

# *When and How Using Structural Information to Improve IR-based Traceability Recovery*



Annibale Panichella<sup>1</sup>, Collin McMillan<sup>2</sup>, Evan Moritz<sup>3</sup>, Davide Palmieri<sup>4</sup>,  
Rocco Oliveto<sup>4</sup>, Denys Poshyvanyk<sup>3</sup>, Andrea De Lucia<sup>1</sup>

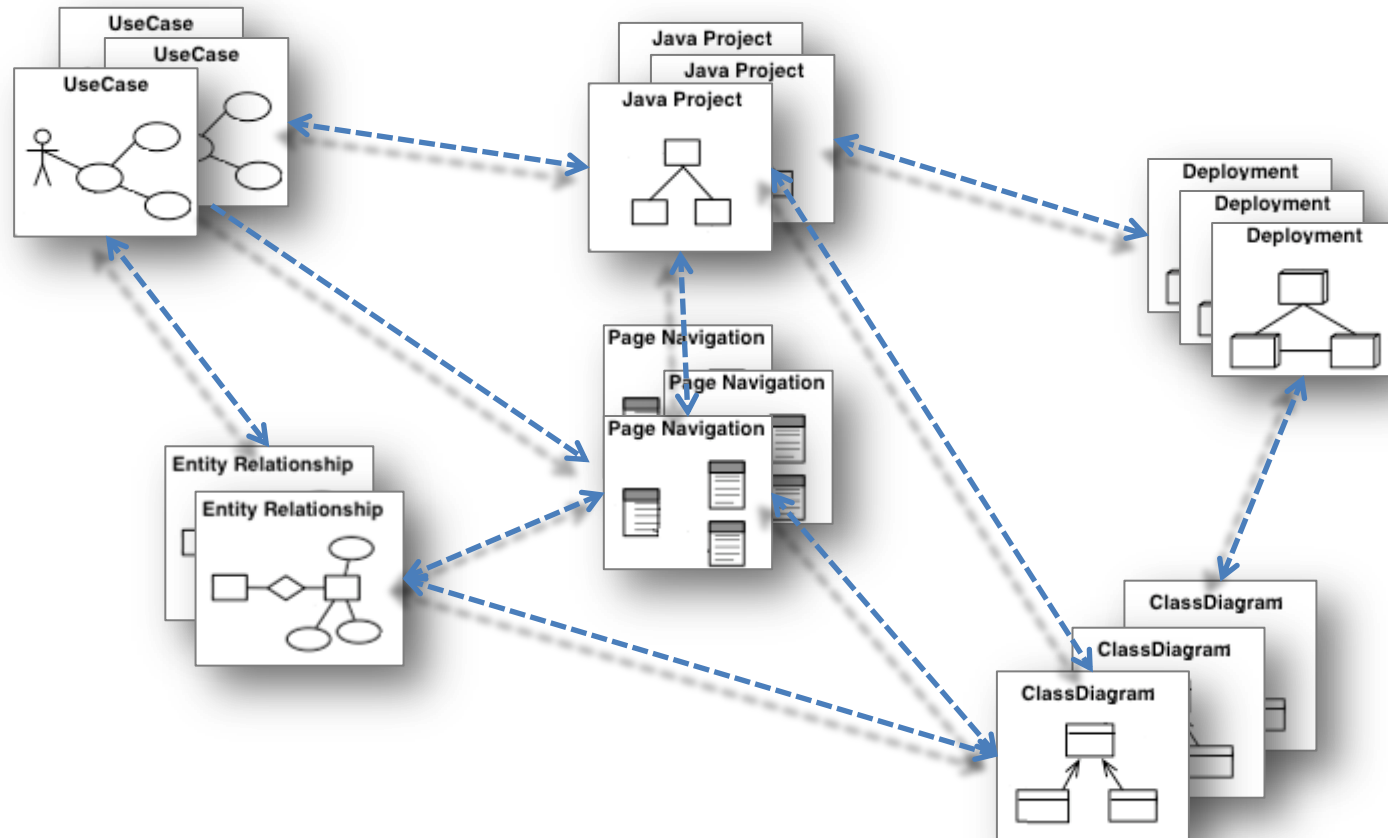
<sup>1</sup> Software Engineering Lab , University of Salerno, Italy

<sup>2</sup> University of Notre Dame, Notre Dame, USA

<sup>3</sup> The College of William and Mary, Williamsburg, USA

<sup>4</sup> University of Molise, Pesche (IS), Italy

*Traceability Recovery is the ability to describe and follow the artifacts life-cycle*

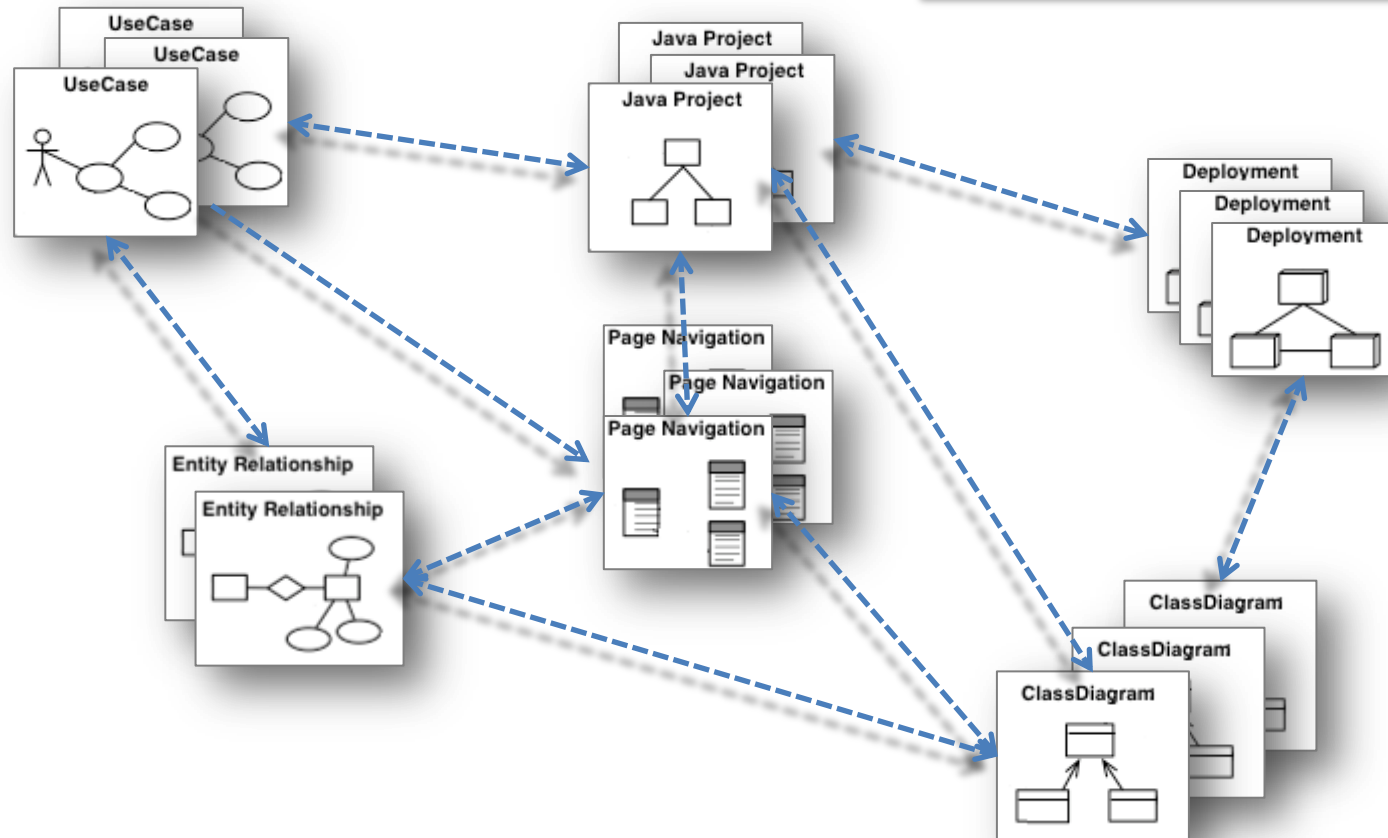


*Traceability Recovery is the ability to describe and follow the artifacts life-cycle*

*Impact Analysis*

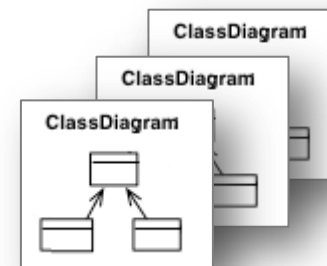
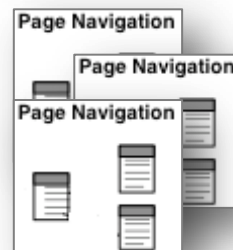
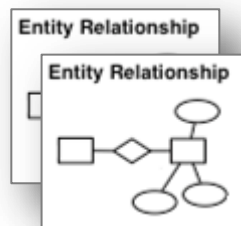
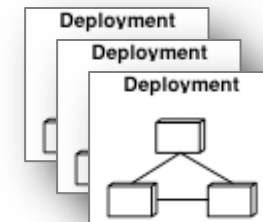
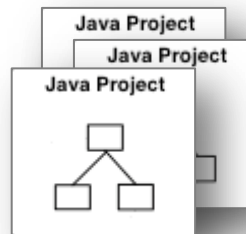
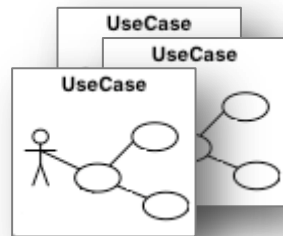
*Program Comprehension*

