



# Speculative Analysis for Comment Quality Assessment

Pooja Rani

Advisor: Prof. Oscar Nierstrasz

Software Composition Group

University of Bern, Switzerland

# Motivation

```
/**
 *A class representing a window on the screen.
 *For example:
 *<pre>
 *Window win = new Window(parent);
 *win.show();
 *</pre>
 *
 * @author Sami Shaio
 * @version 1.13, 06/08/06
 * @see java.awt.BaseWindow
 */

class Window extends BaseWindow{
  ..
}
```

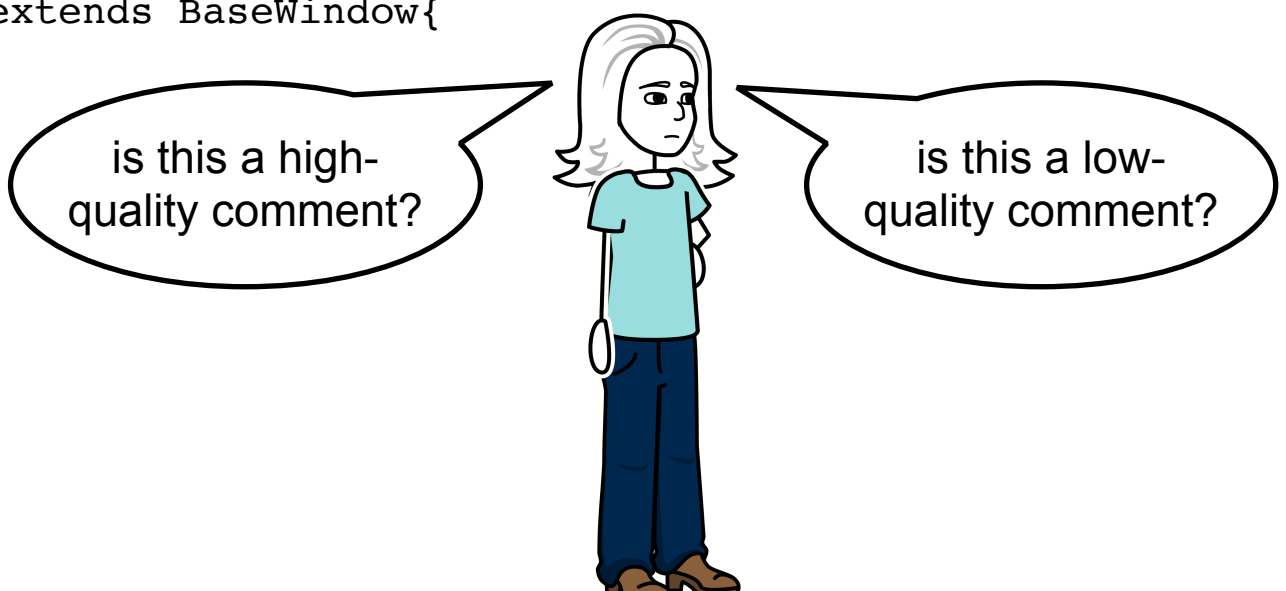
Trustworthy form of documentation

High-quality code comments assist developers

# Problem

```
/**  
 *A class representing a window on the screen.  
 *For example:  
 *<pre>  
 *Window win = new Window(parent);  
 *win.show();  
 *</pre>  
 *  
 * @author Sami Shaio  
 * @version 1.13, 06/08/06  
 * @see java.awt.BaseWindow  
 */
```

```
class Window extends BaseWindow{  
  ..  
}
```



is this a high-quality comment?

is this a low-quality comment?

