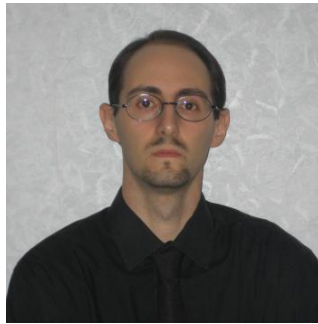# Improving IR-based Traceability Recovery Using Smoothing Filters

Andrea
De Lucia

Massimiliano
Di Penta

Rocco
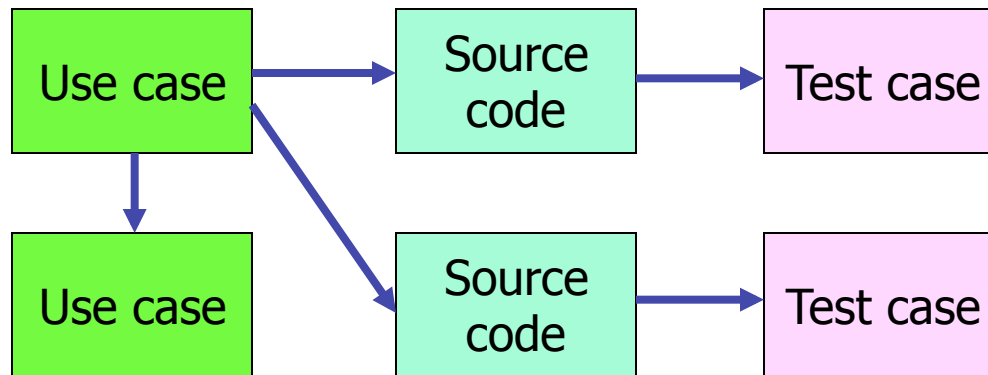Oliveto

Annibale
Panichella

Sebastiano
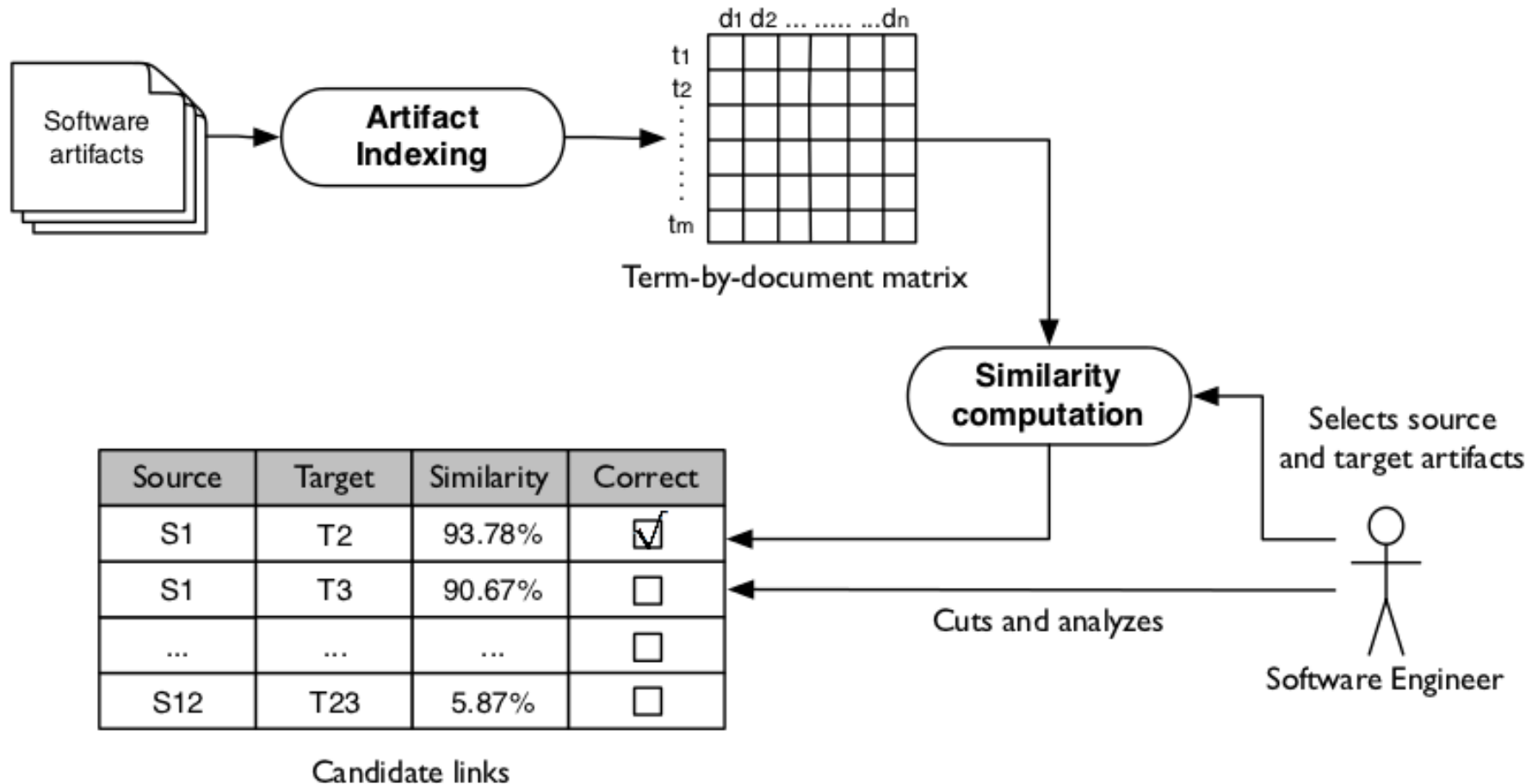Panichella

# Software traceability

"*The degree to which a relationship can be established between two products of a software development process*"
[IEEE Glossary for Software Terminology]



- Important for:
  - program comprehension
  - requirement tracing
  - impact analysis
  - software reuse
  - ...

Up-to-date traceability links rarely exist → need to recover them

# IR-based traceability recovery



Antoniol et al., 2002 (VSM+Probabilistic model)
Marcus and Maletic, 2003 (LSI)

# Traditional IR vs. IR applied to Software Engineering

## Traditional IR

- Deals with heterogeneous documents for what concerns:
  - Linguistic choices
  - Syntax
  - Semantics
- We just live with that differences

## IR applied to SE

- We have sets of homogeneous documents for what concerns
  - Syntax, linguistic choices
- Examples:
  - Use cases, test documents, design documents follow a common template and contain recurrent words

# Problem

- Different kinds of software artifacts require specific preprocessing

```
Test case    Change the date for a visit:

   C51       Version: 0 02 000

Use case     Satisfies the request to modify a visit
             for a patient

UcModVis

Priority     High

....

Test description

Input        Select a visit:

             26/09/2003 11:00  First visit

             Change: 03/10/2003 11:00

Oracle       Invalid sequence: The system does not allow
             to change a booking

Coverage     Valid classes: CE1  CE8  CE14  CE19  CE21

             Invalid classes: None
```

