

Evaluating Code Summarization Techniques: A New Metric and an Empirical Characterization

Antonio Mastropaolo
16/11/2023

SEART

antoniomastropaolo.com 

AntonioMastro2 



**Documenting the code is
extremely important.....**



**Documenting the code is
extremely important.....**

SUC**



**Commented code is
more **understandable**
as compared to
non-commented code**

The Effect of Modularization and Comments on Program Comprehension

S. N. Woodfield

Computer Sciences Department
Arizona State University
Tempe, Arizona 85281

H. E. Dunsmore
V. Y. Shen

Department of Computer Sciences
Purdue University
West Lafayette, Indiana 47907

Abstract

An experiment was conducted to investigate how different types of modularization and comments are related to programmers' ability to understand programs. Forty-eight experienced programmers were given eight different versions of the same program and asked to answer a twenty question quiz used to measure comprehension. These eight different versions were the result of the program being constructed with four types of modularization (monolithic, functional, super, and abstract data type), each with and without comments. Those subjects whose programs contained comments were able to answer more questions than those without comments. Also, those subjects who were given the abstract data type version of the program were able to do significantly better than those with any other type of modularization.

Introduction

One software feature common to almost all large programs is the use of modularization. At first, modularization was motivated by the desire to conserve computer memory. A program was considered properly modularized if it was possible to execute the program using the available memory. Later, as memory became larger, code was modularized to make it easier for people to work with programs. Now, considering the human factors of programming, it has become more difficult to determine or describe proper modularization. Today almost everyone agrees that modularity is an important aspect of reducing the cost of developing software systems. Yet, there remain many questions such as "How large should a module should be?" or "How should modules interface with each other?".

This paper reports on an empirical study of several modularization methodologies to see how different approaches affect the ability of programmers to understand programs. The Yourdon and Constantine [1] methodology attempts to divide a program into functional units, one function per physical module. (A physical module is a section of code whose boundaries are defined by the syntax of the language such as a subroutine in Fortran or a procedure in Pascal). A routine implementing only one function should be easy to understand since all the code is related to that one function. Furthermore, the function (module) is

constructed only once and is invoked in other parts of the program where that function is needed. Since only one function resides in a physical module, the identification of the function's boundaries is simple.

The abstract data type approach to modularization may include several functions in one physical module [2,3]. The module brings together all operations that manipulate a single data type. This abstract data type approach is a special case of physical modules containing several logical modules. (A logical module is a section of code which implements only one function. A logical module need not correspond one-to-one with physical modules).

This paper reports on research conducted to help determine how best to modularize programs to aid in comprehension. As described above, possibilities include physical modules containing exactly one logical module or more than one logical module. It is possible that more than one logical module in a physical module is acceptable if the logical modules are easily identified and easily understood.

To facilitate understanding one means of separating the logical modules in a physical module is to insert a brief comment at the beginning of each logical module. Comments used in this research were designed to describe a logical module's contents. If there were more than one logical module per physical module, the inserted comments were intended to help define the boundaries of the logical modules. Previous work has shown that comment statements in a program are not harmful, but may not be very useful either [4,5]. Such studies were conducted primarily to decide if the prolific use of comments aids in the general comprehension of a program. Our research attempts to determine whether comments used more judiciously (i.e., in the right locations) might be a better approach. Specifically, this research attempts to determine if short comments, inserted just before a logical module, can aid comprehension by briefly describing the function of the logical module and helping to define the logical module's boundaries.

We considered four types of modularization:

- (1) monolithic - (i.e., non-modularization) a program with one physical module.

Quality Analysis of Source Code Comments

Daniela Steidl Benjamin Hummel Elmar Juergens
CQSE GmbH, Garching b. München, Germany
{steidl,hummel,juergens}@cqse.eu

ICPC'13

Abstract—A significant amount of source code in software systems consists of comments, *i. e.*, parts of the code which are ignored by the compiler. Comments in code represent a main source for system documentation and are hence key for source code understanding with respect to development and maintenance. Although many software developers consider comments to be crucial for program understanding, existing approaches for software quality analysis ignore system commenting or make only quantitative claims. Hence, current quality analyzes do not take a significant part of the software into account. In this work, we present a first detailed approach for quality analysis and assessment of code comments. The approach provides a model for comment quality which is based on different comment categories. To categorize comments, we use machine learning on Java and C/C++ programs. The model comprises different quality aspects: by providing metrics tailored to suit specific categories, we show how quality aspects of the model can be assessed. The validity of the metrics is evaluated with a survey among 16 experienced software developers, a case study demonstrates the relevance of the metrics in practice.

I. INTRODUCTION

A significant amount of source code in software systems consists of comments, which document the implementation and help developers to understand the code, *e. g.*, for later modification or reuse: Several researchers have conducted experiments showing that commented code is easier to understand than code without comments [1], [2]. Comments are the second most-used documentary artifact for code understanding, behind only the code itself [3]. In addition, source code documentation is also vital in maintenance and forms an important part of the general documentation of a system. In contrast to external documentation, comments in source code are a convenient way for developers to keep documentation and code consistently up to date. Developers widely agree that poor general documentation leads to misunderstandings [4], and studies have shown that poor documentation significantly lowers the maintainability of software [5]. Although developers commonly agree on the importance of software documentation [3], commenting code is often neglected due to release deadlines and other time pressure during development.