*Requirements tracing*



# Traceability in practice

*Traceability information is still not commonplace in software projects!*



# *Information Retrieval*



# IR-based Traceability Recovery

<b>Use Case</b>	<i>Insert Laboratory Data</i>
<b>Description</b>	<i>The user inserts the data of a specific laboratory</i>
<b>Events</b>	<i>1. The user opens the Laboratory GUI 2. The user inserts the laboratory data . . .</i>

## GUILaboratoryData.java

```
/* *This class implements the GUI for
managing laboratories data */

public class GUI Laboratory Data {
    private JFrame window;
    private JButton insert;
    ...

    public GUI Laboratory Data(){
        window = new JFrame();
        insert = new JButton();
        ...
    }

    ...
}
```

# IR-based Traceability Recovery

<i>Use Case</i>	<i>Insert Laboratory Data</i>
<i>Description</i>	<i>The user inserts the data of a specific laboratory</i>
<i>Events</i>	<i>1. The user opens the Laboratory GUI</i> <i>2. The user inserts the laboratory data</i> <i>.</i> <i>.</i> <i>.</i>

GUILaboratoryData.java

```
/* *This class implements the GUI for
managing laboratories data */

public class GUI Laboratory Data {
    private JFrame window;
    private JButton insert;
    ...

    public GUI Laboratory Data() {
        window = new JFrame();
        insert = new JButton();
        ...
    }

    ...
}
```

# IR-based Traceability Recovery

Use Case	Insert Laboratory Data
Description	The user inserts the data of a specific laboratory
Events	<ol style="list-style-type: none"><li>1. The user opens the Laboratory GUI</li><li>2. The user inserts the laboratory data</li><li>...</li><li>...</li><li>...</li></ol>

Similarity  
42%

GUILaboratoryData.java

```
/* *This class implements the GUI for
managing laboratories data */

public class GUI Laboratory Data {
    private JFrame window;
    private JButton insert;
    ...

    public GUI Laboratory Data() {
        window = new JFrame();
        insert = new JButton();
        ...
    }

    ...
}
```



# IR-based Traceability Recovery

Use Case	Insert Laboratory Data
Description	The user inserts the data of a specific laboratory
Events	<ol style="list-style-type: none"><li>1. The user opens the Laboratory GUI</li><li>2. The user inserts the laboratory data</li><li>...</li><li>...</li><li>...</li></ol>

True Link

GUILaboratoryData.java

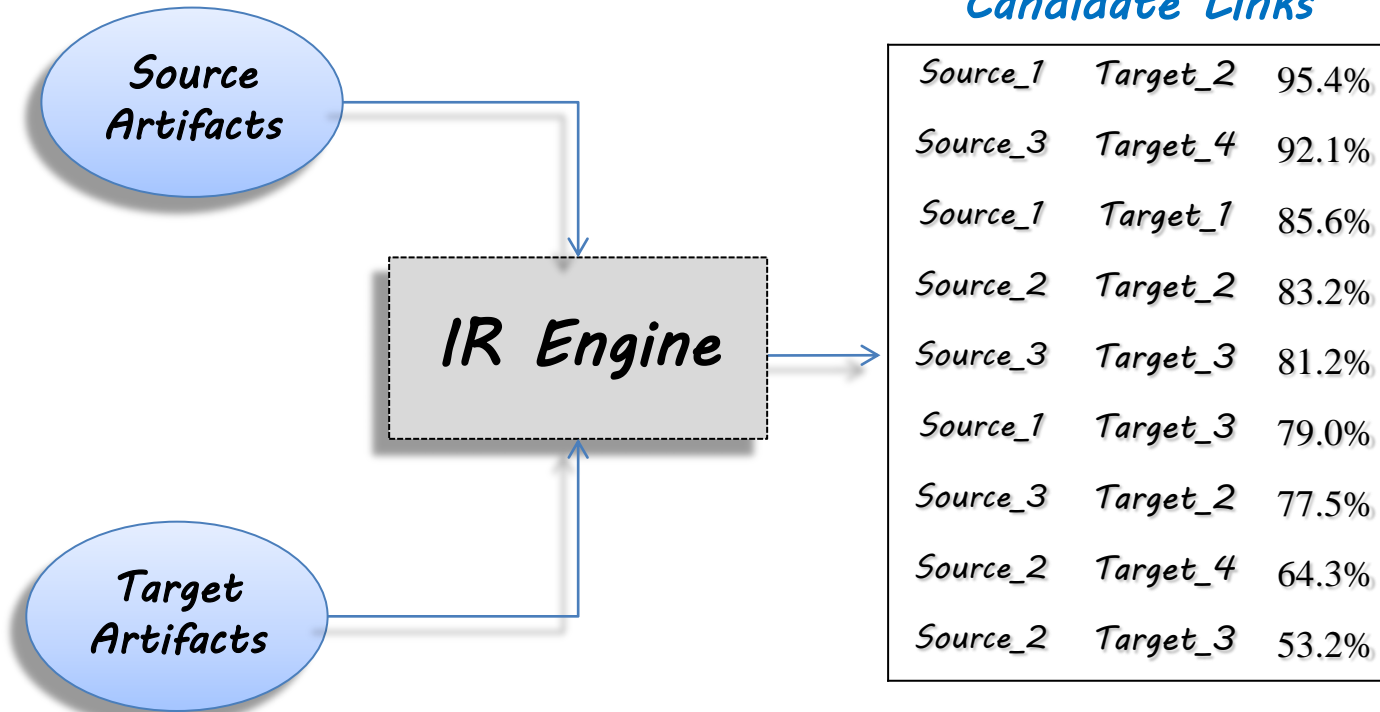
```
/* *This class implements the GUI for
managing laboratories data */

public class GUI Laboratory Data {
    private JFrame window;
    private JButton insert;
    ...

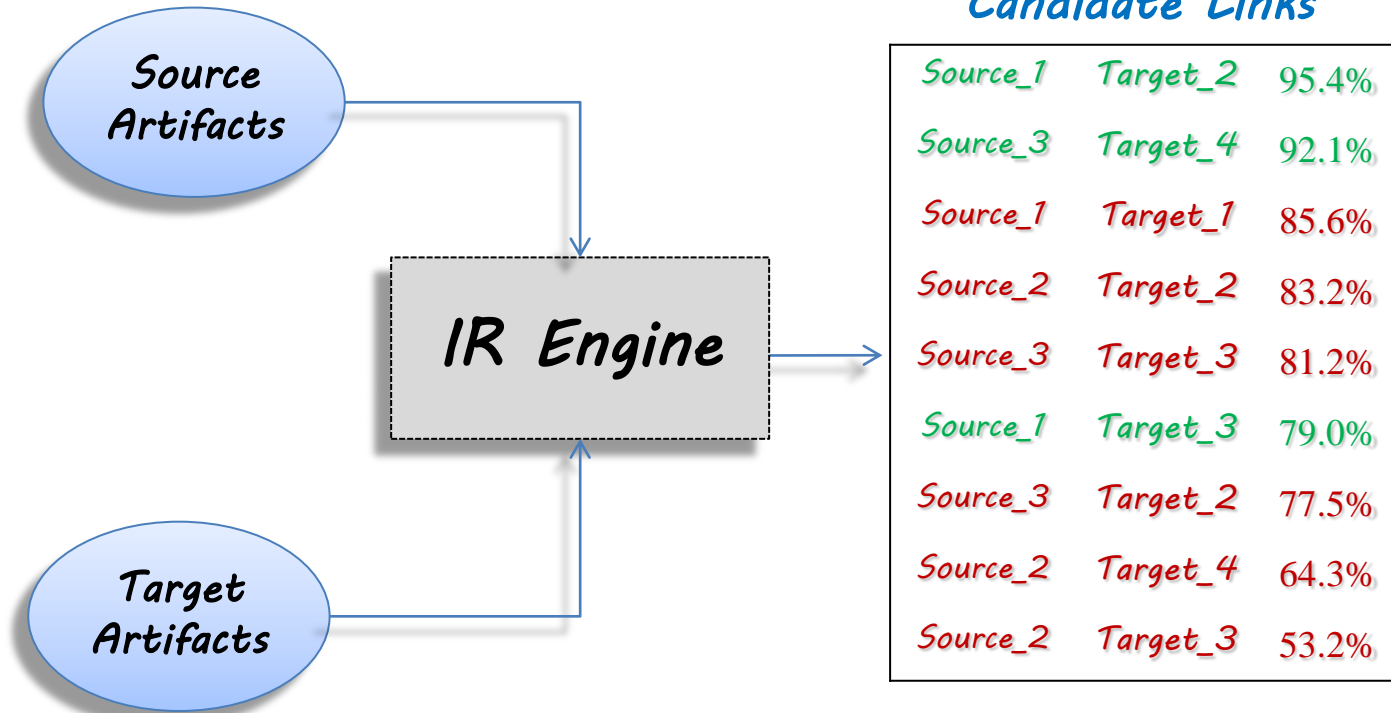
    public GUI Laboratory Data() {
        window = new JFrame();
        insert = new JButton();
        ...
    }

    ...
}
```

# IR-based Traceability Recovery



# IR-based Traceability Recovery



# Structural Information and Traceability

## Combining Textual and Structural Analysis of Software Artifacts for Traceability Link Recovery

Collin McMillan, Denys Poshyvanyk, Meghan Revelle

Department of Computer Science  
The College of William and Mary  
Williamsburg, VA 23185  
{cmc, denys, meghan}@cs.wm.edu

### Abstract

Existing methods for recovering traceability links among software documentation artifacts analyze textual similarities among these artifacts. It may be the case, however, that related documentation elements share little terminology or phrasing. This paper presents a technique for indirectly recovering these traceability links in requirements documentation by combining textual with structural information as we conjecture that related requirements share related source code elements. A preliminary case study indicates that our combined approach improves the precision and recall of recovering relevant links among documents as compared to stand-alone methods based solely on analyzing textual similarities.