# Problem

- Different kinds of software artifacts require specific preprocessing

```
Test case    Change the date for a visit:

  C51        Version: 0 02 000

Use case     Satisfies the request to modify a visit
             for a patient

UcModVis

Priority     High

....

Test description

Input        Select a visit:

             26/09/2003 11:00  First visit

             Change: 03/10/2003 11:00

Oracle       Invalid sequence: The system does not allow
             to change a booking

Coverage     Valid classes: CE1  CE8  CE14  CE19  CE21

             Invalid classes: None
```
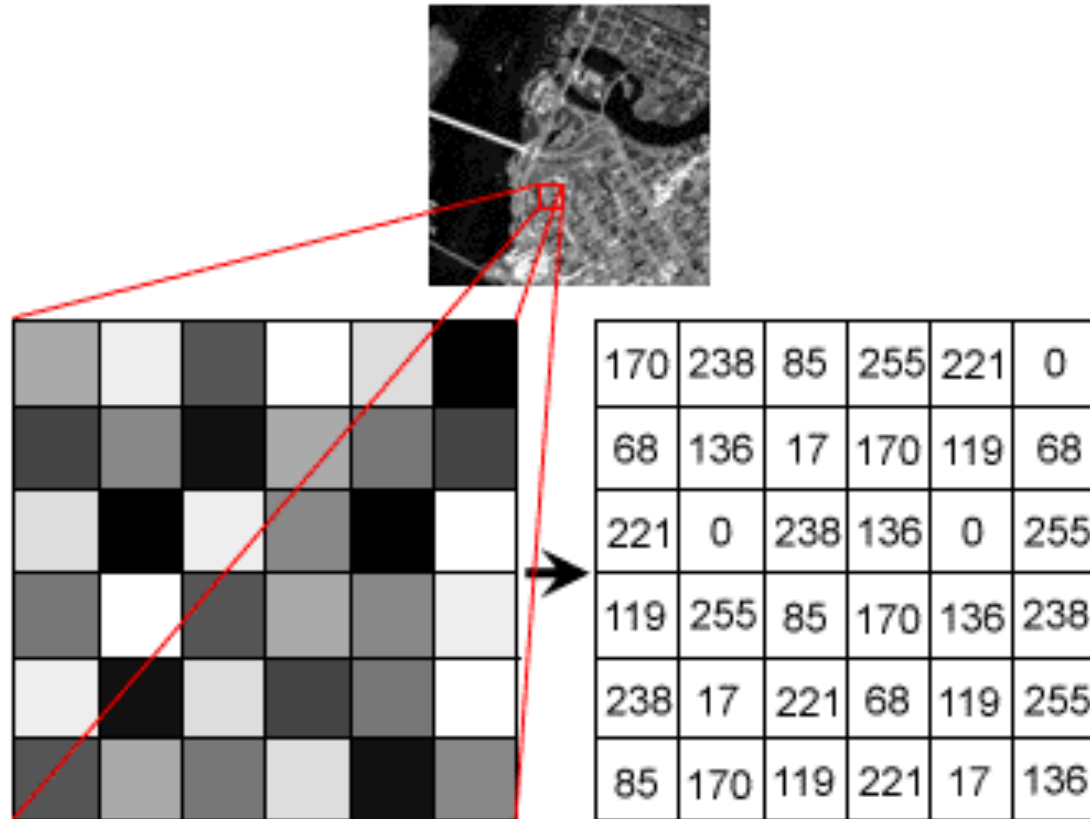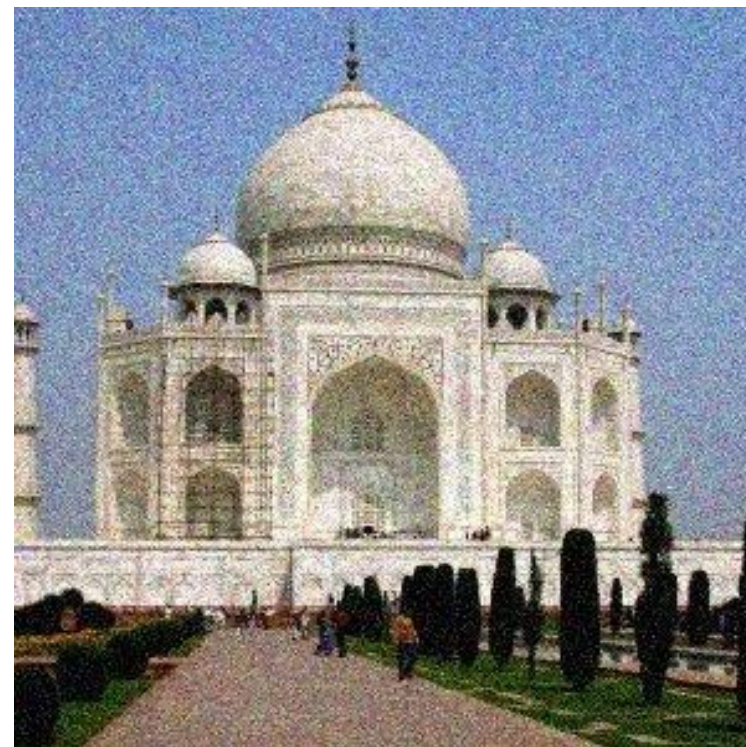
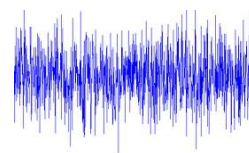Artifact-specific words do not bring useful information
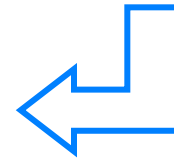
# Noisy images