Previous approaches to analyzing software quality ignore comments, or make only quantitative claims: they evaluate the comment ratio of a system to measure documentation quality [6], [7]. However, this metric is not sufficient: Several comments (copyrights or commented out code) should be excluded

This work was partially funded by the German Federal Ministry of Education and Research (BMBF), grant EvoCon, 01IS12034A. The responsibility for this article lies with the authors.

as they do not enhance system understanding and quantitative measures cannot detect outdated/ useless comments.

Furthermore, a complete model of comment quality does not exist. Coding conventions, *e. g.*, marginally touch on the topic of commenting code but mostly lack depth and precision [8]. So far, (semi-) automatic methods for comment quality assessment have not been developed as comment analysis is a difficult task: Comments comprise natural language and have no mandatory format aside from syntactic delimiters. Hence, algorithmic solutions will be heuristic in nature.

Problem Statement. Current quality analysis approaches ignore system commenting or are restricted to the comment ratio metric only. Hence, a major part of source code documentation is ignored during software quality assessment.

Contribution. Based on comment classification, we provide a semi-automatic approach for quantitative and qualitative evaluation of comment quality.

We present a semi-automatic approach for comment quality analysis and assessment. First, we perform comment categorization both for Java and C/C++ programs based on machine learning to differentiate between different comment types. Comment categorization enables a detailed quantitative analysis of a system's comment ratio and a qualitative analysis tailored to suit each single category. Comment categorization is the underlying basis of our comprehensive quality model. The model comprises quality attributes for each comment category based on four criteria: consistency throughout the project, completeness of system documentation, coherence with source code, and usefulness to the reader. To assess quality attributes, we provide metrics detecting quality defects in comments of specific categories. We evaluate the metrics' validity and relevance separately: Validity is evaluated with a survey among experienced software developers. The survey shows that the metrics can additionally give refactoring recommendations. A case study of five open source projects evaluates the relevance of our approach showing that comment classification provides better insights of the system documentation quality than the simple comment ratio metric and that our metrics reveal quality defects in practice.

II. RELATED WORK

We group the related work into categories of general comment analysis, information retrieval techniques, comment evolution studies, and code recognition approaches.

Developers do not always carefully
comment code, or **update** existing
comments in response to **code**
changes

Do Code and Comments Co-Evolve? On the Relation Between Source Code and Comment Changes

Beat Fluri, Michael Würsch, and Harald C. Gall
s.e.a.l. – software architecture and evolution lab
Department of Informatics
University of Zurich, Switzerland
{fluri,wuersch,gall}@ifi.uzh.ch

Abstract

Comments are valuable especially for program understanding and maintenance, but do developers comment their code? To which extent do they add comments or adapt them when they evolve the code? We examine the question whether source code and associated comments are really changed together along the evolutionary history of a software system. In this paper, we describe an approach to map code and comments to observe their co-evolution over multiple versions. We investigated three open source systems (i.e., ArgoUML, Azureus, and JDT Core) and describe how comments and code co-evolved over time. Some of our findings show that: 1) newly added code—despite its growth rate—barely gets commented; 2) class and method declarations are commented most frequently but far less, for example, method calls; and 3) that 97% of comment changes are done in the same revision as the associated source code change.

1. Introduction

“Comment your code!” The task of commenting one’s source code is often neglected; even though everybody who is writing software knows the value of good comments [20]. Reading code is a fundamental task during software engineering [10]—and code is read more often than it is written. Even books covering best-practices in commenting exist, e.g., *The Elements of Java Style* by Vermeulen *et al.* [21]. Comments allow one to understand the code faster and deeper and to improve its readability [18, 19]. Especially, they are crucial to sustain software maintainability and aid in reverse engineering, for example when applying the *Read All the Code in One Hour* reengineering pattern [5]. Elshoff and Marcotty already stated in the early eighties that comments as well as the structure of the source

code aid in program understanding and therefore reduce maintenance costs [6]. This finding was confirmed by the studies of Ted Tenny [19]. But as the example of Lakhoria shows, sometimes programmers do not care that someone else might want to understand the source code [13].

For maintenance and reverse engineering tasks both, the lack of comments as well as outdated comments are counter-productive. To understand whether the comments are a reason for decreasing maintainability in software projects we address the following research questions in this paper:

1. Does the ratio between comment and source code remain stable over the history of a software project (i.e., is there a recognizable effort made to comment the source code)?
2. Which source code entities are most likely to be commented?
3. Are comments adapted when source code is changed (i.e., are comments kept up-to-date) and when does the adaptation take place—while changing the source code or afterwards?

The contribution of this paper is an approach to map source code entities to comments in the code and a technique to extract comment changes over the history of a software project. We have conducted three experiments on three different open source software systems to answer our research questions: ArgoUML, Azureus, and JDT Core. We describe how comments and code co-evolved over time. Some of our findings show that: 1) newly added code—despite its growth rate—barely gets commented; 2) class and method declarations are commented most frequently but far less, for example, method calls; and 3) that 97% of comment changes are done in the same revision as the associated source code change.