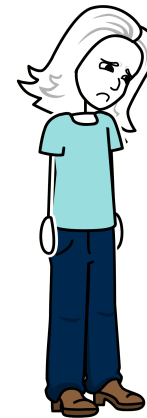
# Challenges

No standard definition of comment quality

No strict syntax and structure conventions

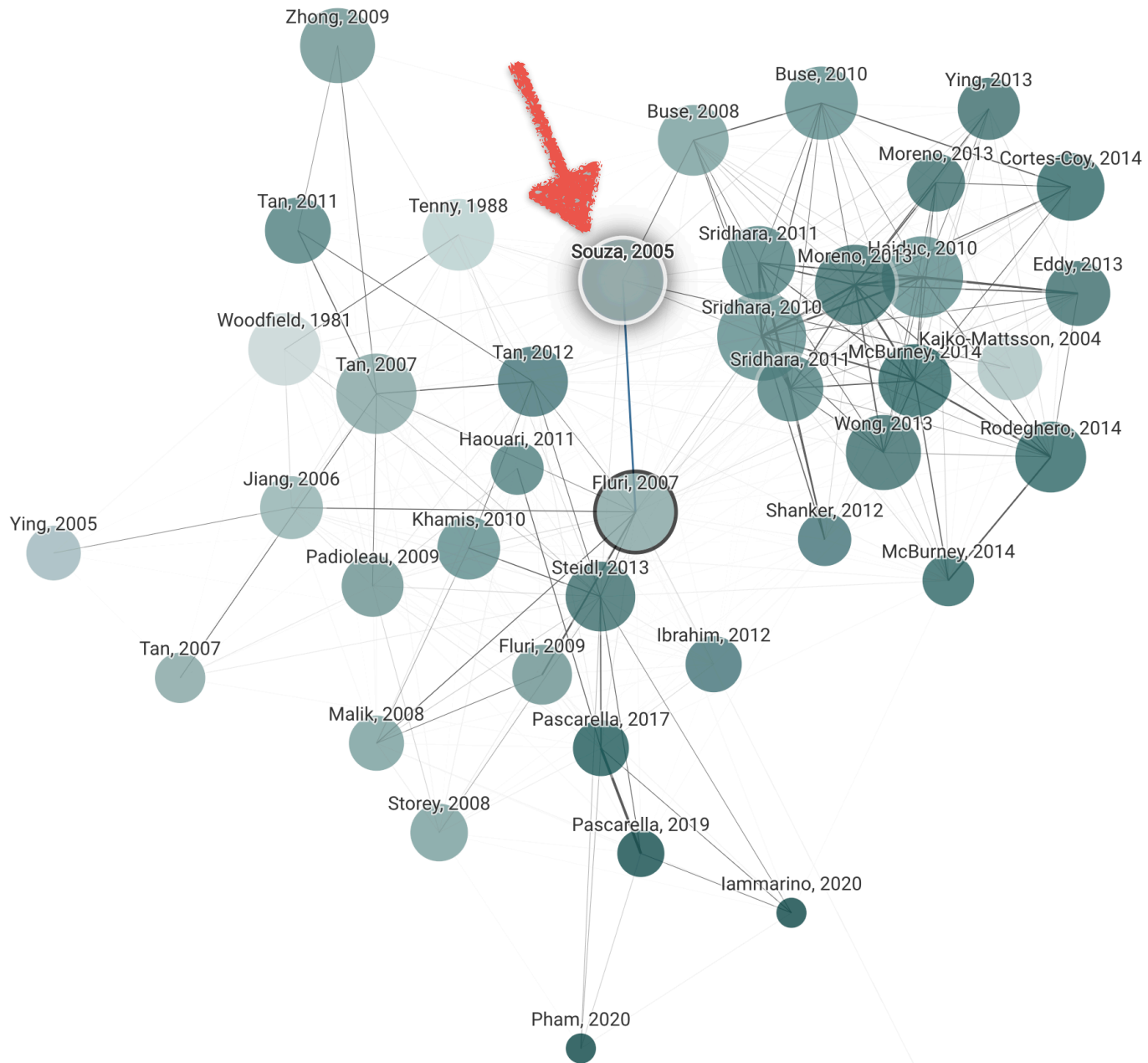
The compiler does not check comments

Lack of quality assessment tools

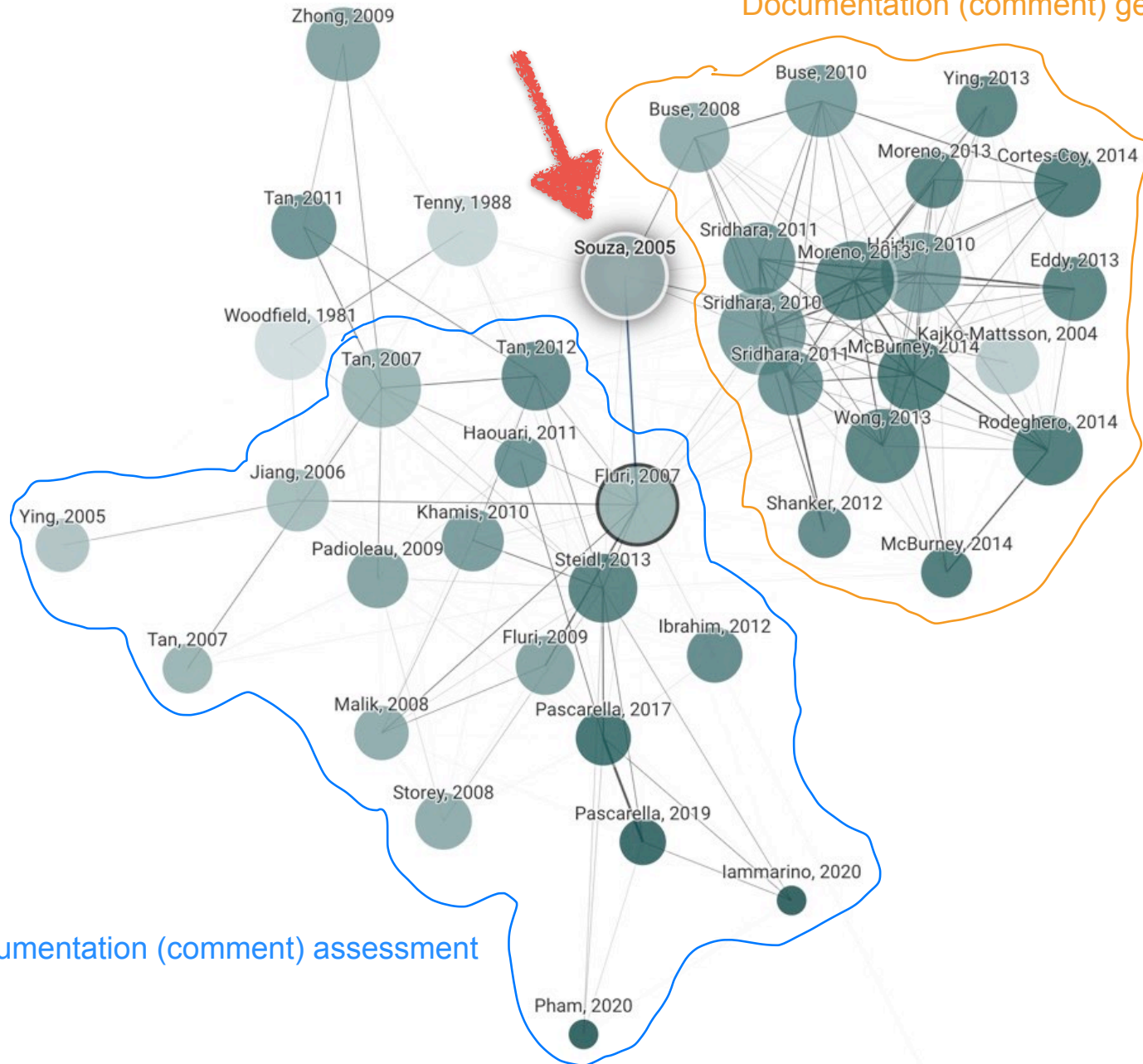


Makes quality assessment a non-trivial problem

# State of the art

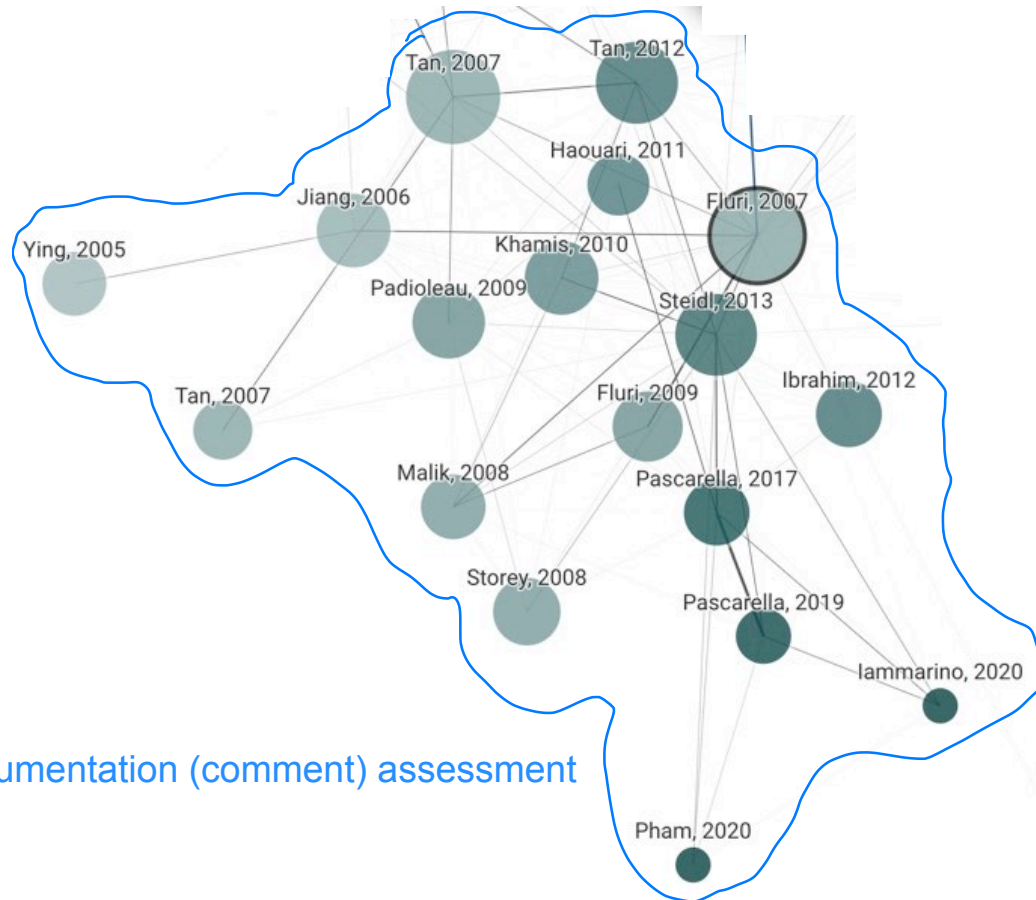


## Documentation (comment) generation



## Documentation (comment) assessment

Not all of them focused mainly on comments or all aspects of comments



Documentation (comment) assessment



Given the increasing use of multi-language environment, we need a deeper understanding of **developer commenting practices** and **concerns** to achieve high-quality comments

