

# Integrating Defect Data, Code Review Data, and Version Control Data for Defect Analysis and Prediction

Tao C. Lee

DS Lab, EPFL  
MS thesis defence, Lausanne, Switzerland  
{tao.lee}@epfl.ch

September 11, 2013

**Honeywell**

  
ÉCOLE POLYTECHNIQUE  
FÉDÉRALE DE LAUSANNE

# A software manager's view of defects

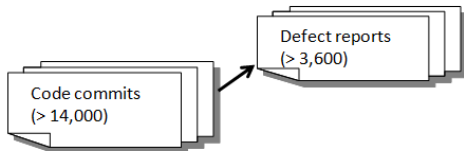
8/1979 projects, Honeywell ACS

Defect reports  
( $> 3,600$ )



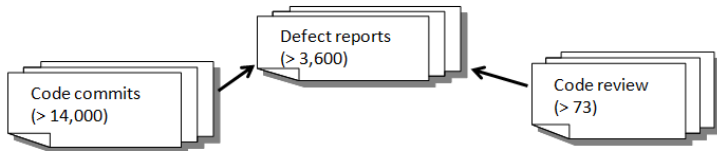
# A software manager's view of defects

8/1979 projects, Honeywell ACS



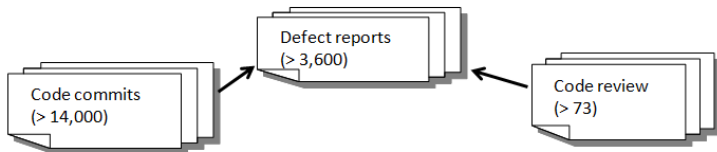
# A software manager's view of defects

8/1979 projects, Honeywell ACS



# A software manager's view of defects

8/1979 projects, Honeywell ACS



**How to make sense of  
this huge amount of  
information?**



# A software manager's view of defects

- Diversified projects: PLs (C, Java), applications (server, control)

# A software manager's view of defects

- Diversified projects: PLs (C, Java), applications (server, control)
- Code reviews are rare

# A software manager's view of defects

- Diversified projects: PLs (C, Java), applications (server, control)
- Code reviews are rare
- Defect types are often unknown



# A software manager's view of defects

- Diversified projects: PLs (C, Java), applications (server, control)
- Code reviews are rare
- Defect types are often unknown
- Prediction of hot files of defects is of interest

# A software manager's view of defects

- Diversified projects: PLs (C, Java), applications (server, control)
- Code reviews are rare
- Defect types are often unknown
- Prediction of hot files of defects is of interest
- Methodology must apply to large defect data of diversified projects

- Prior work

# Outline

- Prior work
- Research contributions

# Outline

- Prior work
- Research contributions
- Methodology of indexing defects by keywords

- Prior work
- Research contributions
- Methodology of indexing defects by keywords
- Applications of indexing defects by keywords

- Prior work
- Research contributions
- Methodology of indexing defects by keywords
- Applications of indexing defects by keywords
- Evaluation

- Prior work
- Research contributions
- Methodology of indexing defects by keywords
- Applications of indexing defects by keywords
- Evaluation
- Conclusions & future work



## Research with different emphases on defect data

Research	Method <sup>a</sup>	Defect type <sup>b</sup>	R <sup>c</sup>	C <sup>d</sup>
Chillarege92	M	function, interface, etc.	Y	N
Chou01	A	null, range, float, lock, etc.	N	Y
Nagappan05	M	N	N	Y
Tan07	A	copy-and-paste	N	Y
Mantyla09	M	functional, evolvable	Y	N
Yin11	A	fixes	N	Y

<sup>a</sup>Methods of defect finding: manual reviews or using automatic tools

<sup>b</sup>Defects are classified by types in research

<sup>c</sup>Code reviews are used in research

<sup>d</sup>Code commits are used in research

## Research with different emphases on defect data

Research	Method <sup>a</sup>	Defect type <sup>b</sup>	R <sup>c</sup>	C <sup>d</sup>
Chillarege92	M	function, interface, etc.	Y	N
Chou01	A	null, range, float, lock, etc.	N	Y
Nagappan05	M	N	N	Y
Tan07	A	copy-and-paste	N	Y
Mantyla09	M	functional, evolvable	Y	N
Yin11	A	fixes	N	Y

<sup>a</sup>Methods of defect finding: manual reviews or using automatic tools

<sup>b</sup>Defects are classified by types in research

<sup>c</sup>Code reviews are used in research

<sup>d</sup>Code commits are used in research

**But these methods cannot solve our problem!**

- Methodology of indexing defects by keywords

- Methodology of indexing defects by keywords
- Applications of indexing defects by keywords
  - What are the defect types?

