



Linking Source Code to Untangled Change Intents

Xiaoyu Liu¹ LiGuo Huang¹ Chuanyi Li^{1,2} Vincent Ng³

- ¹ Dept. of Computer Science and Engineering, Southern Methodist University
 - ² State Key Laboratory for Novel Software Technology, Nanjing University
- ³ Human Language Technology Research Institute, University of Texas at Dallas

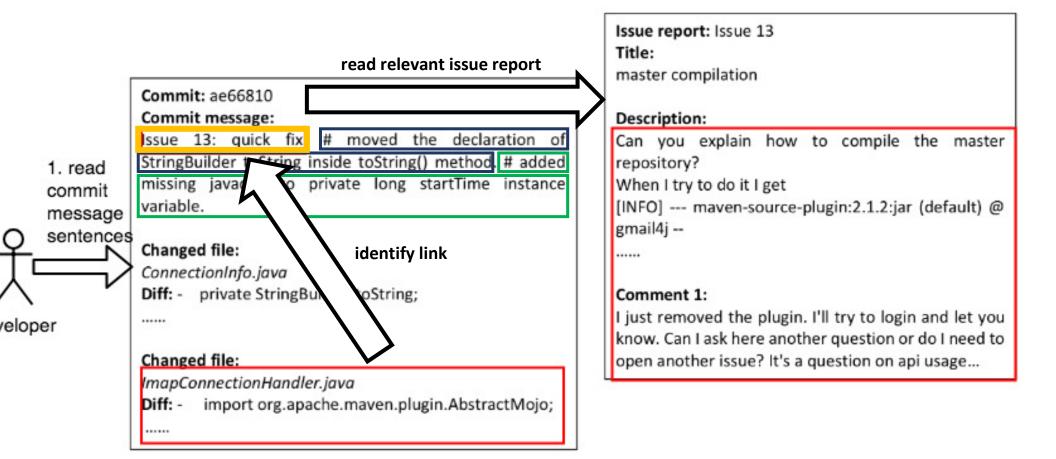
34th IEEE International Conference on Software Maintenance and Evolution Madrid, Spain September 23-29, 2018

@SMU-CSE

/06/2018

Motivating Example





MU.



Key Challenges

Change intents are tangled together in commit message

Manually link changed source code to change intents take lots of human efforts

MU.

UT D

Key Challenges

Change intents are tangled together in commit message

Solution: Untangle change intents

Manually link changed source code to change intents take lots of human efforts

Solution: An automated approach





kisting Software Artifacts Linking Approaches

Manual approaches

Problem: Time-consuming, labor-intensive, and requires a great deal of experience

Information Retrieval (IR)-based approaches

- Use a IR-based model: measures similarity between changed source code and untangled change intent
- Problem: Changed entities extracted from source code could be very different from what is described in commit messages and other related software documents

/06/2018





Goal

Propose approaches to address the task of linking changed source code to untangled change intent **at the sentence level**





Dur Approaches

AutoCILink-P

Automatically identify links by using manually defined patterns

AutoCILink-ML

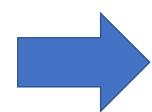
Automatically identify links by applying supervised learning



attern-based Link Identification System (AutoCILink-P)

ntangling change intents

Issue 13: quick fix. # moved the declaration of StringBuilder inside toString() method. # added missing javadoc to provide long startTime instance variable



- Issue 13: quick fix.
- # moved the declaration of StringBuilder inside toString() method.
- # added missing javadoc to provide long startTime instance variable





attern-based Link Identification System (AutoCILink-P)

Enriched untangled change intents

Commit: ae66810

Commit message:

ssue 13: quick fix. # moved the declaration of StringBuilder coString inside toString() method. # added missing javadoc

to private long startTime instance variable.

Issue report: Issue 13

Title:

master compilation

Description:

Can you explain how to compile the master repositor When I try to do it I get

. . .

Text
Preprocessing

Link identification using regular expression

Link identification vocabulary similar



attern-based Link Identification System (AutoCILink-P)

Extracting terms from enriched change intents and changed entities

om changed source code

"# moved the declaration of StringBuilder inside toString() method."



Terms: "declaration", "string", "builders", ...

Changed file:

SimulationRun.java

Diff:

+ eventHandlerMap.put(SimulationEvent.-TYPE_ACQUIRE_JOB_NOTIFICATION_EVENT, new AcquireJobNotificationEventHandler(job-Executor));

.....

Changed entities:

"event", "handler", "map",

"simulation", "type", "job",...