Pixels with peaks of low color intensity

Pixels with peaks of high color intensity

**Noise**

# Reducing noise using smoothing filters



Mean filter

$$g(x,y) = \frac{1}{M} \sum_{f(n,m) \in S} f(n,m)$$

# Image vs. traceability noise

## Image noise:

- Pixels with high or low color intensity
- Pixels are position dependent

## Traceability noise:

- Terms and linguistic patterns occurring in many artifacts of a given category
    - Use cases, test cases..
- Artifacts (columns) are position independent

**Source Documents**

**Target Documents**

$s_1$    $s_2$    $s_3$    ...    $s_k$     $t_1$    $t_2$    $t_3$    ...    $t_z$

$word_1$   $v_{1,1}$   $v_{1,2}$   $v_{1,3}$   ...   $v_{1,k}$

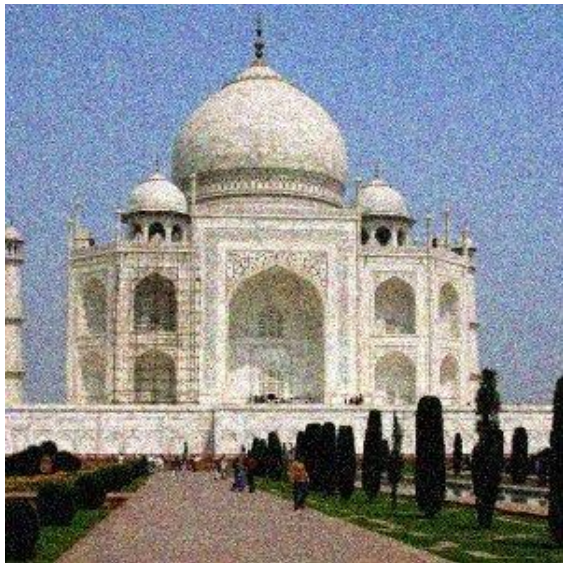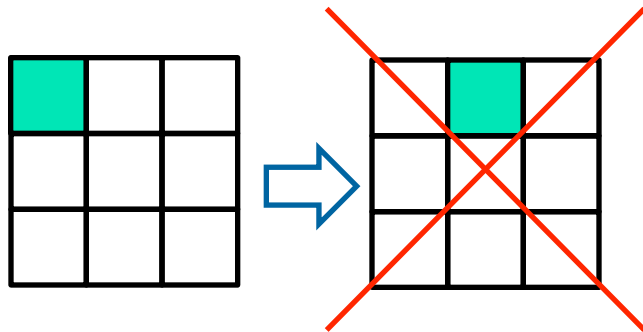$word_2$   $v_{2,1}$   $v_{2,2}$   $v_{2,3}$   ...   $v_{2,k}$