# Thesis statement

*“Understanding the specification of high-quality comments to build effective assessment tools requires a multi-perspective view of the comments”.*

# Thesis statement

*“Understanding the specification of high-quality comments to build effective assessment tools requires a multi-perspective view of the comments. The view can be approached by”*

**RQ1:** What do developers write in comments across languages?

# Thesis statement

*“Understanding the specification of high-quality comments to build effective assessment tools requires a multi-perspective view of the comments. The view can be approached by”*

**RQ1:** What do developers write in comments across languages?

**RQ2:** What do developers ask about code commenting practices?

# Thesis statement

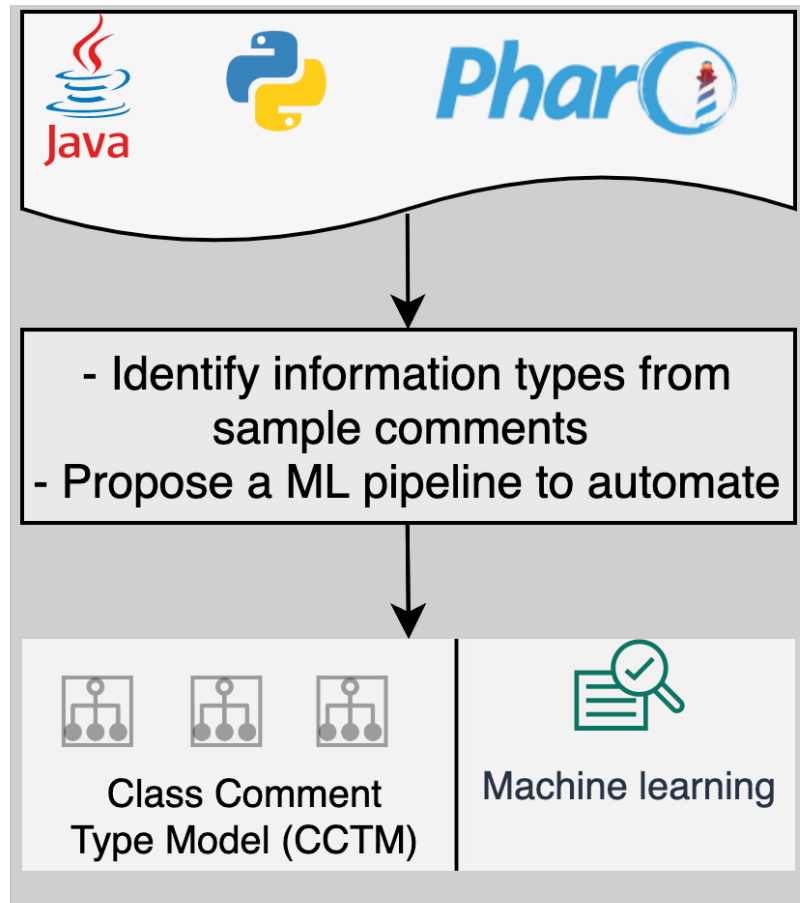
*“Understanding the specification of high-quality comments to build effective assessment tools requires a multi-perspective view of the comments. The view can be approached by”*

**RQ1:** What do developers write in comments across languages?

**RQ2:** What do developers ask about code commenting practices?

**RQ3:** What quality attributes are often considered in assessing comment quality?

# RQ1: What do developers write in comments across languages?



Java, Python, Smalltalk,  
**20 projects**

**Extract** class comments,  
**37,446 comments**

**Classify** sample comments,  
**1,066 comments**

**Output:** a taxonomy, and a classifier

# RQ1: What do developers write in comments across languages?

```
/**
 *A class representing a window on the screen.
 *For example:
 *<pre>
 *Window win = new Window(parent);
 *win.show();
 *</pre>
 *
 *@author Sami Shaio
 *@version 1.13, 06/08/06
 *@see java.awt.BaseWindow
 */
```