### 1. Introduction

Existing methods of traceability link recovery rely on the analysis of textual information derived from software artifacts using Information Retrieval (IR) techniques [1, 2, 7, 11, 13, 18, 19, 21]. Feature and class descriptions in external documentation, for instance, are often likely to share some terms in common with their related source code artifacts due to naming conventions and formats that enforce identifier names, making IR techniques effective at discovering these documentation-to-source traceability links. However, there is no reason why two related documentation artifacts will use similar terminology since they might describe different tasks. IR methods are unlikely to recover useful traceability links in this situation. Therefore, an analysis of textual similarities among related source code artifacts (e.g., classes) may not yield any pertinent links.

Another approach to traceability link recovery is structural analysis, which alleviates this particular problem by examining control and data flow (i.e., coupling) among source code artifacts.

Taken together, textual and structural analysis

provide good methods to find source-to-source and source-to-documentation links, but not documentation-to-documentation links. Structural analysis closely lends itself to source code analysis, whereas using IR methods for traceability link recovery is a well-studied problem [1, 6-9, 11, 13, 15, 21]. In this work, we combine both textual and structural analysis methods on two types of software artifacts (i.e., source code and documentation) as well as on artifacts of one type to another, to create an indirect traceability link recovery scheme for software documentation.

This paper makes the following contributions:

- A method for indirectly recovering traceability links in requirements documentation using a combination of textual and structural information;
- A comparison of the proposed method with existing IR-based approaches for traceability link recovery;
- An investigation into the impacts of choosing various dimensionality reduction factors on the performance of both combined and stand-alone IR methods.

The paper is organized as follows. Section 2 presents background information on existing approaches for traceability link recovery. Section 3 outlines some key details of our approach. Section 4 presents an initial case study comparing the proposed combination with existing methods. Finally, related work and conclusions are outlined in Sections 6 and 7 respectively.

### 2. Background

In this section, we provide essential background information on the tools we use for textual and static analysis of software artifacts.

#### 2.1. Textual analysis

Latent Semantic Indexing (LSI) [10] is an IR technique that is used to determine textual similarities among words and documents in large passages of text. For example, LSI can determine how similar the text of

Combining Textual and  
Structural Analysis of Software  
Artifacts for Traceability Link  
Recovery - Collin McMillan, Denys  
Poshyvanyk, Meghan Revelle  
TEFSE 2009

# Structural Information and Traceability

<b>Use Case</b>	Insert Laboratory Data
<b>Description</b>	The user inserts the data of a specific laboratory
<b>Events</b>	<ol style="list-style-type: none"><li>1. The user opens the Laboratory GUI</li><li>2. The user inserts the laboratory data</li><li>.</li><li>.</li><li>.</li></ol>

linked

## GUILaboratoryData.java

```
/* *This class implements the GUI for
managing laboratories data */

public class GUILaboratoryData {
    private JFrame window;
    private JButton insert;
    ...

    public GUILaboratoryData() {
        window = new JFrame();
        insert = new JButton();
        ...
    }

    ...
}
```

Similarity =  
42%

# Structural Information and Traceability

<b>Use Case</b>	<i>Insert Laboratory Data</i>
<b>Description</b>	<i>The user inserts the data of a specific laboratory</i>
<b>Events</b>	<ol style="list-style-type: none"><li>1. The user opens the <b>Laboratory GUI</b></li><li>2. The user inserts the <b>laboratory data</b></li><li>.</li><li>.</li><li>.</li></ol>

linked

linked

```
public class Laboratory{  
    private String name;  
    private String position;  
    ...  
    public void setName(String pName){  
        this.name=pName;  
    }  
    ...  
}
```

Laboratory.java

GUILaboratoryData.java

```
/* *This class implements the GUI for  
managing laboratories data */  
  
public class GUILaboratoryData {  
    private JFrame window;  
    private JButton insert;  
    ...  
  
    public GUILaboratoryData(){  
        window = new JFrame("GUI Laboratory Data");  
        insert = new JButton("Insert");  
        ...  
    }  
    ...  
}
```

Similarity =  
42%

# Structural Information and Traceability

<b>Use Case</b>	<i>Insert Laboratory Data</i>
<b>Description</b>	<i>The user inserts the data of a specific laboratory</i>
<b>Events</b>	<ol style="list-style-type: none"><li>1. The user opens the <b>Laboratory GUI</b></li><li>2. The user inserts the <b>laboratory data</b></li><li>...</li><li>...</li><li>...</li></ol>

linked

linked

```
public class Laboratory{  
    private String name;  
    private String position;  
    ...  
    public void setName(String name) {  
        this.name=name;  
    }  
    ...  
}
```

Similarity =  
5%

Laboratory.java

GUILaboratoryData.java

```
/* *This class implements the GUI for  
managing laboratories data */  
  
public class GUILaboratoryData {  
    private JFrame window;  
    private JButton insert;  
    ...  
  
    public GUILaboratoryData() {  
        window = new JFrame("Laboratory GUI");  
        insert = new JButton("Insert Data");  
        ...  
    }  
    ...  
}
```

Similarity =  
42%

# Structural Information and Traceability

<b>Use Case</b>	<i>Insert Laboratory Data</i>
<b>Description</b>	<i>The user inserts the data of a specific laboratory</i>
<b>Events</b>	<ol style="list-style-type: none"><li>1. The user opens the <b>Laboratory GUI</b></li><li>2. The user inserts the <b>laboratory data</b></li><li>...</li><li>...</li><li>...</li></ol>

linked

linked

```
public class Laboratory{  
    private String name;  
    private String position;  
    ...  
    public void setName(String Name) {  
        this.name=pName;  
    }  
    ...  
}
```

Similarity =  
5%

Laboratory.java

GUILaboratoryData.java

```
/* *This class implements the GUI for  
managing laboratories data */  
  
public class GUILaboratoryData {  
    private JFrame window;  
    private JButton insert;  
    ...  
    public GUILaboratoryData() {  
        window = new JFrame("GUI Laboratory Data");  
        insert = new JButton("Insert");  
        ...  
    }  
    ...  
}
```

Similarity =  
42%

Structural dependency



# Structural Information and Traceability

<b>Use Case</b>	Insert Laboratory Data
<b>Description</b>	The user inserts the data of a specific laboratory
<b>Events</b>	<ol style="list-style-type: none"><li>1. The user opens the Laboratory GUI</li><li>2. The user inserts the laboratory data</li><li>.</li><li>.</li><li>.</li></ol>

linked

linked

```
public class Laboratory{  
    private String name;  
    private String position;  
    ...  
    public void setName(String Name) {  
        this.name=pName;  
    }  
    ...  
}
```

Similarity =  
5%

Laboratory.java

GUILaboratoryData.java

```
/* *This class implements the GUI for  
managing laboratories data */  
  
public class GUILaboratoryData {  
    private JFrame window;  
    private JButton insert;  
    ...  
  
    public GUILaboratoryData() {  
        window = new JFrame("GUI Laboratory Data");  
        insert = new JButton("Insert");  
        ...  
    }  
    ...  
}
```

Similarity =  
42%

Structural dependency

Transitivity

# Structural Information and Traceability

<i>Use Case</i>	<i>Insert Laboratory Data</i>
<i>Description</i>	<i>The user inserts the data of a specific laboratory</i>
<i>Events</i>	<ol style="list-style-type: none"><li>1. The user opens the Laboratory GUI</li><li>2. The user inserts the laboratory data</li><li>.</li><li>.</li><li>.</li></ol>

linked

linked

```
public class Laboratory{  
    private String name;  
    private String ...  
    ...  
    public void ... (name) {  
        this.name=p  
    }  
    ...  
}
```

Similarity =  
5% + **BONUS**

Laboratory.java

GUILaboratoryData.java

```
/* *This class implements the GUI for  
managing laboratories data */  
  
public class GUILaboratoryData {  
    private JFrame window;  
    private JButton insert;  
    ...  
    public GUILa  
    window = ne  
    insert = new  
    ...  
}  
  
...
```

Similarity =  
42%

Structural dependency

Transitivity

# *Open Issues*

*1) The choice of the bonus value is crucial*

- Different systems require different bonus*
- Different IR methods require different bonus*

# *Open Issues*

*1) The choice of the bonus value is crucial*

- Different systems require different bonus*
- Different IR methods require different bonus*

*2) When applying the bonus?*

- The transitivity property does not always hold*

# Example

Use Case	Insert Laboratory Data
Description	The user <i>inserts</i> the <i>data</i> of a specific laboratory
Events	1. The user opens the Laboratory GUI 2. The user <i>inserts</i> the laboratory data . . .