- Methodology of indexing defects by keywords
- Applications of indexing defects by keywords
  - What are the defect types?
  - Do defects cluster in files?

- Methodology of indexing defects by keywords
- Applications of indexing defects by keywords
  - What are the defect types?
  - Do defects cluster in files?
  - How to predict the hot files of defects?

# Methodology of indexing defects by keywords

Keywords are found in defect reports ...

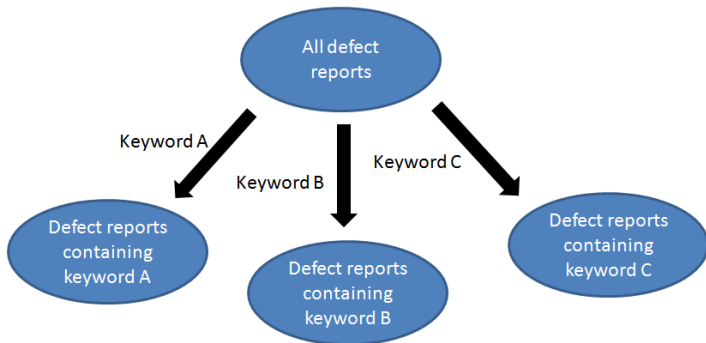
```
<item>
  <title> [AJE-105] Exception in Defect Scorecard </title>
  <summary> Exception in Defect Scorecard </summary>
  <description> ... Currently left = class java.lang.Long, right = class java.lang.String. ...</description>
  <key id=" 215630 "> AJE-105 </key>
  <status id="6"> Closed </status>
  <resolution id="1"> Fixed </resolution>
  <assignee username="eXX">XXX</assignee>
  <reporter username="eYY">YYY</reporter>
  <comments>
    <comment> XXXX</comment>
    <comment> YYYY </comment>
  </comments>
</item>
```



**Keyword String!**

# Methodology of indexing defects by keywords

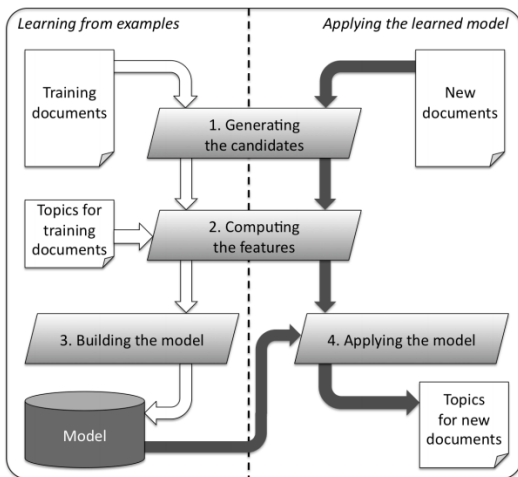
Defect reports can be indexed by keywords ...





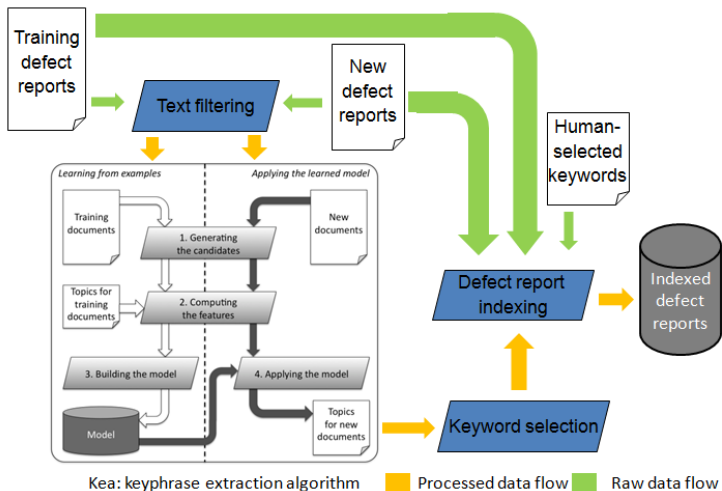
# Methodology of indexing defects by keywords

