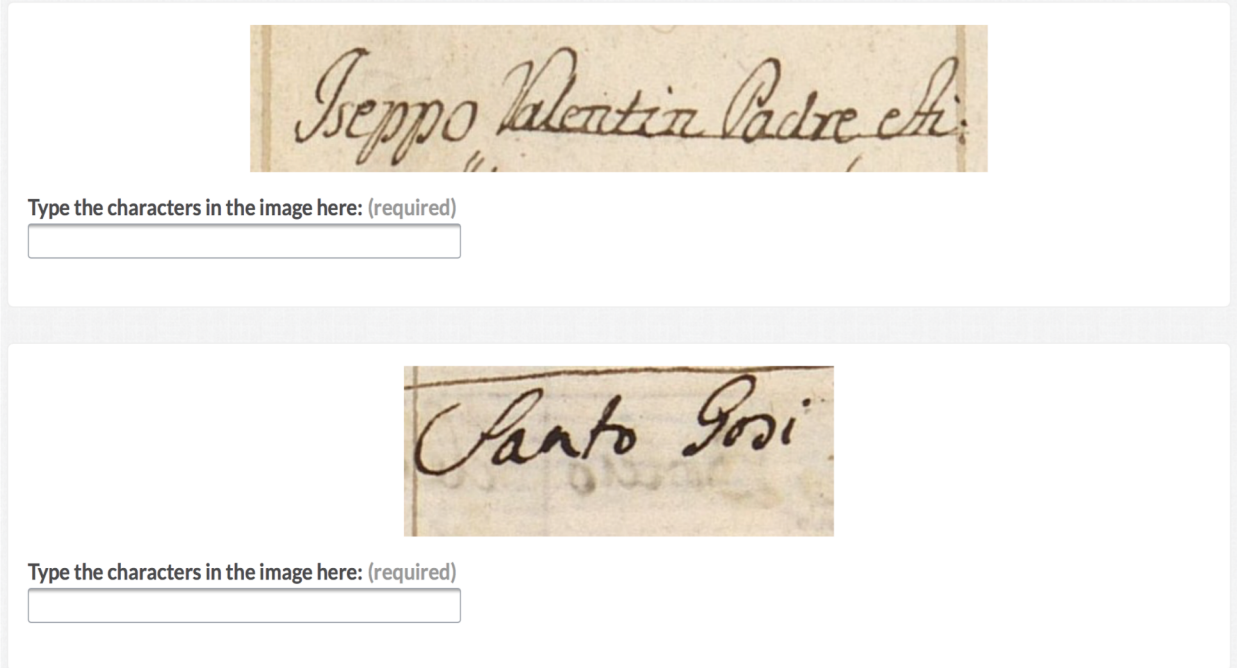# Evaluation of the average performance of Italian-speaking transcribers

Platform : CrowdFlower (now Figure Eight)

Task : Transcribe text in image segment, taking into account capitals and punctuation

Data : 2 317 image segments from test set



Type the characters in the image here: (required)

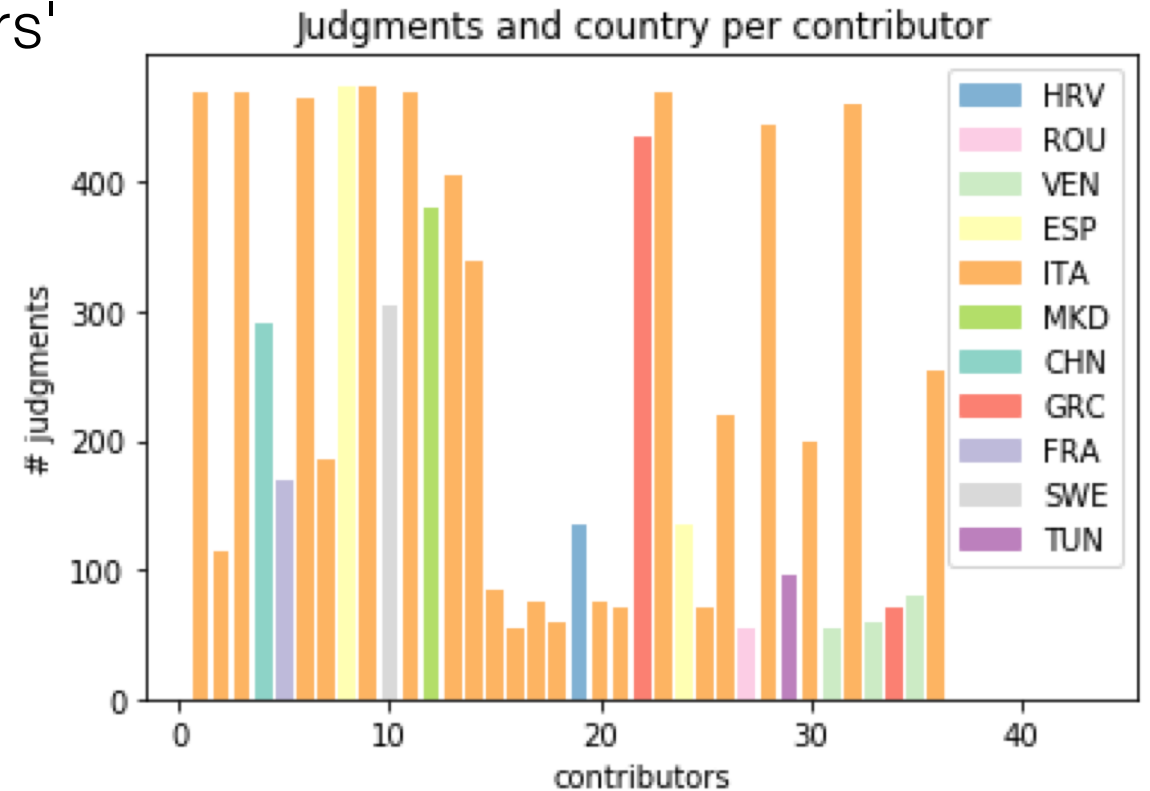Type the characters in the image here: (required)