Summary

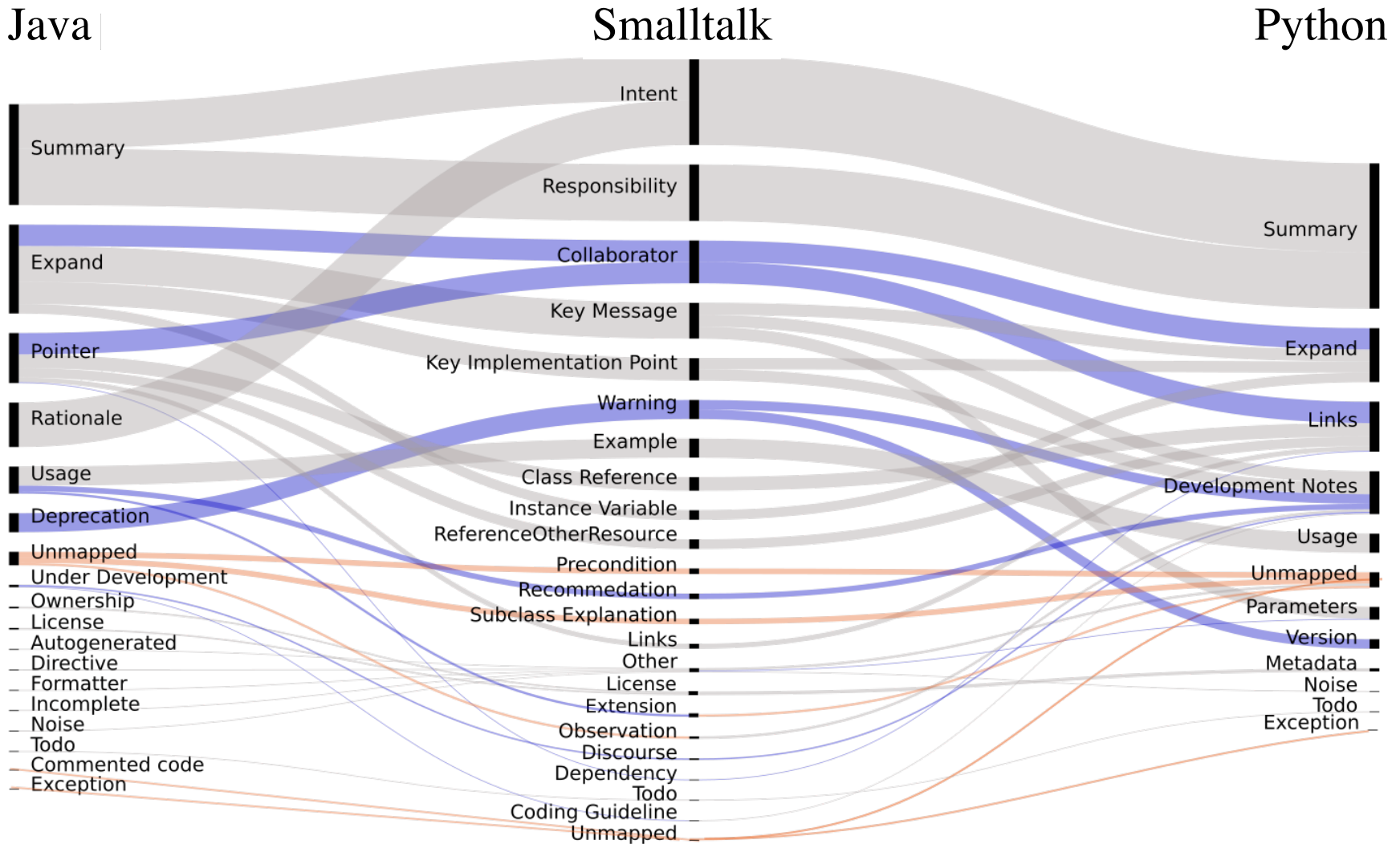
Usage

Ownership

Pointer

```
class Window extends BaseWindow{
 ..
 }
```

# RQ1: What do developers write in comments across languages?





# Automatically identify an information type

```
/**  
 *A class representing a window on the screen.  
 *For example:  
 *<pre>  
 *Window win = new Window(parent);  
 *win.show();  
 *</pre>  
 *  
 *@author Sami Shaio  
 *@version 1.13, 06/08/06  
 *@see java.awt.BaseWindow  
 */
```

Summary

Pointer

```
class Window extends BaseWindow{  
  ..  
}
```

Recurrent natural language patterns in writing a specific type of information

# Automatically identify an information type

```
/**  
 *A class representing a window on the screen.  
 *For example:  
 *

```
  
 *Window win = new Window(parent);  
 *win.show();  
 *
```

  
 *  
 *@author Sami Shaio  
 *@version 1.13, 06/08/06  
 *@see java.awt.BaseWindow  
 */
```

Summary

Represents [something]

[verb]s [noun]

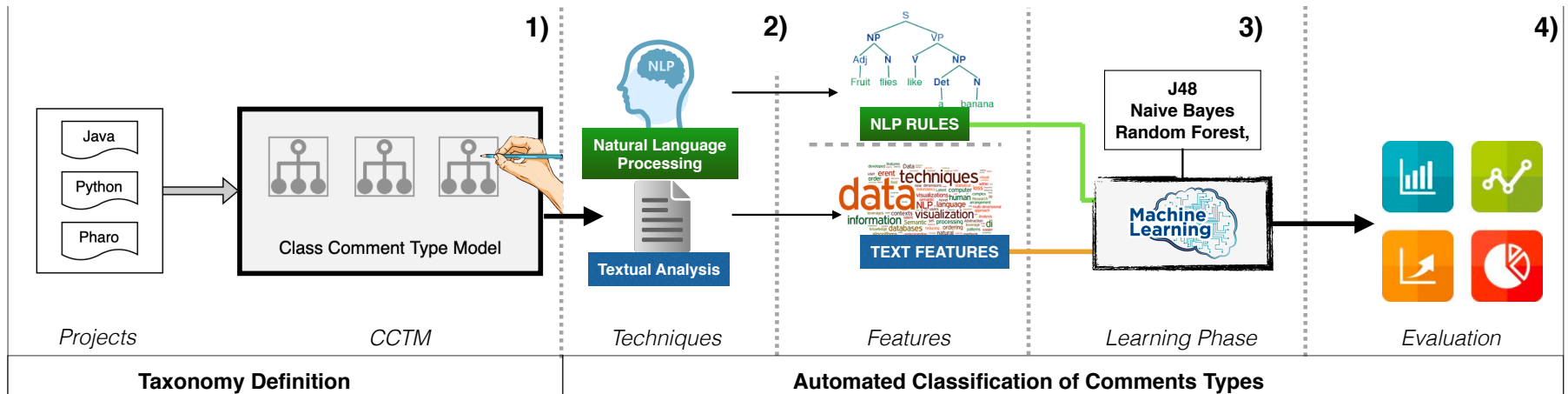
Pointer

Sees [something]

```
class Window extends BaseWindow{  
  ..  
}
```

Recurrent natural language patterns in writing a specific type of information

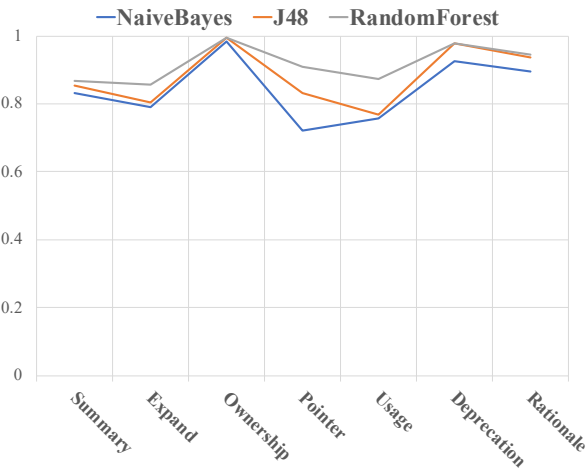
# Automatically identify an information type



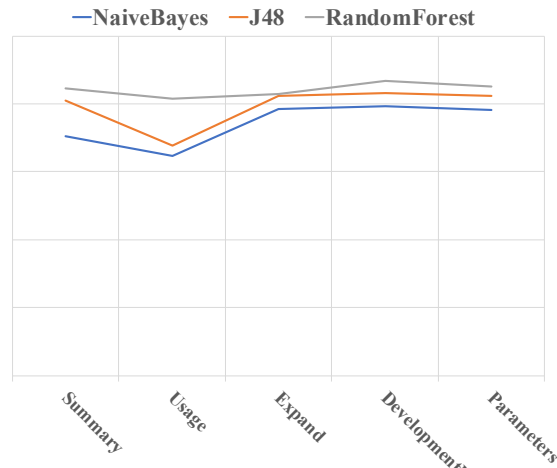
Ground truth: manually analyzed comments

Used recurrent patterns as feature set + text features

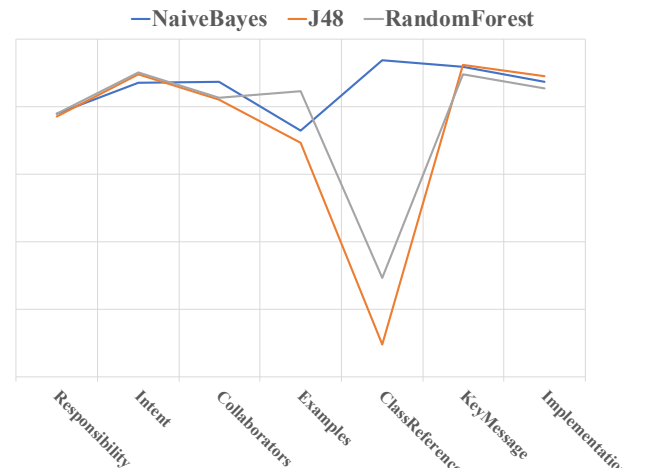
# Automatically identify an information type



Java



Python



Smalltalk

Random Forest technique classifies comments better

# Takeaway

Developers embed at least 10 types of information in comments across languages

Using recurrent natural language patterns as features improves the classification

We are currently testing our classification approach using the deep learning technique (FastText)

# Timeline & Contribution

*What do class comments tell us? An investigation of comment evolution and practices in Pharo Smalltalk*

An investigation of comment evolution, content and their adherence to style guideline