Not  
linked



## GUIDoctorData.java

```
/* *This class implements the GUI for
managing laboratories data */

public class GUIDoctorData {
    private JFrame window;
    private JButton insert;
    ...

    public GUIDoctorData() {
        window = new JFrame();
        insert = new JButton();
        ...
    }
}
```

# Example

<i>Use Case</i>	<i>Insert Laboratory Data</i>
<i>Description</i>	<i>The user inserts the data of a specific laboratory</i>
<i>Events</i>	<div>1. The user opens the Laboratory GUI</div> <div>2. The user inserts the laboratory data</div> <div>.</div> <div>.</div> <div>.</div>

Not  
linked



## GUIDoctorData.java

```
/* *This class implements the GUI for  
managing laboratories data */  
  
public class GUILaboratoryData {  
    private JFrame window;  
    private JButton insert;  
    ...  
  
    public GUILaboratoryData() {  
        window = new JFrame("Laboratory Data");  
        insert = new JButton("Insert");  
        ...  
    }  
  
    ...  
}
```

Similarity =  
42%

# Example

Use Case	Insert Laboratory Data
Description	The user inserts the data of a specific laboratory
Events	1. The user opens the Laboratory GUI 2. The user inserts the laboratory data . . .

Not  
linked



Not linked

```
public class Authorization{  
  
    public void setAuthorization(Doctor  
    pDoctor, Laboratory pLab){  
        ...  
    }  
  
    ...  
}
```

Authorization.java

GUIDoctorData.java

```
/* *This class implements the GUI for  
managing laboratories data */  
  
public class GUILaboratoryData {  
    private JFrame window;  
    private JButton insert;  
    ...  
  
    public GUILaboratoryData() {  
        window = new JFrame("Laboratory Data");  
        insert = new JButton("Insert");  
        ...  
    }  
  
    ...  
}
```

Similarity =  
42%

# Example

Use Case	Insert Laboratory Data
Description	The user inserts the data of a specific laboratory
Events	1. The user opens the Laboratory GUI 2. The user inserts the laboratory data . . .

Not  
linked



Not linked

```
public class LaboratoryAuthorization{  
  
    public void ... Doctor  
    pDoctor){ {  
        ...  
    }  
  
    ...  
}
```

Similarity =  
5%

## GUILaboratoryData.java

```
/* *This class implements the GUI for  
managing laboratories data */  
  
public class GUILaboratoryData {  
    private JFrame window;  
    private JButton insert;  
    ...  
  
    public GUILaboratoryData() {  
        window = new JFrame();  
        insert = new JButton();  
        ...  
    }  
  
    ...  
}
```

Similarity =  
42%

Structural dependency



LaboratoryAuthorization.java



# Example

Use Case	Insert Laboratory Data
Description	The user inserts the data of a specific laboratory
Events	1. The user opens the Laboratory GUI 2. The user inserts the laboratory data . . .

Not  
linked



Not linked

```
public class LaboratoryAuthorization{  
  
    public void ... Doctor  
    pDoctor){ {  
        ...  
    }  
  
    ...  
}
```

Similarity =  
5% + **BONUS**

## GUILaboratoryData.java

```
/* *This class implements the GUI for  
managing laboratories data */  
  
public class GUILaboratoryData {  
    private JFrame window;  
    private JButton insert;  
    ...  
  
    public GUILaboratoryData() {  
        window = new JFrame("GUI Laboratory Data");  
        insert = new JButton("Insert");  
        ...  
    }  
  
    ...  
}
```

Similarity =  
42%

Structural dependency

LaboratoryAuthorization.java

# Example

Use Case	Insert Laboratory Data
Description	The user inserts the data of a specific laboratory
Events	1. The user opens the Laboratory GUI 2. The user inserts the data

linked

## GUILaboratoryData.java

```
/* *This class implements the GUI for  
managing laboratories data */
```

```
public class GUILaboratoryData {  
    private JFrame window;  
    private JButton insert;
```

*Structural information is not  
always useful*

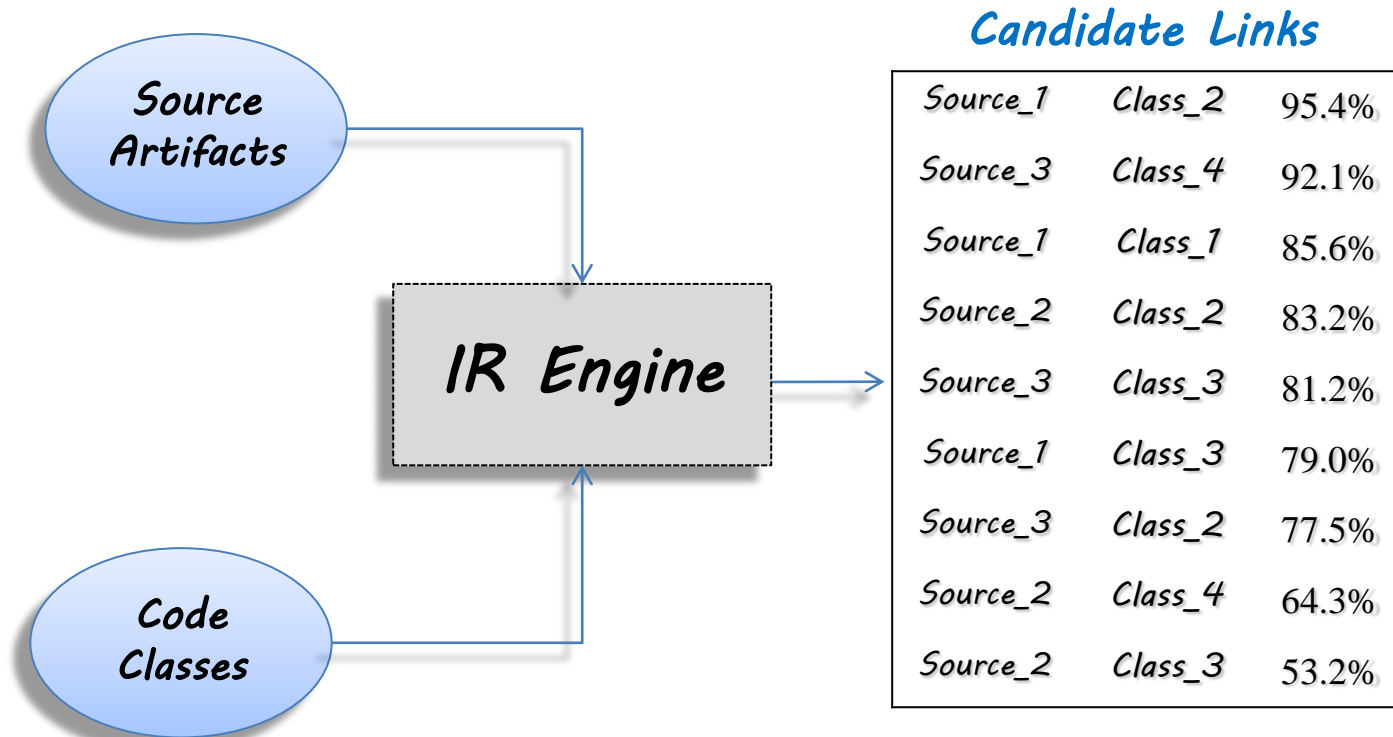
```
public class LaboratoryAuthorization{  
  
    public void setAuthorization(Doctor  
    pDoctor) {  
        ...  
    }  
  
    ...  
}
```

*Structural dependency*

Laboratory.java

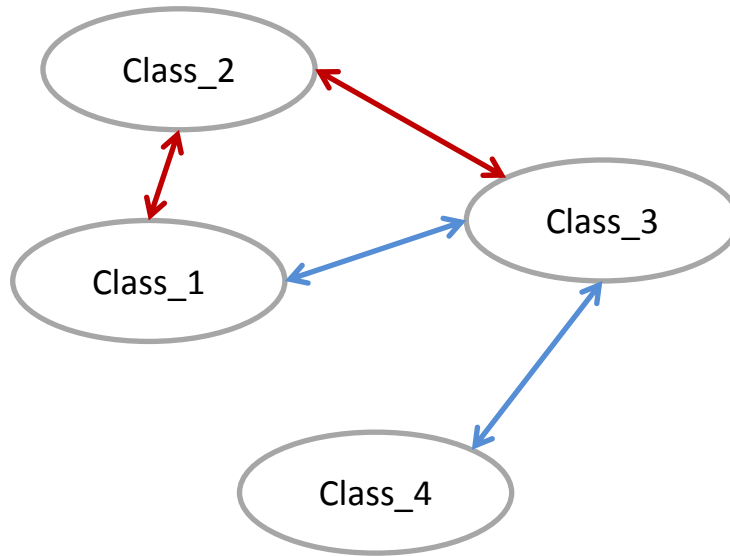
# Approach 1: Optimistic Combination (O-CSTI)