The remainder of the paper is structured as follows. In Section 2 we present our approach to investigate the relation

A Study of the Documentation Essential to Software Maintenance

Sergio Cozzetti B. de Souza
Universidade Católica de Brasília
Brasília, DF, BRAZIL

Nicolas Anquetil
Universidade Católica de Brasília
Brasília, DF, BRAZIL
anquetil@ucb.br

Káthia M. de Oliveira
Universidade Católica de Brasília
Brasília, DF, BRAZIL
kathia@ucb.br

ABSTRACT

Software engineering has been striving for years to improve the practice of software development and maintenance. Documentation has long been prominent on the list of recommended practices to improve development and help maintenance. Recently however, agile methods started to shake this view, arguing that the goal of the game is to produce software and that documentation is only useful as long as it helps to reach this goal.

On the other hand, in the re-engineering field, people wish they could re-document useful legacy software so that they may continue maintain them or migrate them to new platform.

In these two case, a crucial question arises: "How much documentation is enough?" In this article, we present the results of a survey of software maintainers to try to establish what documentation artifacts are the most useful to them.

Categories and Subject Descriptors

D.2.0 [Software Engineering]: General; D.2.7 [Software Engineering]: Distribution, Maintenance, and Enhancement—Documentation

General Terms

software system documentation, empirical study, software maintenance, program understanding

1. INTRODUCTION

Among all the recommended practices in software engineering, software documentation has a special place. It is one of the oldest recommended practices and yet has been, and continue to be, renowned for its absence (e.g. [4]). There is no end to the stories of software systems (particularly legacy software) lacking documentation or with outdated documentation. For years, the importance of documentation has been stressed by educators, processes, quality

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee.
SIGDOC'05, September 21–23, 2005, Coventry, United Kingdom.
Copyright 2005 ACM 1-59593-175-9/05/0009 ...\$5.00.

models, etc. and despite of this we are still discussing why it is not generally created and maintained (e.g. [11]). The topic gained renewed interest with two recent trends:

- Agile methods question the importance of documentation as a development aid;
- The growing gap between "traditional" (e.g. COBOL) and up-to-date technologies (e.g. OO or web-oriented) increased the pressure to re-document legacy software.

Both issues raise a similar question: What documentation would be most useful to software maintenance? If they propose a renewed development paradigm, agile methods do not bring significant changes for software maintenance. They do claim that permanent re-factoring turns maintenance into a normal state of the methods. However, they do not explain how such methods would work over extended periods of time, when a development team is sure to disperse with the knowledge it has of the implementation details. Documentation is still a highly relevant artifact of software maintenance.

Legacy software re-documentation tries to remedy the deficiencies of the past in terms of up-to-date documentation. However, it is a costly activity, difficult to justify to users because it does not bring any visible change for them (at least in the short term).

In this paper we present a survey of software maintainers trying to establish the importance of various documentation artifacts for maintenance. The paper is divided as follows: In Section 2, we review some basic facts about software maintenance and its needs; in Section 3, we summarize the relevant literature on software documentation; in Section 4, we present the survey we conducted; and in Section 5, we comment the result of the survey.

2. SOFTWARE MAINTENANCE

Maintenance is traditionally defined as any modification made on a system after its delivery. Studies show that software maintenance is, by far, the predominant activity in software engineering (90% of the total cost of a typical software [17, 21]). It is needed to keep software systems up-to-date and useful. Any software system reflects the world within which it operates, when this world changes, the software needs to change accordingly. Lehman's first law of software evolution (law of continuing change, [13]) is that "a program that is used undergoes continual change or becomes progressively less useful". Maintenance is mandatory,

A Study of the Document

Sergio Cozzetti B. de Souza
Universidade Católica de Brasília
Brasília, DF, BRAZIL

ABSTRACT

Software engineering has been the practice of software development. The practice of software development has long been documented. Recently however, new practices to improve software maintenance. Recently however, new practices to improve software maintenance. Recently however, new practices to improve software maintenance.

Analyzing the co-evolution of comments and source code

Beat Fluri · Michael Würsch · Emanuel Giger · Harald C. Gall

Published online: 26 March 2009
© Springer Science+Business Media, LLC 2009

Abstract Source code comments are a valuable instrument to preserve design decisions and to communicate the intent of the code to programmers and maintainers. Nevertheless, commenting source code and keeping comments up-to-date is often neglected for reasons of time or programmers' obliviousness. In this paper, we investigate the question whether developers comment their code and to what extent they add comments with source code they evolve the code. We present an approach to associate comments with source code entities to track their co-evolution over multiple versions. A set of heuristics are used to decide whether a comment is associated with its preceding or its succeeding source code entity. We analyzed the co-evolution of code and comments in eight different open source and closed source software systems. We found with statistical significance that (1) the relative amount of comments and source code grows at about the same rate; (2) the type of a source code entity, such as a method declaration or an if-statement, has a significant influence on whether or not it gets commented; (3) in six out of the eight systems, code and comments co-evolve in 90% of the cases; and (4) surprisingly, API changes and comments do not co-evolve but they are re-documented in a later revision. As a result, our approach enables a quantitative assessment of the commenting process in a software system. We can therefore, leverage the results to provide feedback during development to increase the awareness of when to add comments or when to adapt comments because of source code changes.

Keywords Software evolution analysis · Software repositories · Source code changes · Comment changes · Comment quality · Software maintenance

1 Introduction

Documenting software is painful, especially when time is scarce and release deadlines are putting serious pressure on development teams, making it necessary to prioritize tasks.