$word_n$   $v_{n,1}$   ... $v_{n,2}$   ... $v_{n,3}$     $v_{n,k}$

$v_{1,1}$   $v_{1,2}$   $v_{1,3}$   ...   $v_{1,z}$

$v_{2,1}$   $v_{2,2}$   $v_{2,3}$   ...   $v_{2,z}$

$v_{n,1}$   ... $v_{n,2}$   ... $v_{n,3}$     $v_{n,z}$

Linguistic information strictly belonging to source documents

Linguistic information strictly belonging to target documents

**Common Information**
for source documents

**Common Information**
For target documents

# Representing the noise

**Source Documents**                    **Target Documents**



**Mean source vector** $S=$

$$\begin{bmatrix} \frac{1}{k}\sum_{j=1}^{k} v_{1,j} \\ \frac{1}{k}\sum_{j=1}^{k} v_{2,j} \\ \vdots \\ \frac{1}{k}\sum_{j=1}^{k} v_{n,j} \end{bmatrix}$$

**Mean target vector** $T=$

$$\begin{bmatrix} \frac{1}{z}\sum_{j=k+1}^{m} v_{1,j} \\ \frac{1}{z}\sum_{j=k+1}^{m} v_{2,j} \\ \vdots \\ \frac{1}{z}\sum_{j=k+1}^{m} v_{n,j} \end{bmatrix}$$

Common Information
for source documents

Common Information
for target documents

The Mean vectors are like the continuous component of a signal…

# Representing the noise

**Source Documents**

**Target Documents**



$s_1$ $\quad$ $s_2$ $\quad$ $s_3$ $\quad$ ... $\quad$ $s_k$

$$\begin{bmatrix} v_{1,1} & v_{1,2} & v_{1,3} & & v_{1,k} \\ & & & ... & \\ v_{2,1} & v_{2,2} & v_{2,3} & ... & v_{2,k} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ & & & ... & \\ v_{n,1} & ... & v_{n,2} & ... & v_{n,3} \end{bmatrix}$$

$word_1$ $\quad$ $word_2$ $\quad$ $\vdots$ $\quad$ $word_n$

$t_1$ $\quad$ $t_2$ $\quad$ $t_3$ $\quad$ ... $\quad$ $t_z$

$$\begin{bmatrix} v_{1,1} & v_{1,2} & v_{1,3} & & v_{1,z} \\ & & & ... & \\ v_{2,1} & v_{2,2} & v_{2,3} & ... & v_{2,z} \\ \vdots & \ddots & \vdots & \ddots & \vdots \\ & & & ... & \\ v_{n,1} & ... & v_{n,2} & ... & v_{n,3} \end{bmatrix}$$

$-$

$-$

S

T

**(mean source vector)**

**(mean target vector)**

Filtered
Source Set

Filtered
Target Set

# Empirical Study

- **Goal:** analyze the effect of smoothing filter
- **Purpose:** investigating how the filter affects traceability recovery
- **Quality focus:** traceability recovery performance
- **Perspective:**
  - **Researchers:** evaluating the novel technique
- **Context:** artifacts from two systems
  - EasyClinic and Pine

# Context

| | EasyClinic | Pine |
|---|---|---|
| Description | Medical doctor office management | Text-based email client |
| Language | Java | C |
| Files/Classes | 37 | 31 |
| KLOC | 20 | 130 |
| Documents | 113 | 100 |
| Language | Italian | English |
| Artifacts | Use cases<br>Interaction diagrams<br>Source code<br>Test cases | Requirements<br>Use cases |

# Research Questions and Factors

- RQ1: Does the smoothing filter improve the recovery performances of traceability recovery?
- RQ2: How effective is the smoothing filter in filtering out non-relevant words, as compared to stop word removal?