*How to Identify Class Comment Types? A Multi-language Approach for Class Comments Classification*

An investigation of comments in Java, Python, Smalltalk, and develop an approach to automatically identify information types across languages

In the stage of acceptance

## What do class comments tell us? An investigation of comment evolution and practices in Pharo Smalltalk

Pooja Rani · Sebastiano Panichella · Manuel Leuenberger · Mohammad Ghafari · Oscar Nierstrasz

Received: date / Accepted: date

**Abstract** Previous studies have characterized code comments in various programming languages, and have shown how a high quality of code comments is crucial to support program comprehension activities and to improve the effectiveness of maintenance tasks. However, very few studies have focused on the analysis of the information embedded in code com-

ments. None of them has compared developer practices to w

standard guidelines or analyzed these characteristics in the Pharo Smalltalk environment. These class commenting practices have their origins in Smalltalk 80. Over the years, particularly in the Pharo environment. This paper study investigating commenting practices in Pharo Smalltalk class comment evolution over seven Pharo versions. Then, we analyze class comments of the most recent version of Pharo to identify the information types of Pharo comments. Finally, we study the adherence of class comments to the class template over Pharo versions.

The results of this study show that there is a rapid increase in the number of class comments in the latest three Pharo versions, while in subsequent versions develop new and old classes, thus maintaining a similar ratio. In addition, the results of the comments from the latest Pharo version suggests that the information types typically embedded in class comments by developers and that of the latest *Pharo class comment template*. However, the information types in the standard template tend to be present more often than other types. Finally, we find that a substantial proportion of comments follow the standard template in writing these information types, but they are written in a way. This suggests the need to standardize the commenting g

Accepted  
Empirical Software  
Engineering (EMSE),  
2021

## How to Identify Class Comment Types? A Multi-language Approach for Class Comments Classification

Pooja Rani<sup>1</sup>, Sebastiano Panichella<sup>2</sup>, Manuel Leuenberger<sup>3</sup>, Andrea Di Sorbo<sup>4</sup>, Oscar Nierstrasz<sup>5</sup>

<sup>1</sup>Software Composition Group, University of Bern, Switzerland  
<sup>2</sup>Zurich University of Applied Sciences, Switzerland  
<sup>3</sup>Department of Engineering, University of Sannio, Italy

### Abstract

Most software maintenance and evolution tasks require developers to understand the source code of their software systems. Software developers usually inspect class comments to gain knowledge about program behavior, regardless of the programming language they are using. Unfortunately, (i) different programming languages present language-specific code commenting notations and guidelines; and (ii) the source code of software projects often lacks comments that adequately describe the class behavior, which complicates program comprehension and evolution activities. To handle these challenges, this paper investigates the different language-specific class commenting practices of three programming languages: Python, Java, and Pharo. In particular, we systematically analyze the similarities and differences of the information types found in class comments of projects developed in these languages. We propose an approach that leverages two techniques — namely Natural Language Processing and Text Analysis — to automatically identify *class comment types*, i.e., the specific types of semantic information found in class comments. To the best of our knowledge, no previous work has provided a comprehensive taxonomy of class comment types for these three programming languages with the help of a common automated approach. Our results confirm that our approach can classify frequent class comment information types with high accuracy for the Python, Java, and Pharo programming languages. We believe this work can help to monitor and assess the quality and evolution of code comments in different programming languages, and thus support maintenance and evolution tasks.

**Keywords:** Natural Language Processing Technique, Code Comment Analysis, Software Documentation

### 1. Introduction

Software maintenance and evolution tasks require developers to perform program comprehension activities [1, 2]. To understand a software system, developers usually refer to the software documentation of the system [3, 4]. Previous studies have demonstrated that developers trust code comments more than other forms of documentation when they try to answer program

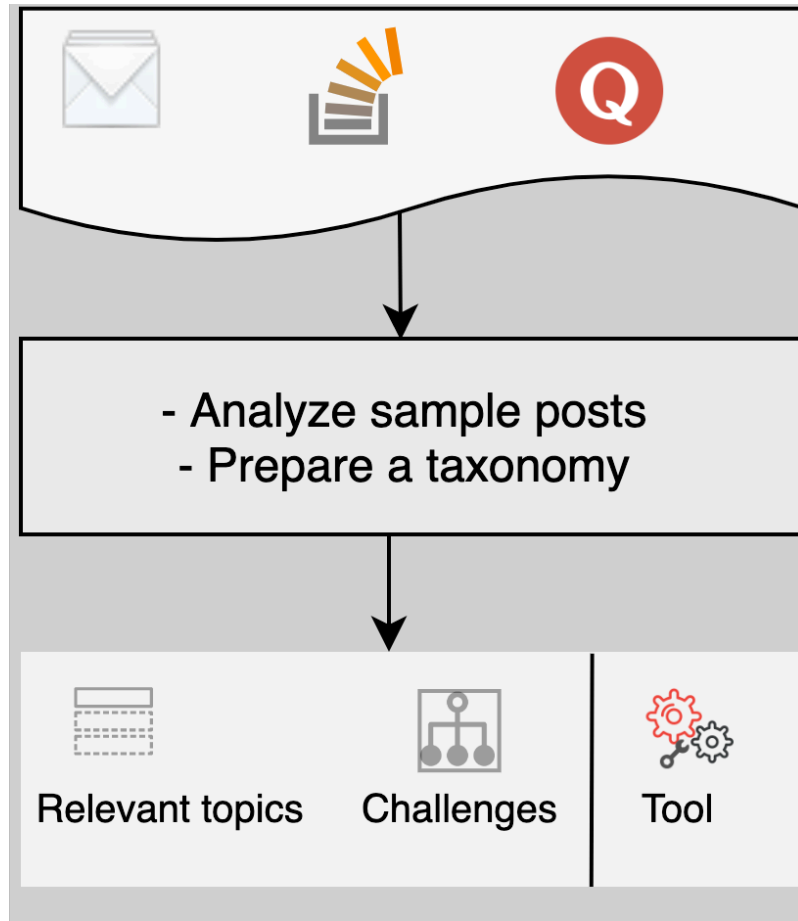
comprehension questions [5, 6, 4]. In addition, recent work has also demonstrated that “code documentation” is the most used source of information for bug fixing, implementing features, communication, and even code review [7]. In particular, well-documented code simplifies software maintenance activities, but many programmers often overlook or delay code commenting tasks [8].

Under minor revision

Journal of Systems and Software  
(JSS), 2021

RQ2: What do developers ask about comments?