B. Fluri (✉) · M. Würsch · E. Giger · H. C. Gall
Department of Informatics, University of Zurich, Zurich, Switzerland
e-mail: fluri@ifi.uzh.ch

A Large-Scale Empirical Study on Code-Comment Inconsistencies

Fengcai Wen, Csaba Nagy, Gabriele Bavota, Michele Lanza
Software Institute, Università della Svizzera italiana (USI), Switzerland
 {fengcai.wen, csaba.nagy, gabriele.bavota, michele.lanza}@usi.ch

Abstract—Code comments are a primary means to document source code. Keeping comments up-to-date during code change activities requires substantial time and attention. For this reason, researchers have proposed methods to detect code-comment inconsistencies (*i.e.*, comments that are not kept in sync with the code they document) and studies have been conducted to investigate this phenomenon. However, these studies were performed at a small scale, relying on quantitative analysis, thus limiting the empirical knowledge about code-comment inconsistencies. We present the largest study at date investigating how code and comments co-evolve. The study has been performed by mining 1.3 Billion AST-level changes from the complete history of 1,500 systems. Moreover, we manually analyzed 500 commits to define a taxonomy of code-comment inconsistencies fixed by developers. Our analysis discloses the extent to which different types of code changes (*e.g.*, change of *selection* statements) trigger updates to the related comments, identifying cases in which code-comment inconsistencies are more likely to be introduced. The defined taxonomy categorizes the types of inconsistencies fixed by developers. Our results can guide the development of tools aimed at detecting and fixing code-comment inconsistencies.

Index Terms—Software Evolution, Code Comments

I. INTRODUCTION

Any code-related activity lays its foundations in program comprehension: before fixing a bug, refactoring a class, or writing new tests, developers first need to acquire knowledge about the involved code components. As recently shown by Xia *et al.* [1], this results in 58% of developers' time spent comprehending code. Besides the code itself, code comments are considered as the most important form of documentation for program comprehension [2]. Indeed, not surprisingly, studies showed that commented code is easier to comprehend than uncommented code [3], [4]. This empirical evidence also pushed researchers to consider code comments as a pivotal factor to study technical debt [5]–[7], or to assess code quality [8], [9].

While the importance of code comments is undisputed, developers do not always have the chance to carefully comment new code and/or to update comments as consequence of code changes [10]. This latter scenario might result in the introduction of code-comment inconsistencies, manifesting when the source code does not co-evolve with the related comments. For example, if a method comment is not updated after major changes to the method's application logic, the comment might provide misleading information to developers comprehending the method, hindering program comprehension rather than fostering it.

Given the potential harmfulness of code-comment inconsistencies, several researchers studied the co-evolution of code and comments [11]–[14], while others proposed techniques and tools able to detect code-comment inconsistencies automatically [15]–[18]. These techniques are able to identify specific types of code-comment inconsistencies. For example, @TCOMMENT [17] detects inconsistencies between Javadoc comments related to null values and exceptions with the behavior implemented in the related method's body, while Fraco [18] focuses on inconsistencies introduced as result of rename refactoring operations. Still, more research is needed in this area to increase the types of code-comment inconsistencies that can be automatically identified. Also, the empirical evidence provided by studies that pioneered the investigation of code-comment evolution [11]–[14] is limited to the analysis of the change history of a few software systems (less than 10).

To raise the knowledge about the co-evolution of code and comments and the introduction/fixing of code-comment inconsistencies, we present a large-scale empirical study quantitatively and qualitatively analyzing these phenomena. We mine the complete change history of 1,500 Java projects hosted on GitHub for a total of 3,323,198 analyzed commits. For each commit, we use GUMTREEDIFF [19] to extract AST operations performed on the files modified in it. In this way, we captured fine-grained changes performed in code (*e.g.*, change of a *selection* statement) as well as update, delete, and insert operations performed in the related comments. Overall, this process resulted in a database of ~476 GB containing ~1.3 Billion AST-level operations impacting code or comments. Using this data, we study the extent to which code changes impacting different code constructs (*e.g.*, *literals*, *iteration statements*) trigger the update of the related code comments (*e.g.*, the developer adds a try statement and updates the method comment to “document” the changed code behavior).

Then, we manually analyze 500 commits identified, via a keywords-matching mechanism, as likely related to the fixing of code-comment inconsistencies. The output of this analysis is a taxonomy of code comment-related changes implemented by developers, from which we present relevant cases related to code-comment inconsistencies, and discuss implications for researchers and practitioners.

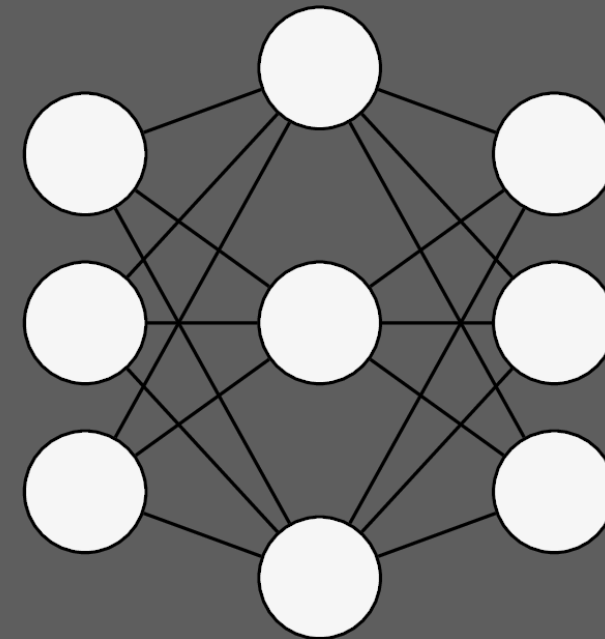
As a contribution to the research community, we make the database of fine-grained code changes publicly available. This enables the replication of this work, making also other types of investigations possible.