- Factors:
  - Use of filter: YES, NO
  - Technique: VSM, LSI
  - Artifact: Req., UC, Int. Diagrams, Code, TC
  - System: Easyclinic, Pine

- **Performances evaluated by precision and recall:**

$$precision = \frac{|correct \cap retrieved|}{|retrieved|} \qquad recall = \frac{|correct \cap retrieved|}{|correct|}$$

- **We statistically compare the # of false positives of different methods for each correct link identified**
  - Wilcoxon Rank Sum test
  - Cliff's delta effect size

M1     M2

| M1 | M2 |
|----|----|
|    | 0  |
| 2  | 2  |
|    | 3  |

- We replace stop word filtering by one of the following treatments:
  1. Standard stop word removal
  2. Manually customized stop word removal
  3. Smoothing filter
  4. Standard stop word removal + filter
  5. Customized stop word removal + filter

- …and compare the performances

# Results

# EasyClinic: Use cases into source (VSM)



[-60, -74]% of false positives for recall<80%

# EasyClinic: Use cases into source (LSI)

[-60, -77]% of false positives for recall<80%

# EasyClinic: Test cases into source (LSI)



Test cases are:
- Short documents
- Limited vocabulary
- Mostly consistent with source code

# Pine: Use cases into requirements (LSI)



[-42, -62]% of false positives for recall<80%

# Statistical Comparison

| Data set | Traced Artifacts | VSM | | LSI | |
|---|---|---|---|---|---|
| | | p-value | Effect size | p-value | Effect size |
| EasyClinic | UC→Code | **<0.01** | 0.50 (**large**) | **<0.01** | 0.50 (**large**) |
| | Int. Diag. → Code | **<0.01** | 0.52 (**large**) | **<0.01** | 0.34 (**medium**) |
| | TC → Code | 1.00 | - (**negligible**) | 1.00 | - (**negligible**) |
| Pine | Req. → UC | **<0.01** | 0.58 (**large**) | **<0.01** | 0.58 (**large**) |

# RQ2 – Summary of results

| Comparison | | EasyClinic | | | Pine |
|---|---|---|---|---|---|
| | | UC➜CC | ID➜CC | TC➜CC | HLR➜ UC |
| Smoothing filter | Standard list | YES (small) | YES (small) | NO (large) | YES (large) |
| Smoothing filter | Cust. list | YES (small) | YES (small) | NO (large) | YES (large) |
| Standard list+ Smoothing filter | Cust. list | YES (large) | YES (large) | NO (medium) | YES (large) |
| Standard list+ Smoothing filter | Cist list + Smoothing filter | NO (small) | - | YES (medium) | YES (small) |

# Link precision improvement



Login Patient vs. Person
Poor vocabulary overlap (10%)

# Threats to validity

- **Construct validity**
  - Mainly related to our oracle
  - Provided by developers and for EasyClinic also peer-reviewed
- **Internal validity**
  - Improvements could be due to other reasons…
  - However, we compared different techniques (VSM, LSI)
  - The approach works well regardless of stop word removal, stemming and use of tf-idf
- **Conclusion validity**
  - Conclusions based on proper (non-parametric) statistics
- **External validity**
  - We considered systems with different characteristics and artifacts
  - … but further studies are desirable

# Conclusions

## Representing the noise

**Source Documents**

**Target Documents**

| | $s_1$ | $s_2$ | $s_3$ | L | $s_k$ |
|---|---|---|---|---|---|
| $word_1$ | $v_{1,1}$ | $v_{1,2}$ | $v_{1,3}$ | L | $v_{1,k}$ |
| $word_2$ | $v_{2,1}$ | $v_{2,2}$ | $v_{2,3}$ | L | $v_{2,k}$ |
| M | M O | M O | M O | | M |
| $word_n$ | $v_{n,1}$ | L $v_{n,2}$ | L $v_{n,3}$ | L | $v_{n,k}$ |

| | $t_1$ | $t_2$ | $t_3$ | L | $t_z$ |
|---|---|---|---|---|---|
| | $v_{1,1}$ | $v_{1,2}$ | $v_{1,3}$ | L | $v_{1,z}$ |
| | $v_{2,1}$ | $v_{2,2}$ | $v_{2,3}$ | L | $v_{2,z}$ |
| | M O | M O | M O | | M |
| | $v_{n,1}$ | L $v_{n,2}$ | L $v_{n,3}$ | L | $v_{n,z}$ |