## RQ2: What do developers ask about comments?



Mailing lists, Stack  
Overflow, Quora

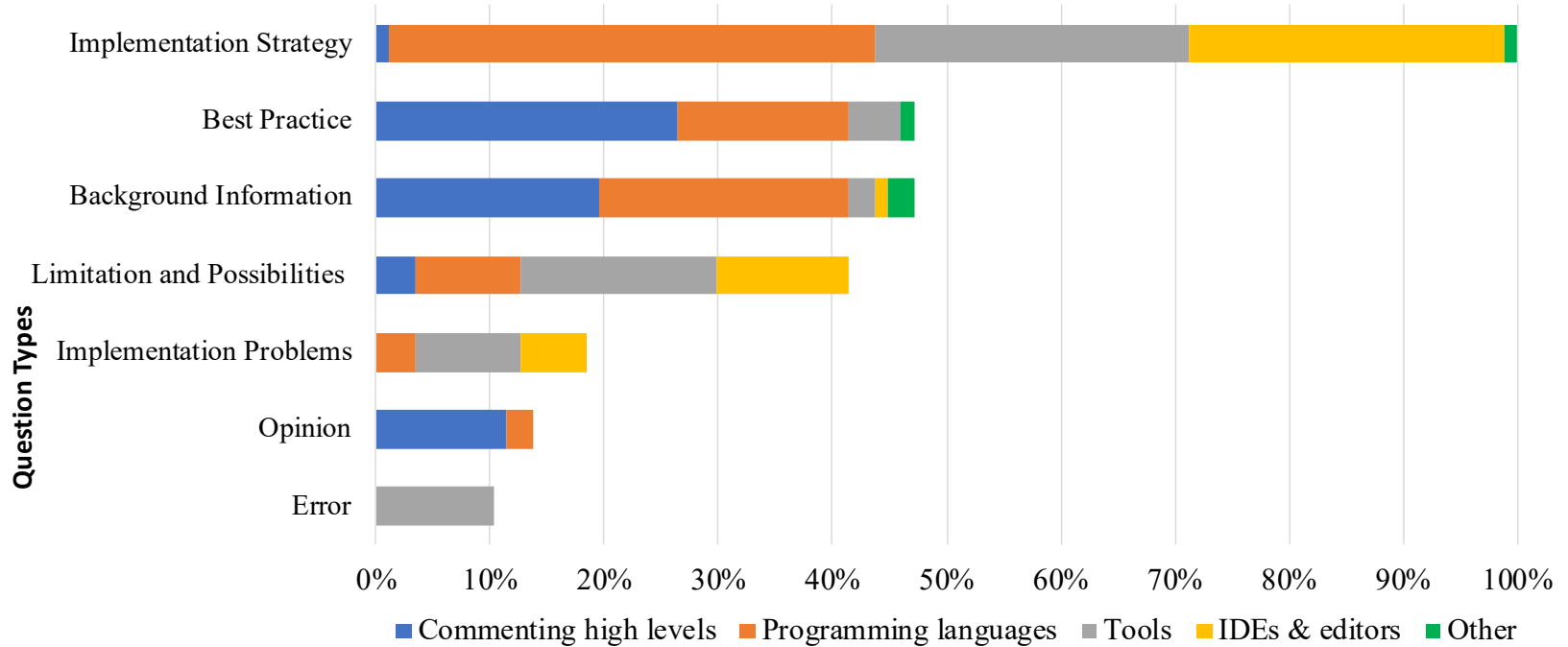
**Extract** discussion,  
**23,631 posts**

**Classify** sample discussion,  
**1,400 posts**

**Output:** a taxonomy,  
guidelines, tool



## RQ2: What do developers ask about comments?



# Takeaway

Developers seek commenting conventions to write quality comments but often find it hard to locate them



To what extent developers follow comment conventions in writing comments?

# Timeline & Contribution

## What do Developers Discuss about Comment Conventions?

An investigation of developer concerns related to comments on various online platforms

## Makar: A Framework for Multi-source Studies based on Unstructured Data

A tool to conduct study on multiple sources such as emails, stack overflow, GitHub

In the stage of acceptance

In the stage of submission

### What do Developers Discuss about Code Comment Conventions?

BLINDED

**Abstract**—Code comments are important for program comprehension, development, and maintenance tasks. Given the unstructured or semi-structured nature, and varying standards of code comments, developers get easily confused (especially novice developers) about which convention(s) to follow while writing code documentation. Thus, they post related questions on external online sources to seek better commenting practices. In this paper, we analyze comment convention discussions on Stack Overflow, Quora, and mailing lists, to shed some light on the questions developers ask about commenting convention practices.

To achieve this goal, we use two approaches. The first semi-automated approach uses Latent Dirichlet Allocation (LDA) to identify emerging topics concerning code comment conventions. In the second approach we manually analyze a statistically significant sample set of posts extracted from the aforementioned sources to derive a taxonomy that provides an overview of the main developer questions about commenting conventions. Our results highlight that nearly 40% of the sampled SO questions specifically mention *implementation strategies* for writing comments in documentation tools and environments. On the other hand, on Quora developer questions focus on best practices (or provide opinions) for writing code comments. On mailing lists, we find that developers do not discuss commenting conventions. Hence, in evaluating our approaches we found that (1) to learn best comment convention practices, developers refer to Quora often, while to know how to best implement comment conventions developers rely on SO; (2) future mining studies should prioritize the usage of specific sources to mine best comment convention practices and implementation strategies.

**Index Terms**—Mining developer sources, Stack Overflow, Quora, Mailing lists, Comment Analysis, Software documentation

by a programming language's grammar nor checked by its compiler. Consequently, developers follow various comment conventions in writing code comments [8]. Therefore, writing good comments and maintaining them in projects is a major developer concern [5], [9].

To resolve potential confusion, and to learn best practices, developers post questions on various Q&A forums. Stack Overflow (SO) is one of the most popular Q&A forums, enabling developers to ask questions to experts and other developers.<sup>2</sup> Similarly, Quora<sup>3</sup> and mailing lists are widely adopted by developers to discuss development and documentation aspects [10], [11]. Previous works have utilized these sources to understand developer needs and challenges about various software aspects [6], [12], [11]. However, no prior work has investigated the complementarity and usefulness of various online sources in supporting specific comment convention discussions.

The goal of this discussion on SO, concerns specific goal, we formulate

- 1) *RQ<sub>1</sub>*: What commenting
- 2) *RQ<sub>2</sub>*: What *opers discuss*

### In preparation Source Code Analysis and Manipulation (SCAM), 2021

### Makar: A Framework for Multi-source Studies based on Unstructured Data

Mathias Birrer\*, Pooja Rani<sup>†</sup>, Sebastiano Panichella<sup>†</sup>, Oscar Nierstrasz\*

\*Software Composition Group, University of Bern, Bern, Switzerland

✉ scg.unibe.ch/staff

<sup>†</sup>Zürich University of Applied Science (ZHAW) panic@zhaw.ch

**Abstract**—To perform various development and maintenance tasks, developers frequently seek information on various sources such as mailing lists, Stack Overflow (SO), and Quora. Researchers analyze these sources to understand developer information needs in these tasks. However, extracting and preprocessing unstructured data from various sources, building and maintaining a reusable dataset is often a time-consuming and iterative process. Additionally, the lack of tools for automating this data analysis process complicates the task to reproduce previous results or datasets.

To address these concerns we propose *Makar*, which provides various data extraction and preprocessing methods to support researchers in conducting reproducible multi-source studies. To evaluate *Makar*, we conduct a case study that analyzes code comment related discussions from SO, Quora, and mailing lists. Our results show that *Makar* is helpful for preparing reproducible datasets from multiple sources with little effort, and for identifying the relevant data to answer specific research questions in a shorter time compared to the manual investigation, which is of critical importance for studies based on unstructured data. Tool webpage: <https://github.com/macthub/makar>

**Index Terms**—Mining developer sources, Code comments, Stack Overflow, Mailing lists



Fig. 1. Developers seek various sources during software development

to lack of automated techniques pose various challenges in conducting reproducible studies [4], [5], [3]. To gain a deeper understanding of these challenges, we surveyed the literature that focuses on studying developers information needs from different external sources (see section II).