Code Summarizer

INPUT

```
public ConnectionConsumer createConnectionConsumer  
    (final Destination destination)  
    throws JMSEException{  
  
    if(LOGGER.isTraceEnabled())  
    {  
        ActiveMQRALogger.LOGGER.  
            trace("Create connectionConsumer");  
    }  
    else {  
        throw new IllegalStateException(ISE);  
    }  
}
```

DL-Model



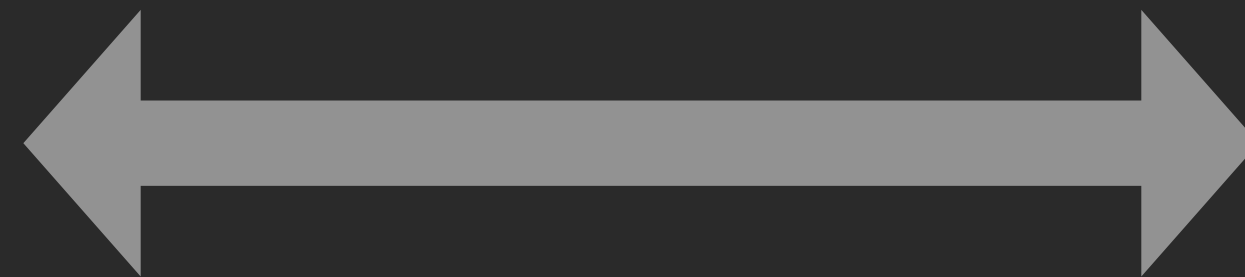
PREDICTION

**Create a connection to
the Consumer.**

MODEL EVALUATION

GROUND TRUTH

**Create a connection
to the consumer.**



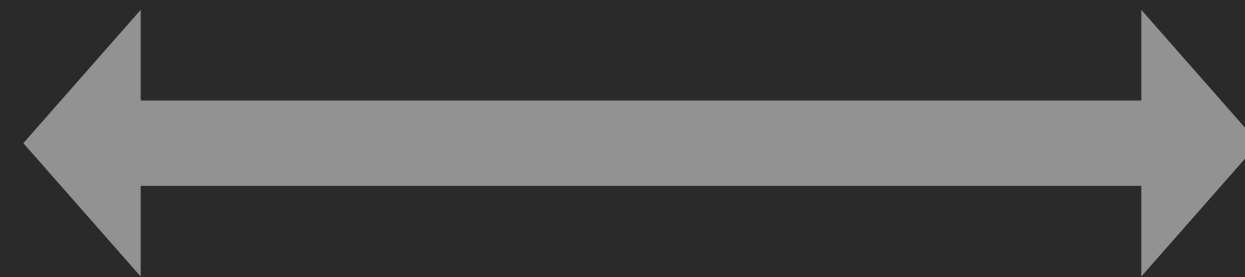
MODEL'S PREDICTION

**Create a connection
to the consumer.**

MODEL EVALUATION

GROUND TRUTH

**Create a connection
to the consumer.**



MODEL'S PREDICTION

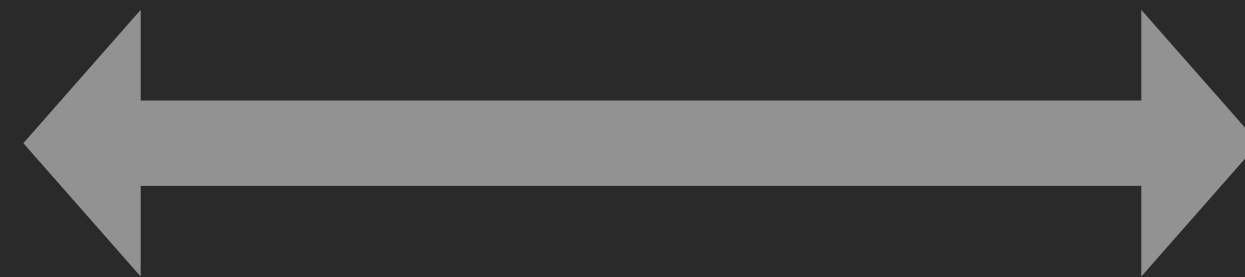
**Create a connection
to the consumer.**



MODEL EVALUATION

GROUND TRUTH

**Create a connection
to the consumer.**



MODEL'S PREDICTION

Create a logger .



NLP (Natural Language Processing) Related Metrics

BLEU SCORE

ROUGE SCORE

BERT SCORE³

...