− → S *(mean target vector)*

− → T *(mean target vector)*

*Filtered Source Set*

*Filtered Target Set*

## EasyClinic: Use cases into source (LSI)

[-60, -77]% of false positives for recall<80%

Filtered

Not Filtered

Precision / Recall

## EasyClinic: Test cases into source (LSI)

Test cases are:
- Short documents
- Limited vocabulary
- Mostly consistent with source code

Filtered

Not Filtered

Precision / Recall

## RQ2 – Summary of results

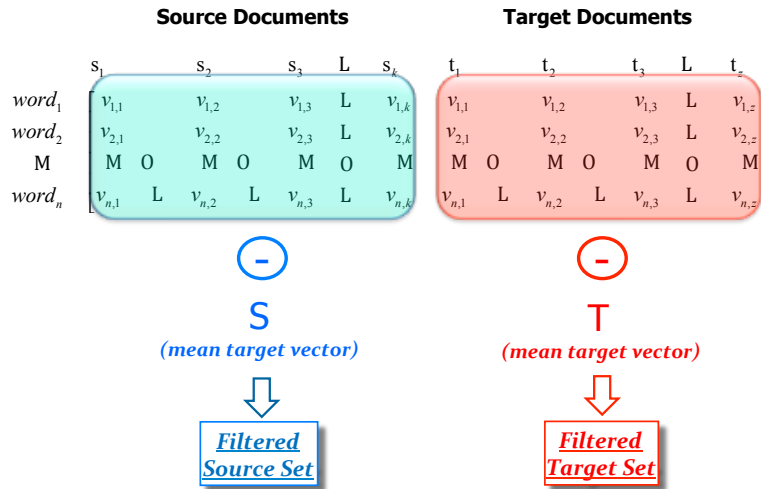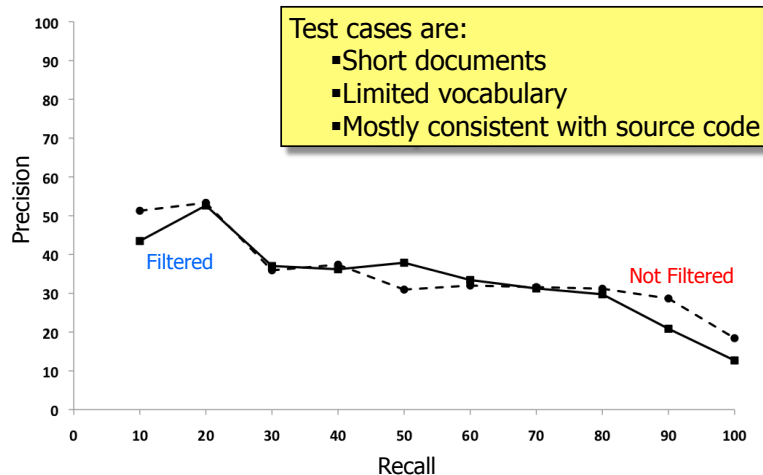| Comparison | | EasyClinic | | | Pine |
|---|---|---|---|---|---|
| | | UC➔CC | ID➔CC | TC➔CC | HLR➔UC |
| Smoothing filter | Standard list | YES (small) | YES (small) | NO (large) | YES (large) |
| Smoothing filter | Cust. list | YES (small) | YES (small) | NO (large) | YES (large) |
| Standard list+ Smoothing filter | Cust. list | YES (large) | YES (large) | NO (medium) | YES (large) |
| Standard list+ Smoothing filter | Cist list + Smoothing filter | NO (small) | - | YES (medium) | YES (small) |

# Work-in-progress

- Study replication
  - Different systems and artifacts
  - Use of relevance feedback
- More sophisticated smoothing technique
  - Non-linear filters
- Use in other applications of IR to software engineering
  - impact analysis
  - feature location