Text Preprocessing

Link identification using regular expression

If no link

Link identification vocabulary similar

1



attern-based Link Identification System (AutoCILink-P)

nerating regular expressions

```
^.*\s*(<verb>)(.*?)<entity>(.*)
```

<verb>: collected from change types
(e.g., move)

<entity>: noun words in changed sou
code identifiers, comments or string
literals (e.g., builder)

Ext reprocessing regular expression Link identification using regular expression vocabulary similarity



attern-based Link Identification System (AutoCILink-P)

```
nerating regular expressions
```

```
^.*\s*(<verb>)(.*?)<entity>(.*)
```

^.*\s*move.*?builder.*

"moved the declaration of stringbuild inside toString() method"

ext reprocessing Link identification using regular expression

If no link

Link identification using vocabulary similarity

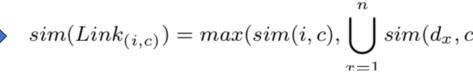


attern-based Link Identification System (AutoCILink-P)

cabulary similarity between enriched untangled change intent terms d changed source code entities

$$im(i,c) = \frac{\sum_{t \in V, e \in V} \omega(t,i,V) \times \omega(e,c,V)}{\sqrt{\sum_{t \in V} \omega(t,i,V)^2} \times \sqrt{\sum_{e \in V} \omega(e,c,V)^2}}$$

$$m(d,c) = \frac{\sum_{t \in V, e \in V} \omega(t, d, V) \times \omega(e, c, V)}{\sqrt{\sum_{t \in V} \omega(t, d, V)^2} \times \sqrt{\sum_{e \in V} \omega(e, c, V)^2}}$$



ext reprocessing Link identification using regular expression

If no link

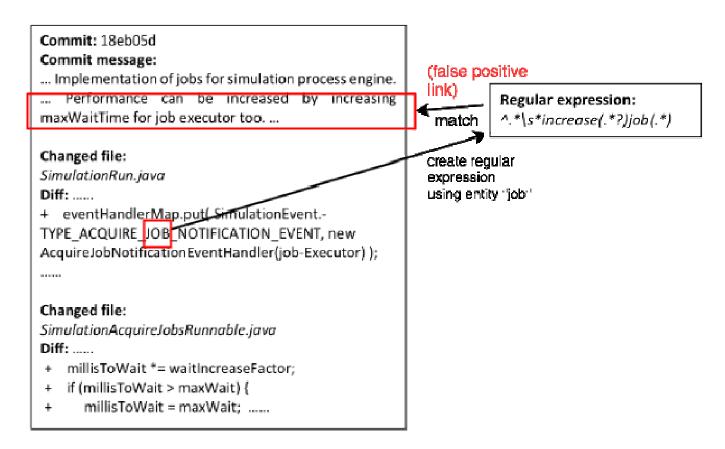
Link identification using vocabulary similarity



The Unit

AutoClLink-P's Key Weakness

Regular expressions are not always precise





AutoClLink-P's Key Weakness

Impreciseness of regular expressions

• Solution: employs a learning-based link classification system to weigh the importance of matched regular expressions



upervised learning-based link classification system (AutoCILink-N

Create Training
Instances

Train
Link Classifier

Classify Linked or

Not-linked in Test

- Create training instances for each changed source code and enriched untangled change intent pair in the training dataset.
- Each instance annotated as *linked* or *not linked*depending on whether there is a link
- Each instance is represented using 6 features

Set



The Univ

eatures

Type 1: Regular Expression Features encode the *presence* or *absence* of the regular expression in the training set

Type 2-4: Three types of Vocabulary Features

The Univ

ee Types of Vocabulary Features

pe 2: Vocabulary Pair Features

be 3: Vocabulary Similarity

atures

be 4: Term Unmatched Features

"Implementation of jobs for simulation

process engine"

Changed file:

SimulationRun.java

Diff:

+ eventHandlerMap.put(SimulationEv

TYPE_ACQUIRE_JOB_NOTIFICATION_EV

AcquireJobNotificationEventHandler(job

......

(process, job)



ee Types of Vocabulary Features

e 2: Vocabulary Pair Features

e 3: Vocabulary Similarity

tures

e 4: Term Unmatched Features

Encode the vocabulary simila scores between an enriched untangled change intent and changed source code



ee Types of Vocabulary Features

e 2: Vocabulary Pair Features

e 3: Vocabulary Similarity

tures

e 4: Term Unmatched Features

```
Terms: "perform", "increas", "max", "wait", "tim
"job", "executor"

Entities: "job", "notif", "time",
"event", "executor", "handle",...
```

3 out of 7 (42.9%) matched, so the percentage of not terms unmatched in changed entities (1-42.9%=57 falls in the range of [50%,60%)



UT The Uni

eatures

Type 5: Code Import Features encode the percentage of terms unmatched each imported code module changed entities

Type 6: Untangled Change Intent Count Features encode the number of untangled change intents in corresponding commit message



ervised learning-based link classification system (AutoCILink-ML

Create Training
Instances

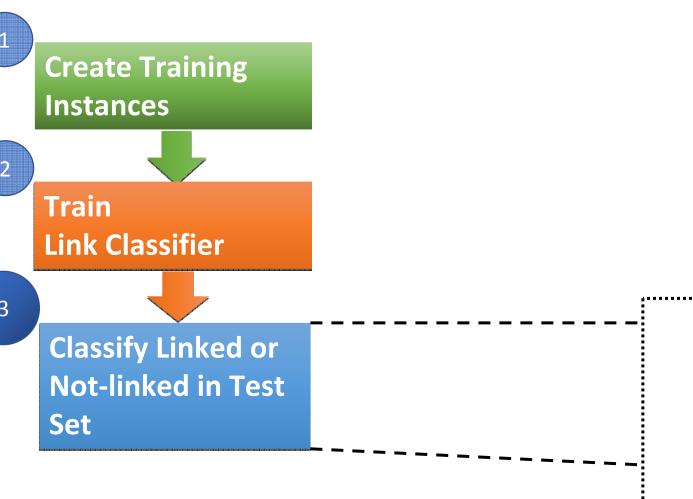
Train
Link Classifier

AutoClLink-ML uses the SVM learning algorithm to train a link classifier.