1. BLEU: A Method for Automatic Evaluation of Machine Translation —*Papinev et al.*—

2. Rouge: A package for automatic evaluation of summaries —*Lin.*—

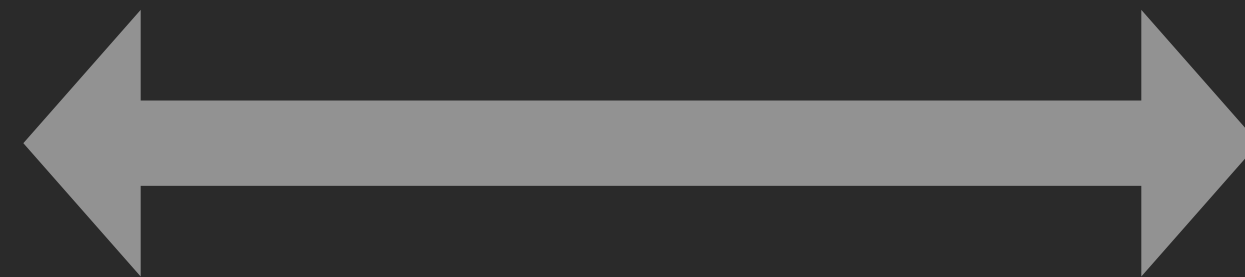
3. Bertscore: Evaluating text generation with bert —*Zhang et al.*—



MODEL EVALUATION

GROUND TRUTH

Create a connection
to the consumer.



MODEL'S PREDICTION

Create a logger .

BLEU SCORE: 0.26



Reads the contents of this source as a string.

PR

Get the textual information from this source and represent it as a string.

GT

```
public String read() throws IOException {
    Closer closer = Closer.create();
    try {
        Reader reader = closer.register(openStream());
        return CharStreams.toString(reader);
    } catch ( Throwable e ) {
        throw closer.rethrow (e);
    } finally { closer.close(); }
}
```

Reads the contents of this source as a string.

PR

Get the textual information from this source and represent it as a string.

GT

```
public String read() throws IOException {  
    Closer closer = Closer.create();  
    try {  
        Reader reader = closer.register(openStream());  
        return CharStreams.toString(reader);  
    } catch ( Throwable e ) {  
        throw closer.rethrow (e);  
    } finally { closer.close(); }  
}
```

SEMANTICALLY
EQUIVALENT
CODE SUMMARIES

Reads the contents of this source as a string.

PR

Get the textual information from this source and represent it as a string.

GT

```
public String read() throws IOException {
    Closer closer = Closer.create();
    try {
        Reader reader = closer.register(openStream());
        return CharStreams.toString(reader);
    } catch (Throwable e) {
        throw closer.rethrow(e);
    } finally { closer.close(); }
}
```

BLEU SCORE: 0.21

**SEMANTICALLY
EQUIVALENT
CODE SUMMARIES**

Reads the contents of this source as a string.

PR

Get the textual information from this source and represent it as a string.

GT

```
public String read() throws IOException {  
    Closer closer = Closer.create();  
    try {  
        Reader reader = closer.register(openStream());  
        return CharStreams.toString(reader);  
    } catch (Throwable e) {  
        throw closer.rethrow(e);  
    } finally { closer.close(); }  
}
```

U SCORE: 0.21

SEMANTICALLY
EQUIVALENT
CODE SUMMARIES



Where

is

the

code

going?

SIDE (Summary alignment to coDe sEmantic)

Get the textual information from this source and represent it as a string.