## Step 1: traditional IR process



# Approach 1: Optimistic Combination (O-CSTI)

Step 2: applying bonus to all the links



Candidate Links

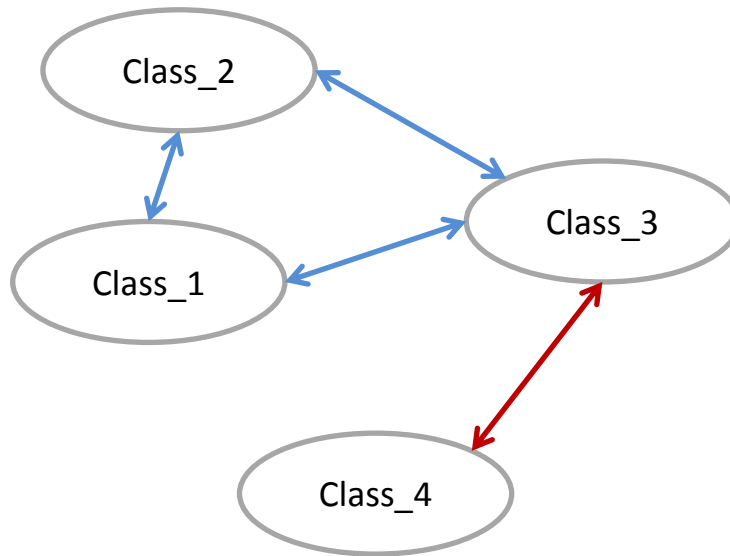
Source_1	Class_2	95.4%
Source_3	Class_4	92.1%
Source_1	Class_1	85.6%
Source_2	Class_2	83.2%
Source_3	Class_3	81.2%
Source_1	Class_3	79.0%
Source_3	Class_2	77.5%
Source_2	Class_4	64.3%
Source_3	Class_4	53.2%

+ Bonus

+ Bonus

# Approach 1: Optimistic Combination (O-CSTI)

Step 2: applying bonus to all the links



Candidate Links

Source_1	Class_2	95.4%
Source_3	Class_4	92.1%
Source_1	Class_1	85.6%
Source_2	Class_2	83.2%
Source_3	Class_3	81.2%
Source_1	Class_3	79.0%
Source_3	Class_2	77.5%
Source_2	Class_4	64.3%
Source_3	Class_4	53.2%

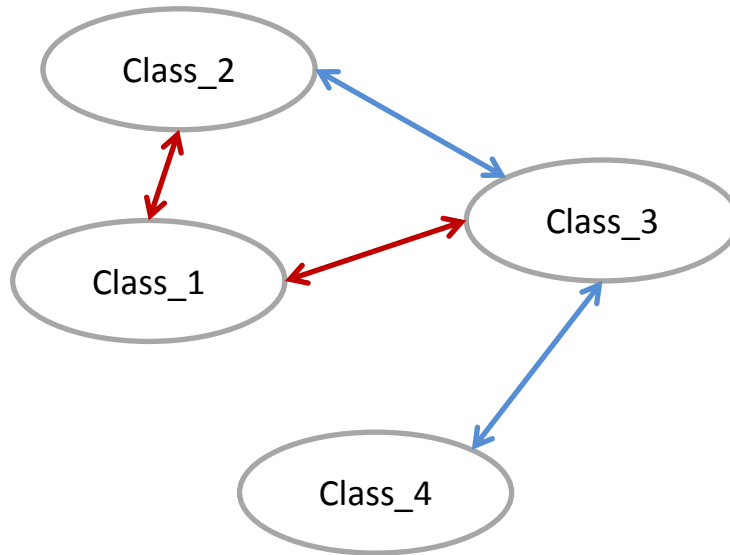
+ Bonus

+ Bonus

+ Bonus

# Approach 1: Optimistic Combination (O-CSTI)

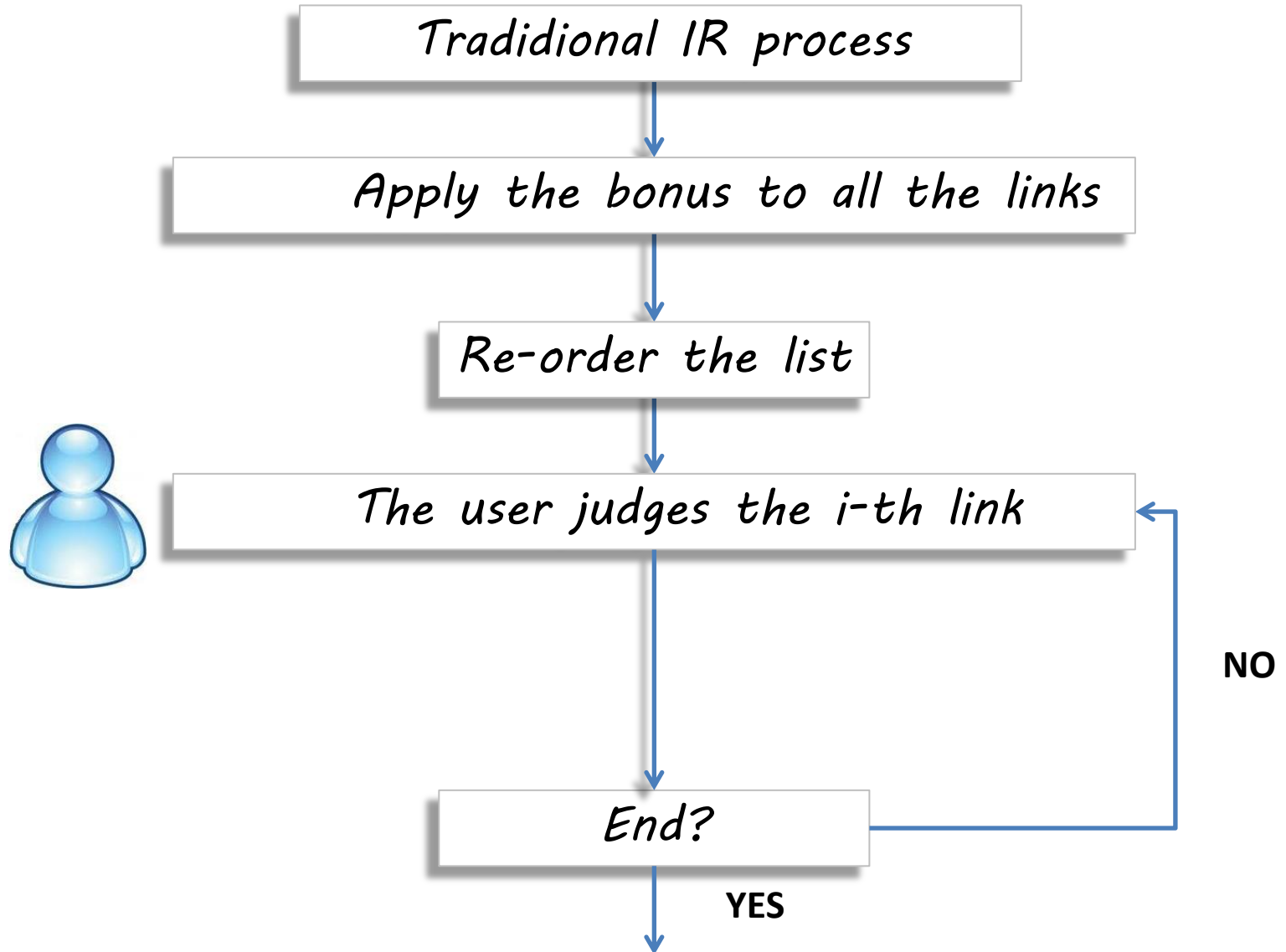
Step 2: applying bonus to all the links



Candidate Links

Source_1	Class_2	95.4%	+ Bonus
Source_3	Class_4	92.1%	
Source_1	Class_1	85.6%	+ Bonus
Source_2	Class_2	83.2%	
Source_3	Class_3	81.2%	
Source_1	Class_3	79.0%	+ 2Bonus
Source_3	Class_2	77.5%	
Source_2	Class_4	64.3%	
Source_3	Class_4	53.2%	+ Bonus