Classify Linked or Not-linked in Test Set



pervised learning-based link classification system (AutoCILink-M



AutoCILink-ML determines whether the given untangled change intent an changed source code file ar "linked" or "not linked".

MU. opirical Evaluation



sets: Open Source Java Projects from GitHub

.9 Projects from GitHub, 572 untangled change intents, 2739 changed source code fil 025 "linked" code-intent pairs (70.1%), and 1288 "not linked" code-intent pairs (29.9

ive-Fold Cross Validation

rics

ecall, Precision, F1-score

Accuracy: percentage of code-intent pairs correctly classified

line Systems

R-based Systems: LSI, VSM, Association-based

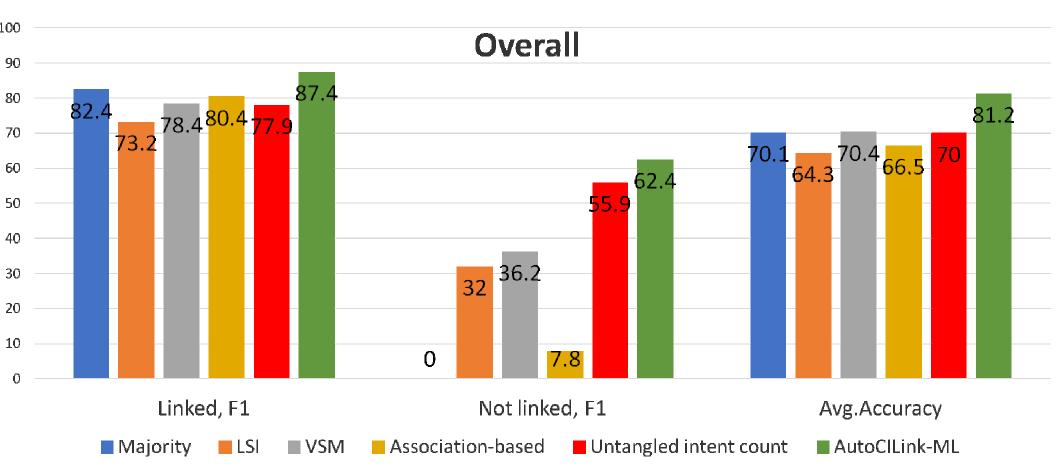
Majority Classifier: a greedy approach simply classify into the majority class (i.e., "link

Intangled Intent Count Classifier: classify using threshold based on untangled chatents count



verall Performance

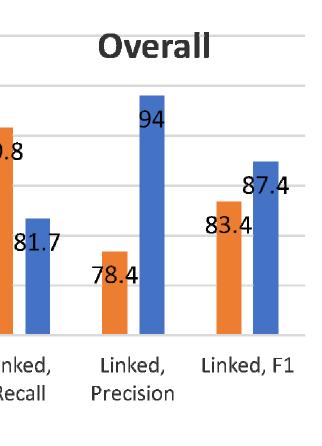
1. How effective is AutoCILink in linking changed code to untangled change intents

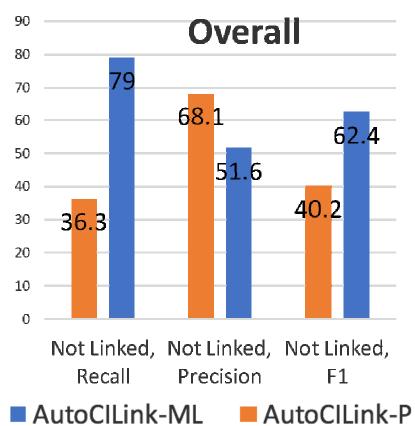


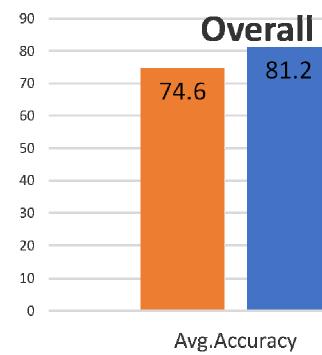
The Univ

itoCILink-ML vs. AutoCILink-P

2. Which system is more accurate in linking changed code to untangled change ents, AutoCILink-ML or AutoCILink-P







UT D

ature Ablation

. Which feature types have the largest impact on the performance of CILink-ML?

The most important features are *Untangled Change Intent Count Features* and *Term Unmatched Features*, as their removal results in a 18.1% and 12.2% drop in accuracy, respectively.

/06/2018

@SMU-CSE