## Kea, a Keyphrase Extraction Algorithm [Medelyan09]



# Methodology of indexing defects by keywords

## Integration of Kea into the framework ...



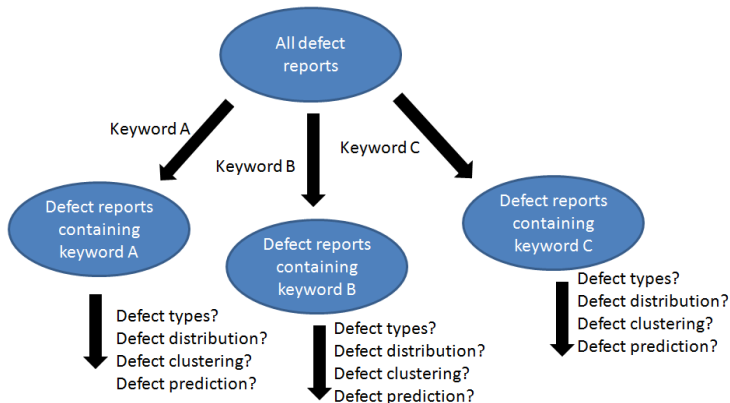
# Methodology of indexing defects by keywords

An example of defects indexed by the keyword **memory**

Defects indexed by keywords	
Keyword	Defect summary
<b>memory</b>	Bad response to HVAC Counter
	Rheem Furnace Model and Serial Number are cleared ...
	ECM Blower problem after User menu update
	no D4 alarm when memory card is invalid and data ...
	d1 alarm active even if memory card with correct ...
	Overflow of array index
	ESD Failure of Memory Card Connector
	RAM usage
	Bad SSD display when is fault clearing

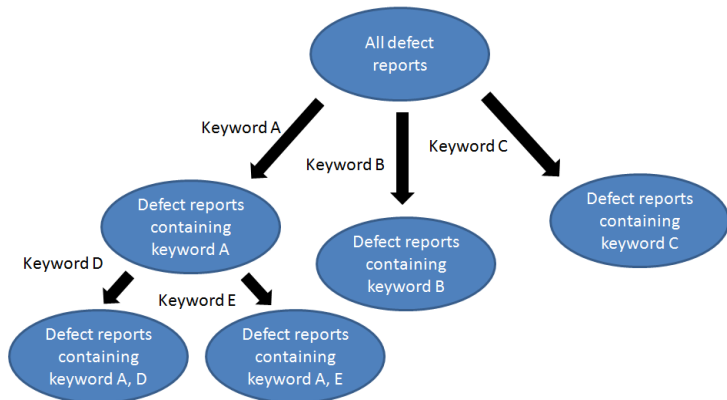
# Applications of indexing defects by keywords

After defect indexing ...



# What are the defect types?

Defect types are represented by (multiple) keywords ...



# What are the defect types?

---

**Algorithm 5** Inferring defect types by keywords

**Input:** a list of defect reports and a list of defect patterns

**Output:** a map of defect types with the associated list of defect reports

---

**Require:**  $x$  is a list of filtered defect reports,  $p$  is a list of defect patterns

```
1: function INFERRINGDEFECTTYPESBYKEYWORDS( $x, p$ )
2:    $z \leftarrow$  new Map for storing indexed defects
3:   for  $x_i : x$  do
4:      $maxFreq \leftarrow 0$ 
5:      $leastFirstPos \leftarrow -1$ 
6:      $defectType \leftarrow null$ 
7:     for  $p_j : p$  do
8:       if  $ContainAllKeywords(x_i, p_j)$  then
9:          $fs \leftarrow$  ComputeGeneralizedFrequencyScore( $p_j, x_i$ )
10:         $fps \leftarrow$  ComputeGeneralizedFirstPositionScore( $p_j, x_i$ )
11:        if  $maxFreq \leq fs$  then ▷ determine the dominant pattern
12:           $maxFreq \leftarrow fs$ 
13:           $defectType \leftarrow p_j$ 
14:           $leastFirstPos \leftarrow fps$ 
15:        else if  $fs = maxFreq$  then
16:          if  $fps \leq leastFirstPos$  then
17:             $defectType \leftarrow p_j$ 
18:             $leastFirstPos \leftarrow fps$ 
19:        if  $defectType \neq null$  then
20:           $z \leftarrow z.put(defectType, x_i)$ 
21:   return  $z$  ▷  $z$  is a map of defect types with the associated list of defect reports
```

---

$$\textit{Fraction80Metric} = \frac{\text{number of files containing 80\% of defects}}{\text{total number of files}}$$

number of files are counted from the files with the most defects.