## *Approach 1: Optimistic Combination (O-CSTI)*



## Approach 2: User Driven Combination (U-CSTI)



Only the user can say  
If a link is **correct** or  
**not**



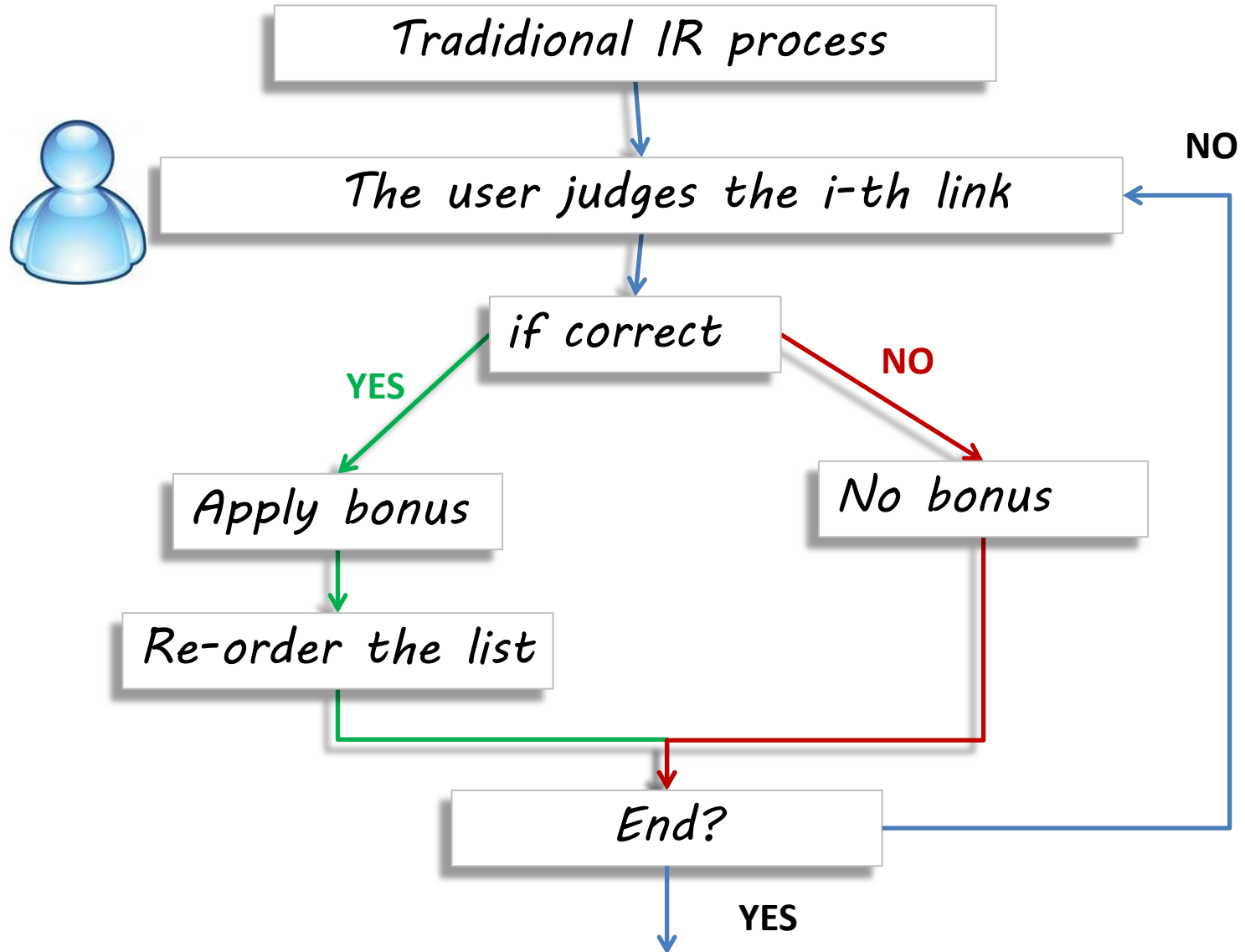
## Approach 2: User Driven Combination (U-CSTI)



Only the user can say  
If a link is **correct** or  
**not**

Only the user can say  
If applying the  
**combination** or **not**

## Approach 2: User Driven Combination (U-CSTI)



# Adaptive bonus

$$\text{Sim}(\text{Source1}, \text{Class1}) = \text{Sim}(\text{Source1}, \text{Class1}) + \delta * \text{Sim}(\text{Source1}, \text{Class1})$$

# Adaptive bonus

$$\text{Sim}(\text{Source1}, \text{Class1}) = \text{Sim}(\text{Source1}, \text{Class1}) + \delta * \text{Sim}(\text{Source1}, \text{Class1})$$

$$\delta = \frac{\max(\text{Sim}) - \min(\text{Sim})}{2} \approx \text{ranked list variability}$$

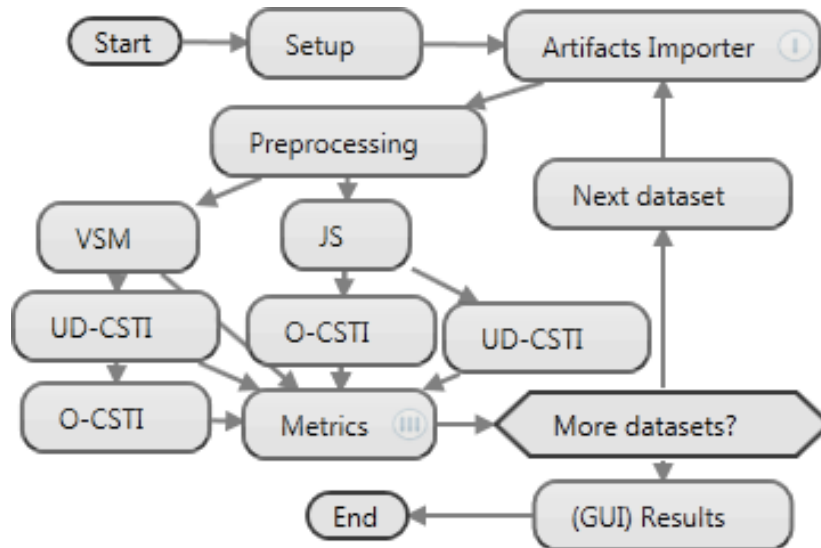
## Candidate Links

Source_1	Class_2	95.4%	← $\max(\text{Sim})$
Source_3	Class_4	92.1%	
Source_1	Class_1	85.6%	
Source_2	Class_2	83.2%	
Source_3	Class_3	81.2%	
Source_1	Class_3	79.0%	
Source_3	Class_2	77.5%	
Source_2	Class_4	64.3%	
Source_2	Class_3	53.2%	← $\min(\text{Sim})$

↑ variability ↓

# Implementation

## TraceLab Components

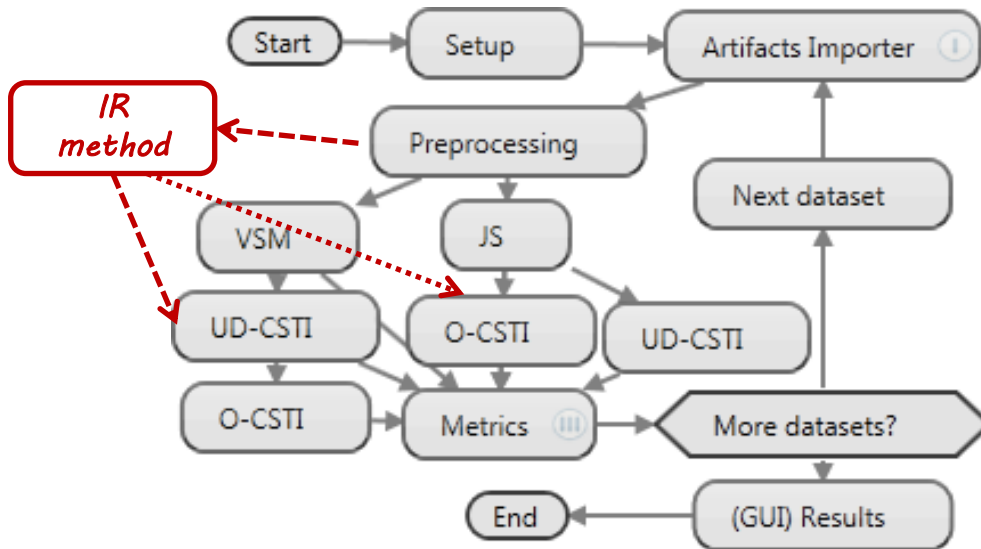


- 1) Adaptive Bonus
- 2) Optimistic Combination (O-CSTI)
- 3) User Driven Combination (U-CSTI)
- 4) Different IR Methods:
  - Vector Space Model
  - Jensen-Shannon

We provide the experiments and datasets for download at <http://www.cs.wm.edu/semeru/data/csmr13/>

# Implementation

## TraceLab Components



- 1) Adaptive Bonus
- 2) Optimistic Combination (O-CSTI)
- 3) User Driven Combination (U-CSTI)
- 4) Different IR Methods:
  - Vector Space Model
  - Jensen-Shannon

We provide the experiments and datasets for download at <http://www.cs.wm.edu/semeru/data/csmr13/>

# *Empirical Evaluation*



# Context

System	Description	KLOC	Source Artifact (#)	Target Artifact (#)	Correct links
EasyClinic	A system used to manage a doctor's office	20	UC (30)	CC (37)	93
			UML (20)	CC (37)	69
			TC (63)	CC (37)	204
eTour	An electronic touristic guide developed by students.	45	UC (58)	CC (174)	366
SMOS	A system used to monitor high school students	23	UC (67)	CC (100)	1,044
UC: Use case, TC: Test case, CC: Code class					