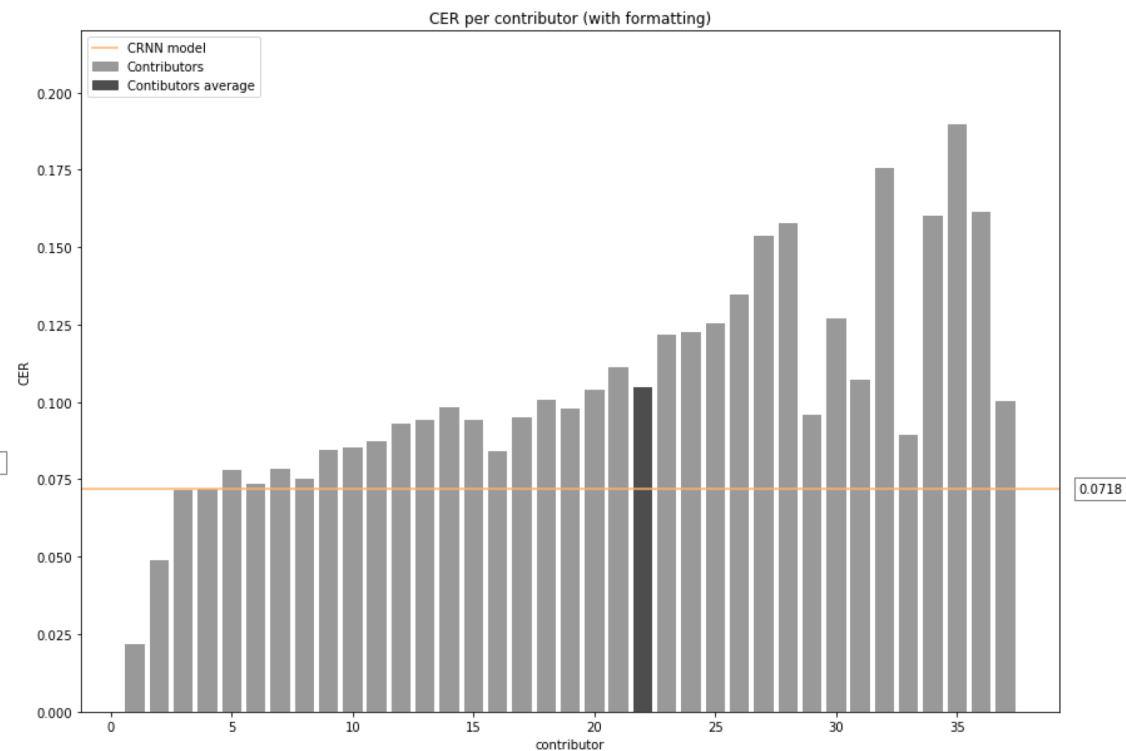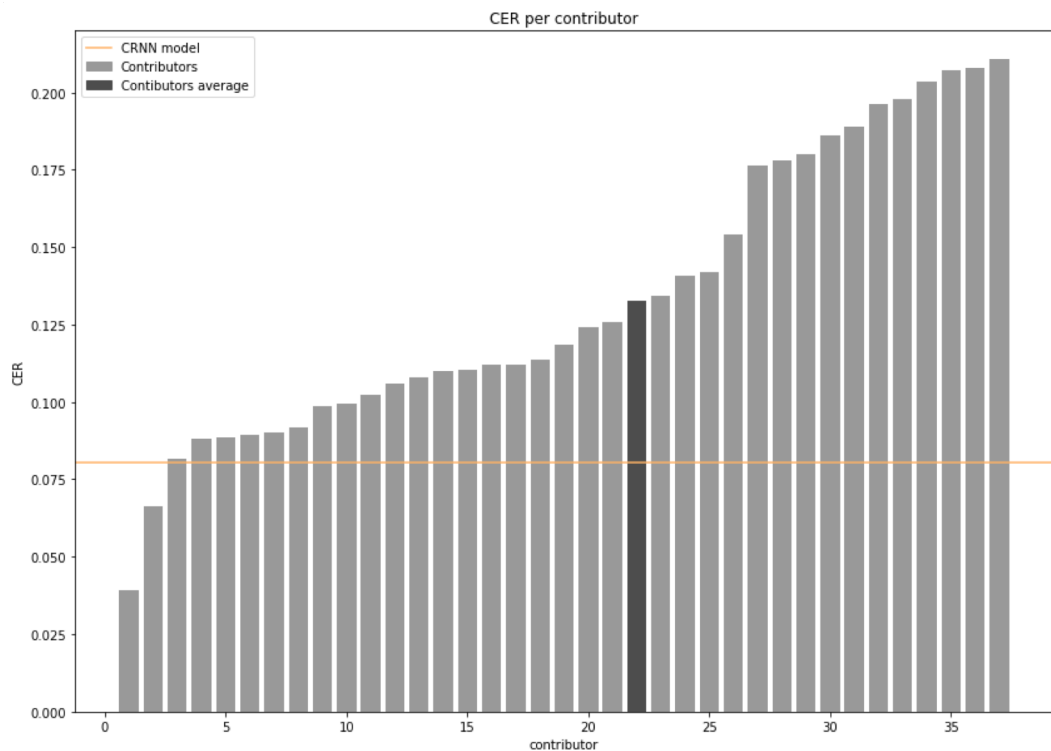Evaluation of the reliability of transcribers' answers during the experiment :