# How to predict the hot files of defects?

$$BugPredMetric = \sum_{i=1}^N \frac{1}{1 + e^{slope \cdot t_i}} \text{ [Google11]}$$

where

*slope* is a configurable parameter with default value 12.

*N* is the number of commits of a file.

*t<sub>i</sub>* is the normalized timestamp of the *i*th commit,

with the oldest timestamp = 1 and now = 0.

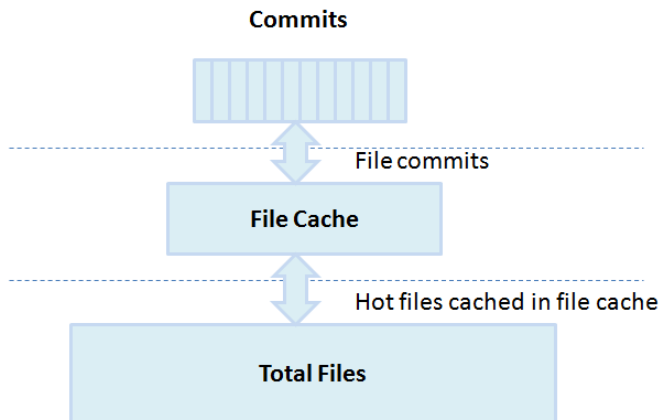
$$\text{defect hit rate} = \text{hit} / (\text{hit} + \text{miss})$$

*hit/miss* is the number of times a committed file is (/not) in the file cache



# How to predict the hot files of defects?

## File cache simulation



# How to predict the hot files of defects?

---

**Algorithm 10** Defect hot spots prediction in files

**Input:** a list of file commits, a file cache size, a slope parameter

**Output:** computed defect hit rate of the file cache replacement policy using a known bug prediction metric, a list of hot files in the file cache

---

**Require:**  $x$  is a list of file commits,  $c$  is the size of file cache,  $s$  is the slope parameter

```
1: function DEFECTHOTSPOTSPREDICTIONINSOURCEFILES( $x, c, s$ )
2:    $hit \leftarrow 0$ 
3:    $miss \leftarrow 0$ 
4:    $defectHitRate \leftarrow 0.0$ 
5:    $cache \leftarrow$  New Set for storing files
6:    $x \leftarrow$  SortCommitsInChronologicalOrder( $x$ )
7:   for  $x_i : x$  do
8:      $f \leftarrow$  ListOfCommittedFiles( $x_i$ )
9:     for  $f_j : f$  do
10:      if  $cache$  contains  $f_j$  then
11:         $hit \leftarrow hit + 1$ 
12:      else
13:         $miss \leftarrow miss + 1$ 
14:       $h \leftarrow$  UpdateCommittedFiles( $f_j$ )
15:       $r \leftarrow$  RankFilesByBugPredScore( $h, s$ )  $\triangleright$  Compute BugPredMetric for
        each file, and sort the files in descending order
16:       $cache \leftarrow$  UpdateFileCacheByScore( $r, c$ )  $\triangleright$  Keep the top  $c$  files in cache
17:    $defectHitRate \leftarrow \frac{hit}{hit+miss}$ 
18:   Output( $cache$ )  $\triangleright$  output the files in the file cache
19:   return  $hitRate$ 
```

---

Project	Description	PL	Defects	Commits
Project X	a JIRA server extension	Java	423	3235
Project Y	Gas ignition controls	C	300	2977

# Project X

Project	Description	PL	Defects	Commits
Project X	a JIRA server extension	Java	423	3235
Project Y	Gas ignition controls	C	300	2977

Table: Fraction of defects by defect types for project X (total 423 defects)

Defect type (keyword)	Fraction of defects
<b>layout</b>	1.42%
<b>string</b>	1.89%
<b>javascript</b>	1.89%
<b>reproduce</b>	4.49%
<b>message</b>	5.91%
<b>configuration</b>	6.15%
<b>display</b>	9.93%
<b>panel</b>	11.58%
<b>dialog</b>	13.71%
<b>report</b>	15.13%
Unindexed defects	27.19%

Table: *Fraction80Metric* by defect types for project X

Defect type (keyword)	Number of files	<i>Fraction80Metric</i>
All defects	1375	42.76%
<b>layout</b>	17	35.29%
<b>string</b>	9	55.56%
<b>javascript</b>	15	20.00%
<b>reproduce</b>	17	29.41%
<b>message</b>	35	28.57%
<b>configuration</b>	16	31.25%
<b>display</b>	65	24.62%
<b>panel</b>	42	30.95%
<b>dialog</b>	51	41.18%
<b>report</b>	56	48.21%