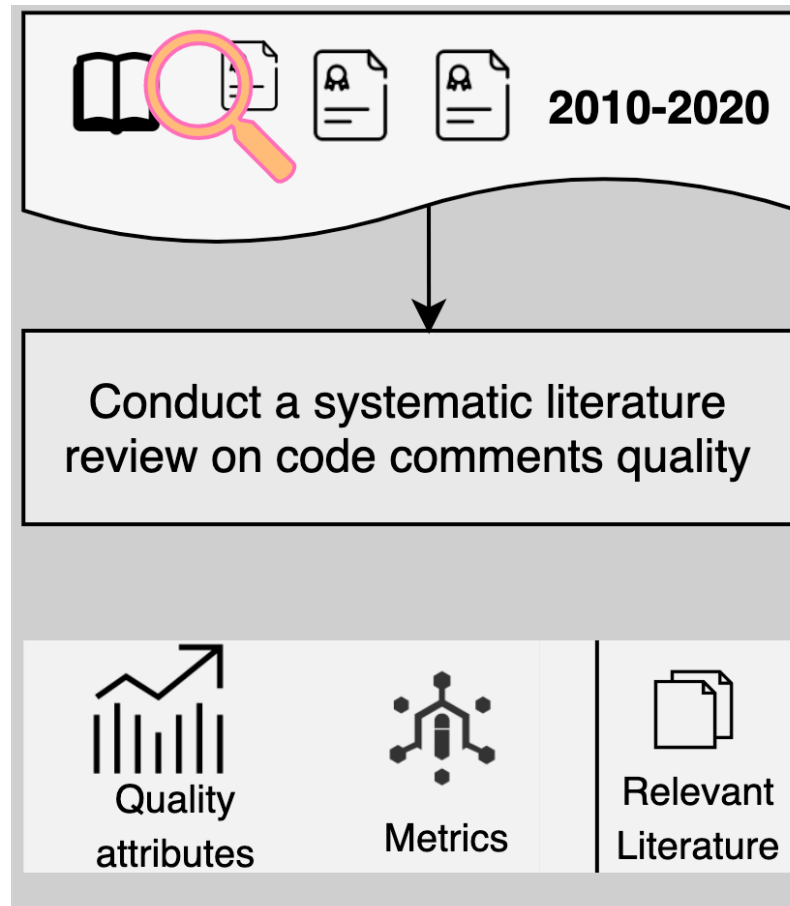
Prior works have raised and identified the crucial factors affecting the reproducibility of the mining studies such as data retrieval methodology, data processing steps, or dataset availability [6], [5], [4]. Chen *et al.* pointed out that 50% of articles do not report whether word stemming, a common text

Accepted

Software Analysis, Evolution and Reengineering (SANER), 2021

RQ3: What quality attributes are often considered in assessing comment quality?

# RQ3: What quality attributes are often considered in assessing comment quality?



SE conferences & journals  
**195 venues**

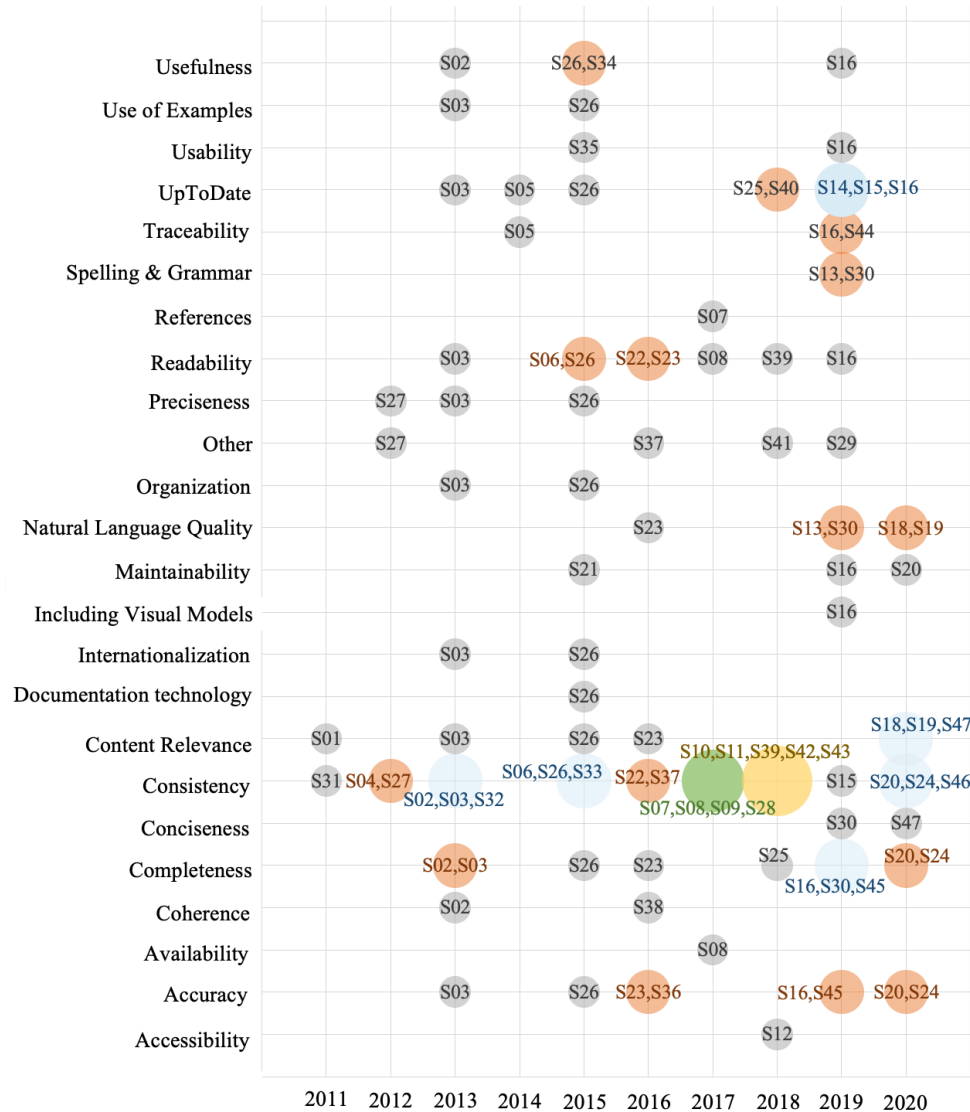
Extract proceedings,  
**332 proceedings**