  103 evaluation units
  0.6 accuracy required

36 transcribers remained after selection
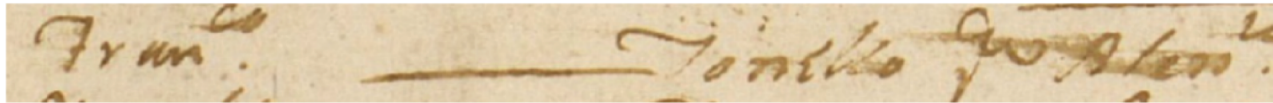
8 674 valid transcriptions to analyze



Judgments and country per contributor

CER per contributor

CER per contributor (with formatting)

|  | CER system | CER amateur | WER system | WER amateur |
|---|---|---|---|---|
| No formatting | 0.0804 | 0.1328 | 0.2709 | 0.4318 |
| With formatting* | 0.0718 | 0.1047 | 0.2551 | 0.3507 |

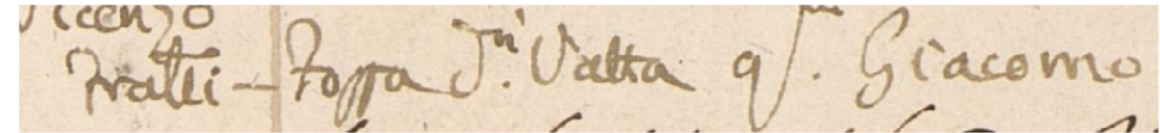*formatting = replace capital letters by lowercase and remove punctuation

# On going work

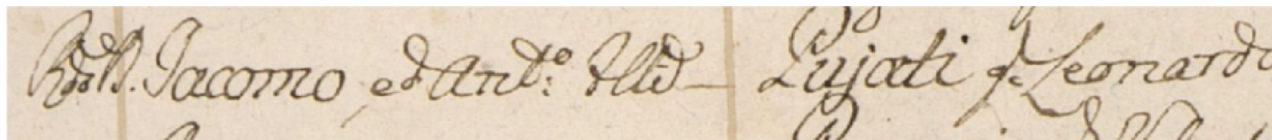Extend the recognition to **abbreviation** symbols

P: francesco Zonello quondam alessandro
GT: francesco Tonello quondam alessandro



P: fratelli Fappa detta Vattin quondam Giacomo
GT: fratelli Foppa detti Vatta quondam Giacomo



P: reverendo don Iacomo, ed antonio fratelli Picati quondam Leonardo
GT: reverendo don Iacomo, ed antonio fratelli Pujati quondam Leonardo



P: 1794 primo maggio
GT: 1794 primo maggio

The system has lower CER and WER than amateur transcribers' average on 18th century Venetian script

→ Sufficiently reliable to use for searching purposes

→ New prospects for analyzing and study large collections of documents

github.com/solivr/tf-crnn

**Venice Time Machine**
vtm.epfl.ch

**Digital Humanities Laboratory**
dhlab.epfl.ch

**Sofia Ares Oliveira**
sofia.oliveiraares@epfl.ch

**Frederic Kaplan**
frederic.kaplan@epfl.ch