Table: Top-10 hot files for project X

atlassian-plugin.xml
ContourProjectManagerImpl.java
component-wise-issue-chart.js
changeHistory.js
ComponentWiselssueReportResource.java
defect-resolution-time-chart.js
DefectResolutionTimeReportResource.java
chart.js
ComponentWiselssueReport.java
ChangeHistoryTableImpl.java

**Table:** Defect hit rate by defect types for project X

file cache size = 20, and slope = 12		
Defect type (keyword)	Files	Defect hit rate
All defects	1375	30.66%
<b>layout</b>	17	100.00%
<b>string</b>	9	100.00%
<b>javascript</b>	15	100.00%
<b>reproduce</b>	17	100.00%
<b>message</b>	35	95.61%
<b>configuration</b>	16	100.00%
<b>display</b>	65	90.60%
<b>panel</b>	42	90.08%
<b>dialog</b>	51	86.23%
<b>report</b>	56	91.10%



# Project X

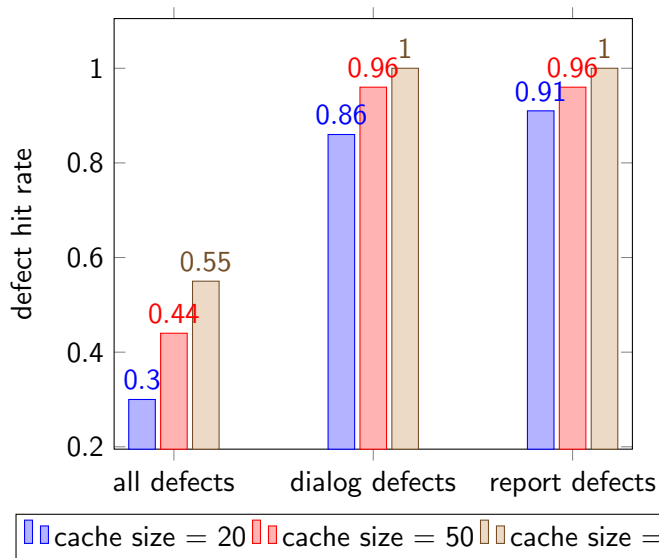


Figure: Defect hit rate with different file cache size and slope = 12 for project X

# Project Y

Project	Description	PL	Defects	Commits
Project X	a JIRA server extension	Java	423	3235
Project Y	Gas ignition controls	C	300	2977

**Table:** Fraction of defects by defect types for project Y (total 300 defects)

Defect type (keyword)	Fraction of defects
<b>print</b>	0.33%
<b>string</b>	1.00%
<b>memory</b>	1.67%
<b>fan</b>	2.67%
<b>cool</b>	3.00%
<b>power</b>	5.67%
<b>lock</b>	6.33%
<b>circulator</b>	7.00%
<b>alarm</b>	11.67%
<b>heat</b>	22.33%
Unindexed defects	38.33%

Table: *Fraction80Metric* by defect types for project Y

Defect type (keyword)	Number of files	<i>Fraction80Metric</i>
All defects	2736	7.27%
<b>string</b>	25	56.00%
<b>fan</b>	30	73.33%
<b>cool</b>	50	64.00%
<b>power</b>	18	33.33%
<b>lock</b>	53	52.83%
<b>circulator</b>	76	44.74%
<b>alarm</b>	64	34.38%
<b>heat</b>	104	55.77%

Table: Top-10 hot files for project Y

application.c
ign_interface.c
osdata.h
IgnSupport.c
RN_ObjectProc.c
LED_UserInterface.c
error.c
appscheduler.c
ModelDataCheck.c
RN_FaultsSetup.h

Table: Defect hit rate by defect types for project Y

file cache size = 20, and slope = 12		
Defect type (keyword)	Files	Defect hit rate
All defects	1375	22.37%
<b>string</b>	25	88.25%
<b>memory</b>	3	100.00%
<b>fan</b>	30	86.59%
<b>cool</b>	50	90.04%
<b>power</b>	18	100.00%
<b>lock</b>	53	85.79%
<b>circulator</b>	76	81.32%
<b>alarm</b>	64	87.78%
<b>heat</b>	104	69.71%