Connect to the server and return the status

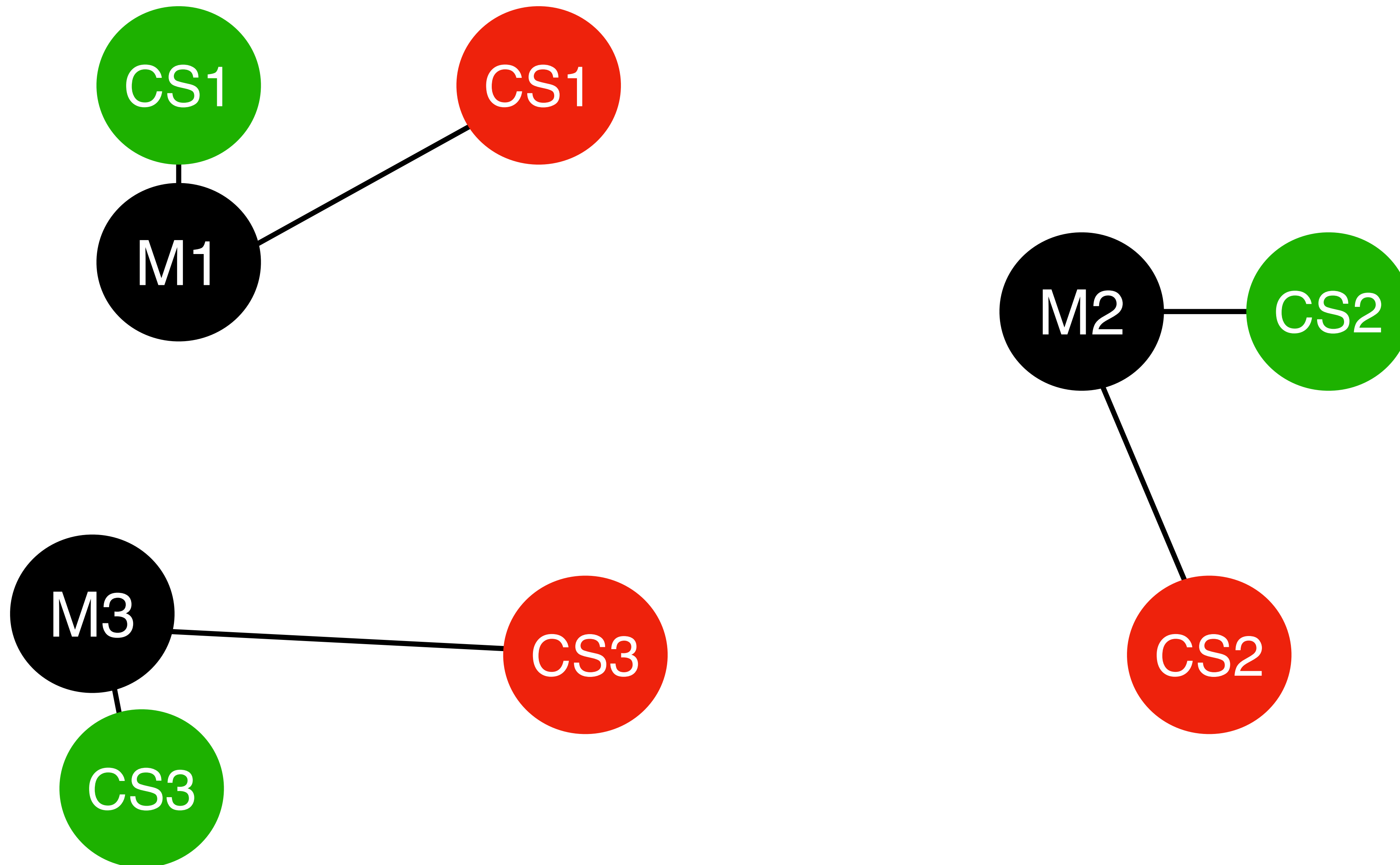


```
public String read() throws IOException {
    Closer closer = Closer.create();
    try {
        Reader reader = closer.register(openStream());
        return CharStreams.toString(reader);
    } catch (Throwable e) {
        throw closer.rethrow(e);
    } finally { closer.close(); }
}
```