Identify candidate papers  
**2,353 papers**

Select relevant papers,  
**47 papers**

**Output:** quality attributes,  
metrics

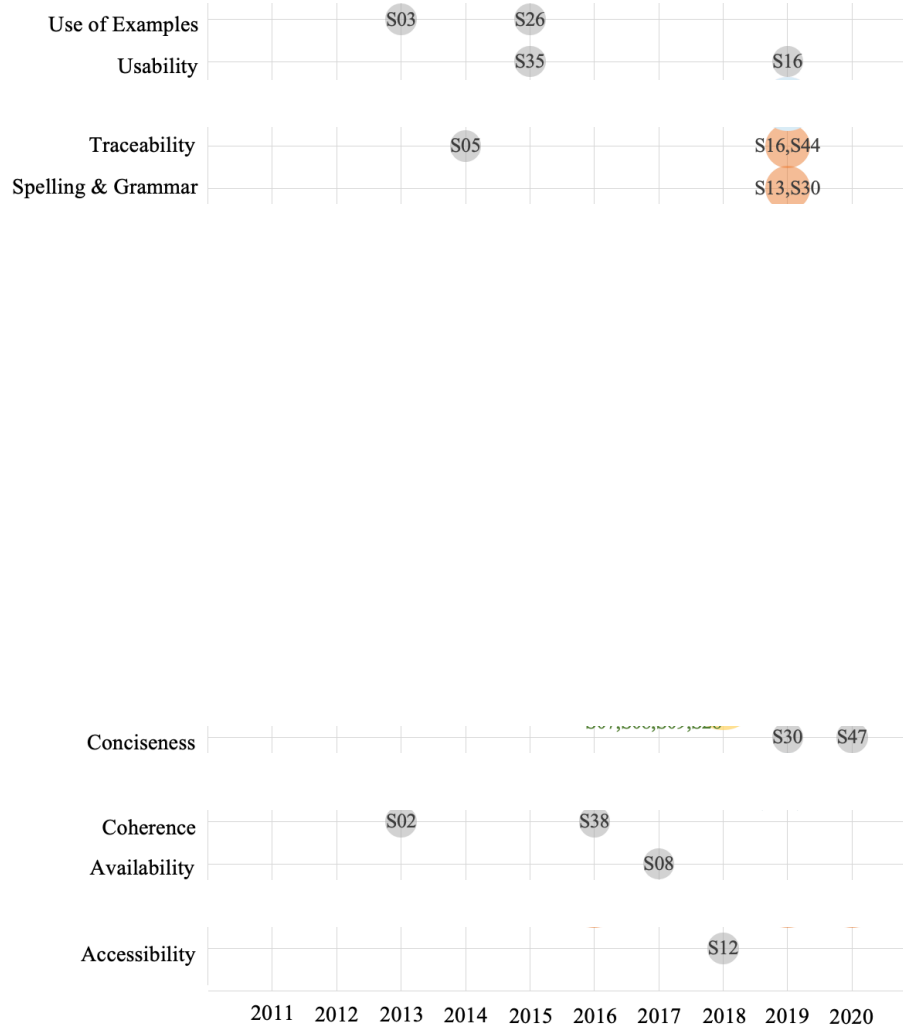
# RQ3: What quality attributes are often considered in assessing comment quality?



# RQ3: What quality attributes are often considered in assessing comment quality?



# RQ3: What quality attributes are often considered in assessing comment quality?





# Takeaway



Researchers are focusing often on handful of quality attributes in assessing comment quality

Numerous studies suggest the benefits of usability, or accessibility of comments. We need to focus on assessing comments w.r.t them.

# Timeline & Contribution



## *A SLR about comments*

An investigation of quality attributes, and techniques used for comment assessment

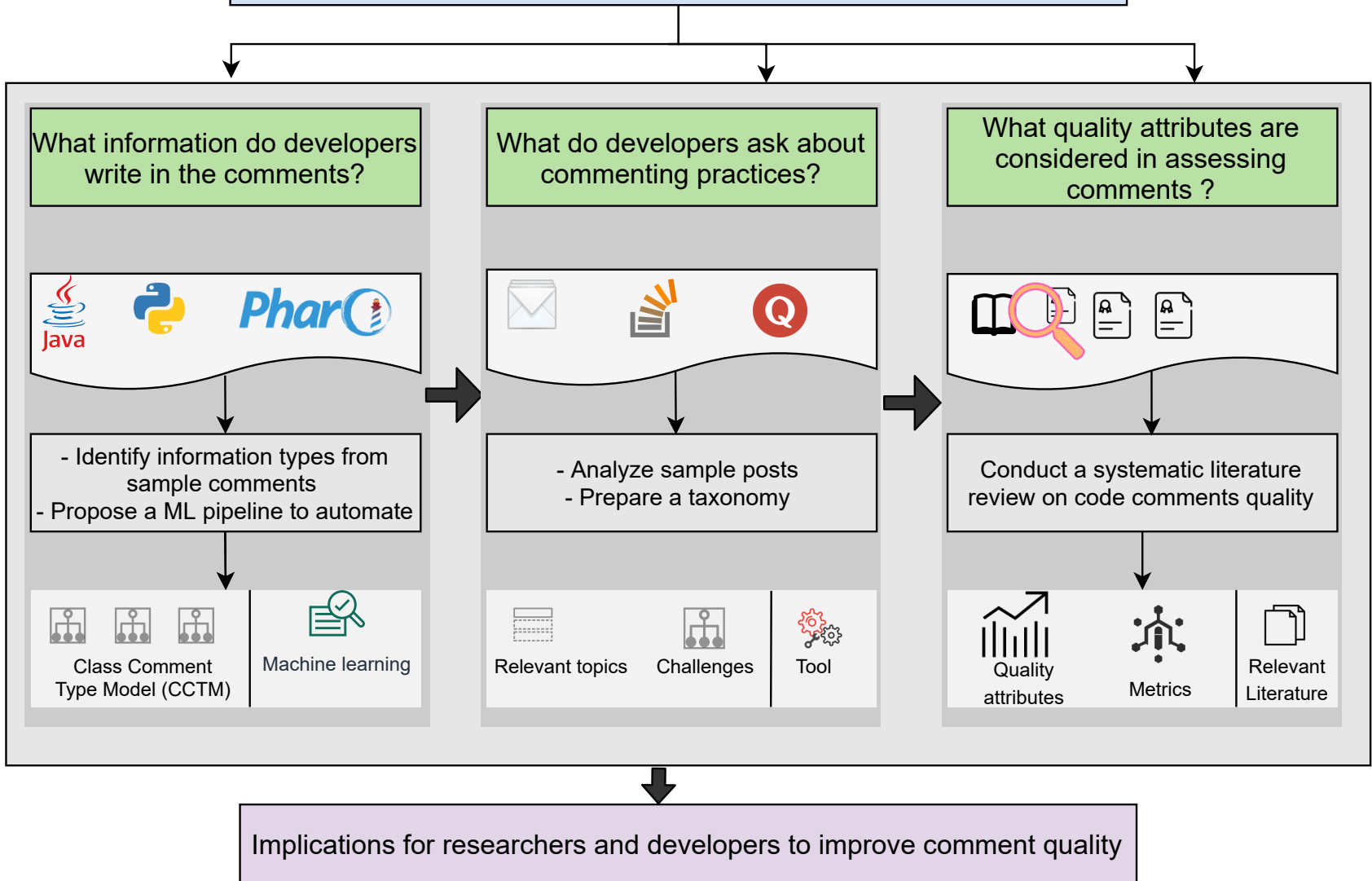
Results are submitted to a software engineering venue

## *Do Developers Follow Comment Conventions?*

An investigation of comment adherence to comment conventions

In preparation

# Speculative Analysis of Code Comments



The ultimate goal of automatically assessing comments is still far away...

# Future work

Which quality attributes do developers find important?

Which information types do developers find important?

How do various information types support developers?

An IDE plugin to support automatic assessment of comments

# Takeaway

Developer seek and follow comment conventions.

Having support to assess various aspects of comments can help maintaining high-quality comments



<https://twitter.com/poojaruhal>



<http://scg.unibe.ch/staff/Pooja-Rani>