# Project Y

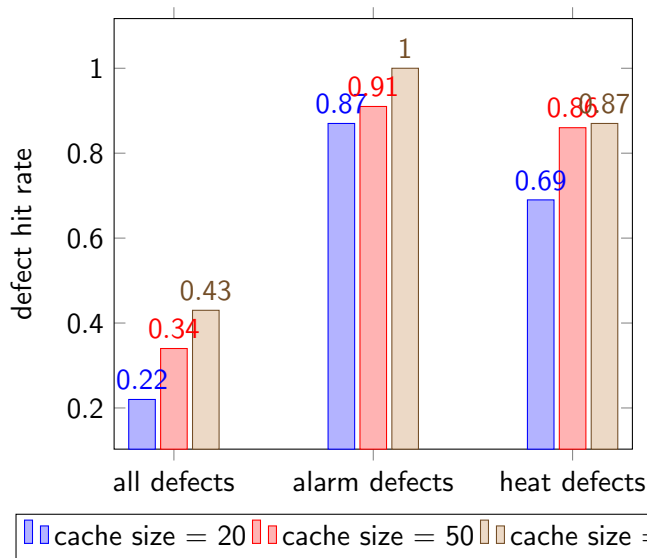


Figure: Defect hit rate with different file cache size and slope = 12 for project Y

What are the defect types? Infer them by keywords.

Project	Defect types (keywords)
Project X	<b>layout, string, javascript, reproduce, message, configuration, display, panel, dialog, report,</b> <10% non-functional defects indexed, 27% defects unindexed.
Project Y	<b>print, string, memory, fan, cool, power, lock, circulator, alarm, heat,</b> <10% non-functional defects indexed, 38% defects unindexed.



Do defects cluster in files? YES!

Project	Defect clustering
Project X	For all defect types 20%-50% of files contained more than 80% of defects.
Project Y	For most defect types 20%-50% of files contained more than 80% of defects.

What are the top files to review?

**Table:** Top files to review for projects X and Y

Project	Hot files of defects
Project X	atlassian-plugin.xml
	ContourProjectManagerImpl.java
	component-wise-issue-chart.js
	changeHistory.js
	ComponentWiselssueReportResource.java
Project Y	application.c
	ign_interface.c
	osdata.h
	IgnSupport.c
	RN_ObjectProc.c

Methodology of indexing defects by keywords

- New approach to defect classification

## Methodology of indexing defects by keywords

- New approach to defect classification
- Scales to large dataset

## Methodology of indexing defects by keywords

- New approach to defect classification
- Scales to large dataset
- Applies to diversified projects (PLs, applications)

## Methodology of indexing defects by keywords

- New approach to defect classification
- Scales to large dataset
- Applies to diversified projects (PLs, applications)

## Applications of indexing defects by keywords

- Effective inference of a wide range of defect types

## Methodology of indexing defects by keywords

- New approach to defect classification
- Scales to large dataset
- Applies to diversified projects (PLs, applications)

## Applications of indexing defects by keywords

- Effective inference of a wide range of defect types
- Type-based clustering analysis of defects

## Methodology of indexing defects by keywords

- New approach to defect classification
- Scales to large dataset
- Applies to diversified projects (PLs, applications)

## Applications of indexing defects by keywords

- Effective inference of a wide range of defect types
- Type-based clustering analysis of defects
- Type-based prediction of the hot files of defects



## Advanced indexing techniques

- Dictionary-based, semantics-aware indexing?

## Advanced indexing techniques

- Dictionary-based, semantics-aware indexing?
- How to integrate user tagging?

## Advanced indexing techniques

- Dictionary-based, semantics-aware indexing?
- How to integrate user tagging?

## Towards dynamic defect analysis

- Would the dominant defect type vary through time?

## Advanced indexing techniques

- Dictionary-based, semantics-aware indexing?
- How to integrate user tagging?

## Towards dynamic defect analysis

- Would the dominant defect type vary through time?
- What would be the implications of defect evolution to testing?

## Advanced indexing techniques

- Dictionary-based, semantics-aware indexing?
- How to integrate user tagging?

## Towards dynamic defect analysis

- Would the dominant defect type vary through time?
- What would be the implications of defect evolution to testing?

## Integrating defect data online: JIRA plugin

- How to optimize the performance of the tool?

## Advanced indexing techniques

- Dictionary-based, semantics-aware indexing?
- How to integrate user tagging?

## Towards dynamic defect analysis

- Would the dominant defect type vary through time?
- What would be the implications of defect evolution to testing?

## Integrating defect data online: JIRA plugin

- How to optimize the performance of the tool?
- How to collect more complete defect dataset?

# Thank you, questions please.

{tao.lee}@epfl.ch

**Honeywell**