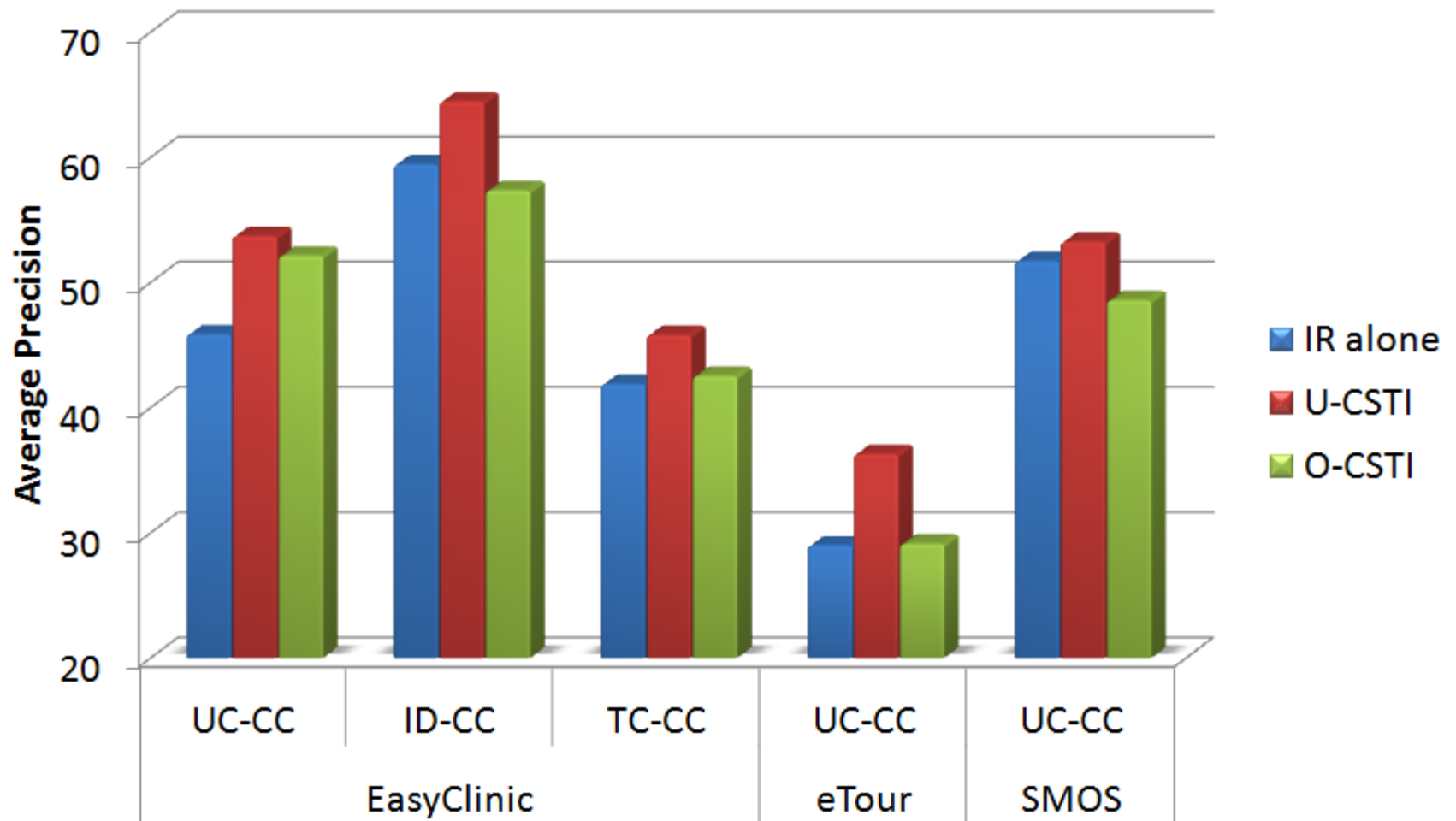
*We compared three IR-based processes:*

- 1) IR process alone*
- 2) O-CSTI (optimistic combination)*
- 3) U-CSTI (user driven combination)*



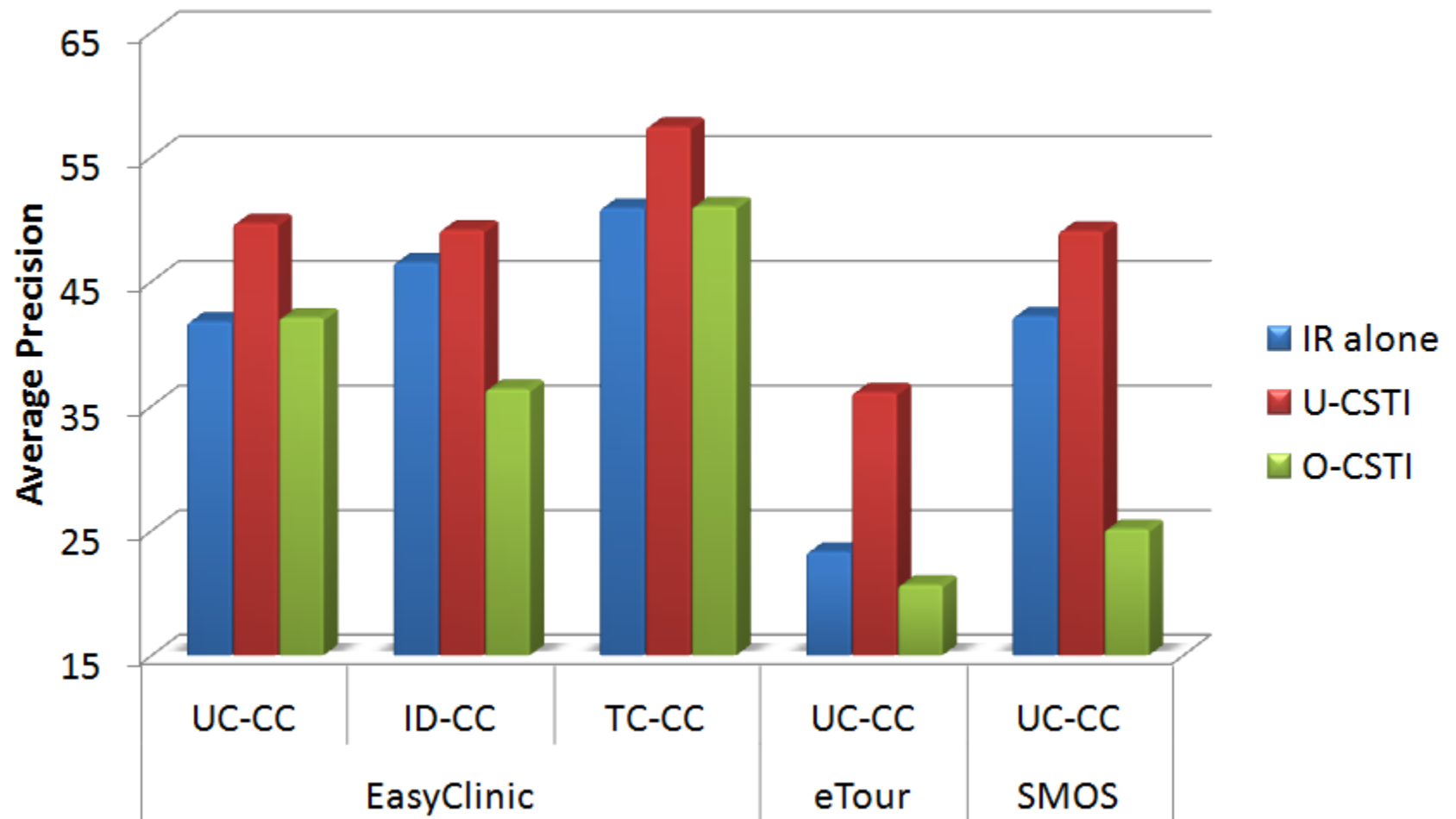
# Results

## *Vector Space Model*



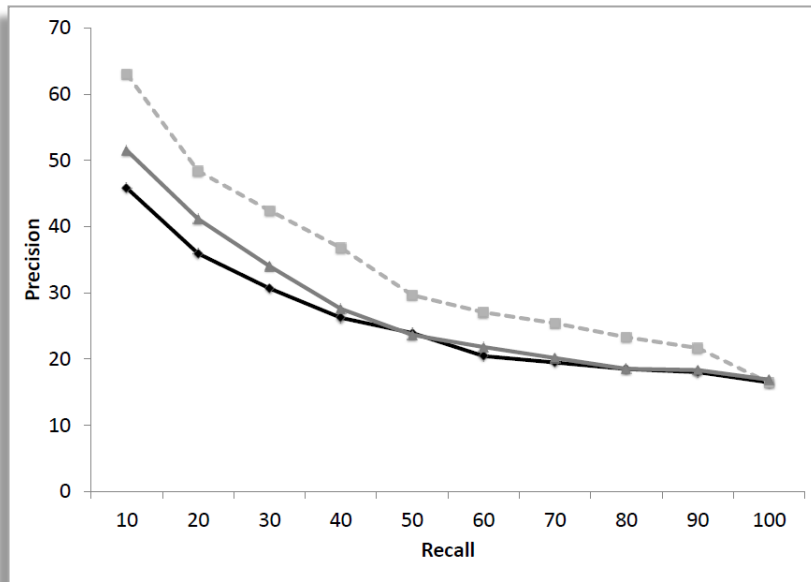
# Results

## *Jensen-Shannon Divergence*



# Results

*Tracing Use Cases onto Code Classes on SMOS*

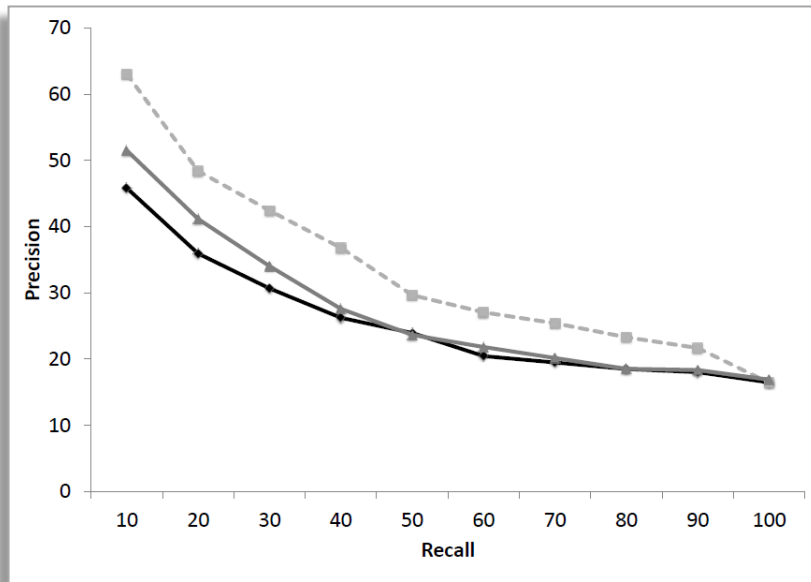


*Vector Space Model*

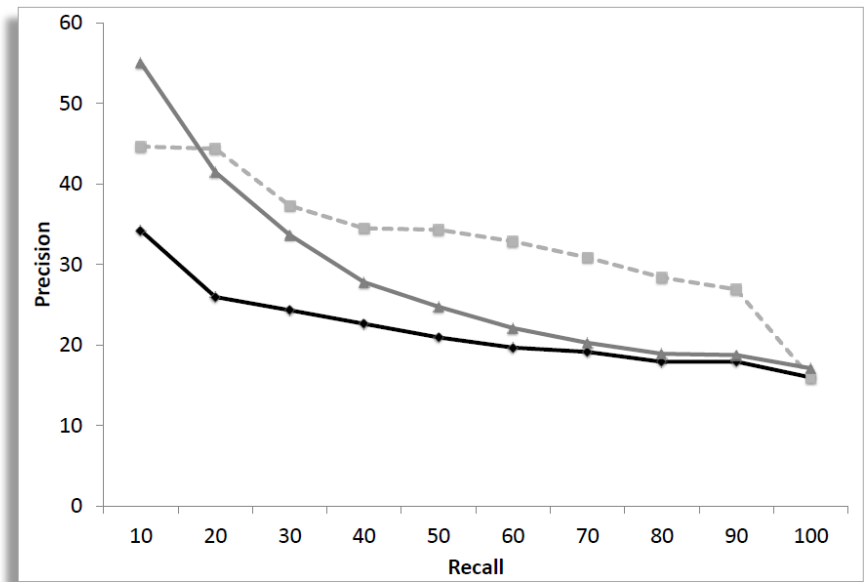
—◆— IR alone    —■— UD-CSTI    —▲— O-CSTI

# Results

*Tracing Use Cases onto Code Classes on SMOS*



*Vector Space Model*



*Jensen-Shannon Divergence*

—◆— IR alone    —■— UD-CSTI    —▲— O-CSTI

# *Optimality of the Adaptive Bonus*

*Adaptive Bonus vs Fixed Bonus*

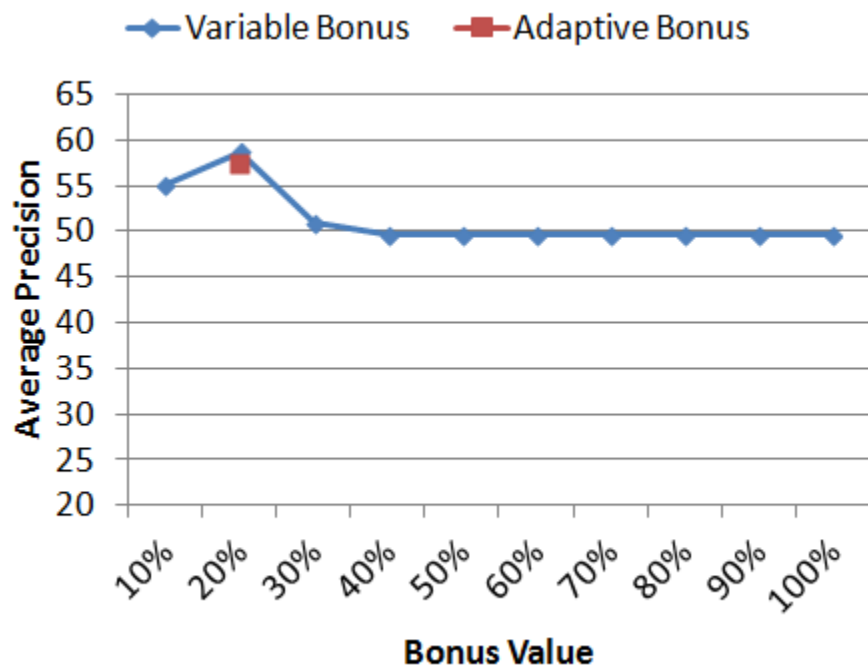
- We used different fixed bonus values*

# Optimality of the Adaptive Bonus

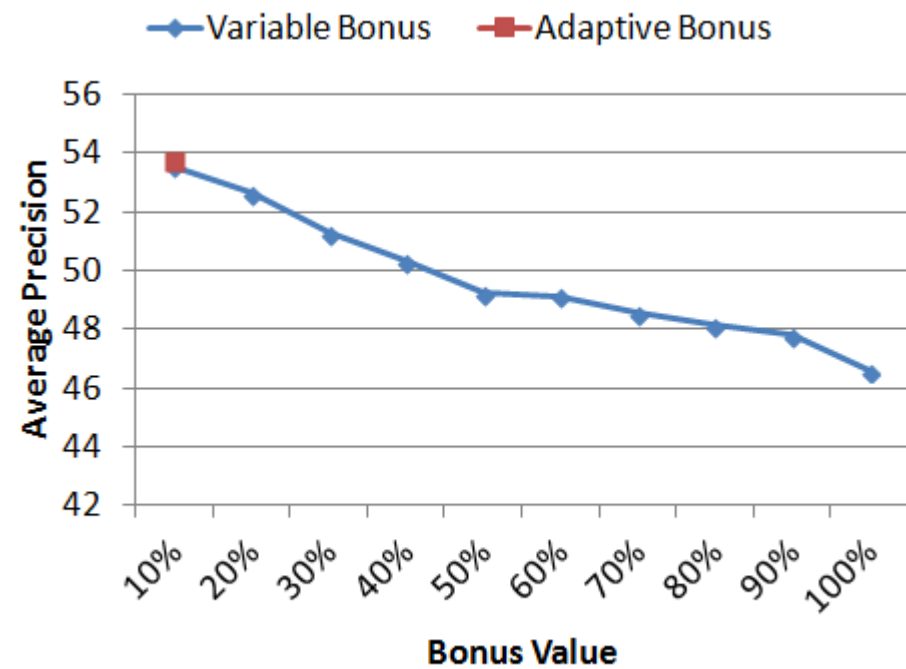
## Adaptive Bonus vs Fixed Bonus

- We used different fixed bonus values

EasyClinic UC-CC with VSM

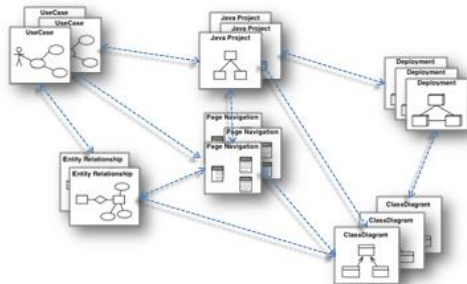


EasyClinic TC-CC with JS



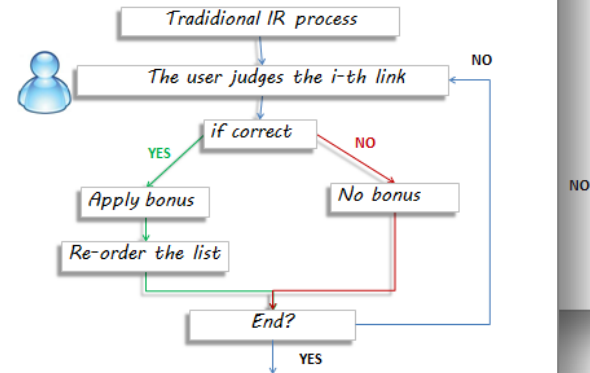
# Conclusions

*Traceability Recovery* is the ability to describe and follow the artifacts life-cycle



*Approache 1: Optimistic Combination (O-CSTI)*

*Approache 2: User Driven Combination (U-CSTI)*



*Adaptive bonus*

$$Sim(SourceI, ClassI) = Sim(SourceI, ClassI) + \delta * Sim(SourceI, ClassI)$$

$$\delta = \frac{\max(Sim) - \min(Sim)}{2} \approx \text{ranked list variability}$$

*Candidate Links*

Source_1	Class_2	95.4%	← max(Sim)
Source_3	Class_4	92.1%	
Source_1	Class_7	85.6%	
Source_2	Class_2	83.2%	
Source_3	Class_3	81.2%	
Source_1	Class_3	79.0%	
Source_3	Class_2	77.5%	
Source_2	Class_4	64.3%	
Source_2	Class_3	53.2%	← min(Sim)

← max(Sim)

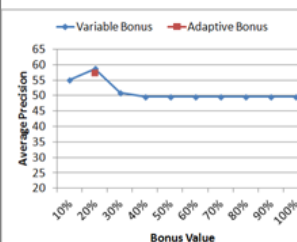
variability

← min(Sim)

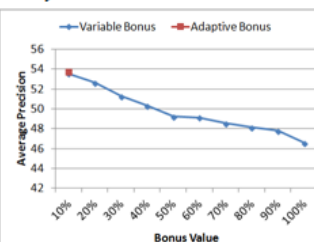
*Results*

*Optimality of the Adaptive Bonus*

*EasyClinic UC-CC with VSIM*



*EasyClinic TC-CC with JS*



alone  
-CSTI  
-CSTI

THANKS