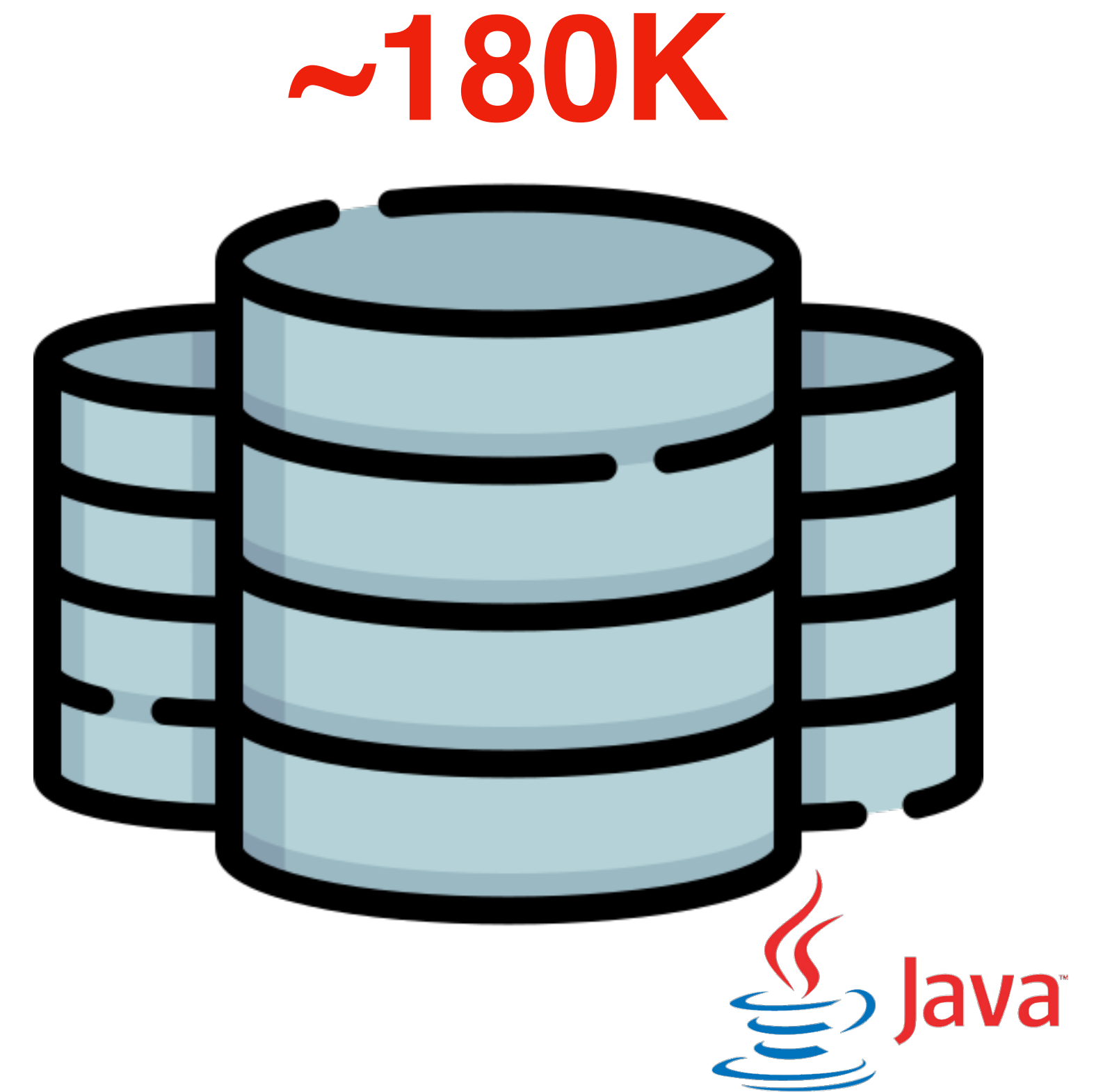
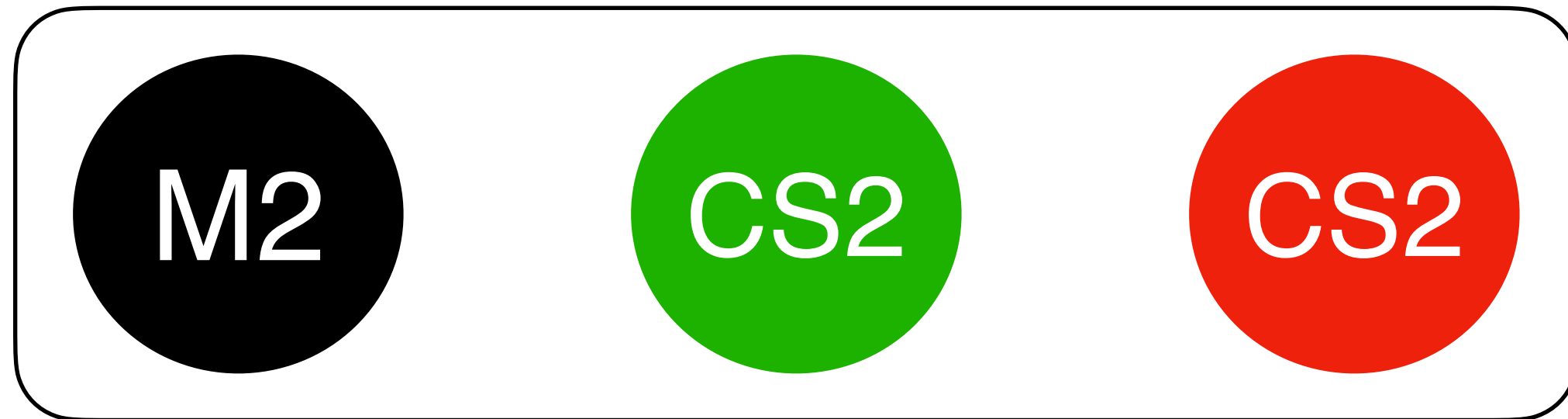
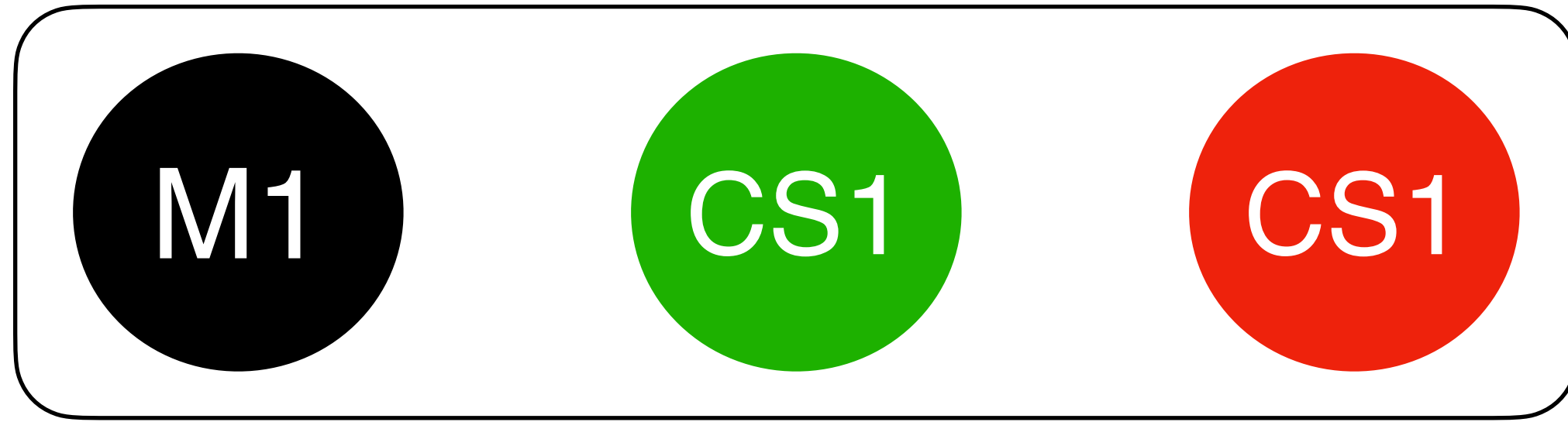
A Contrastive Learning-based Method



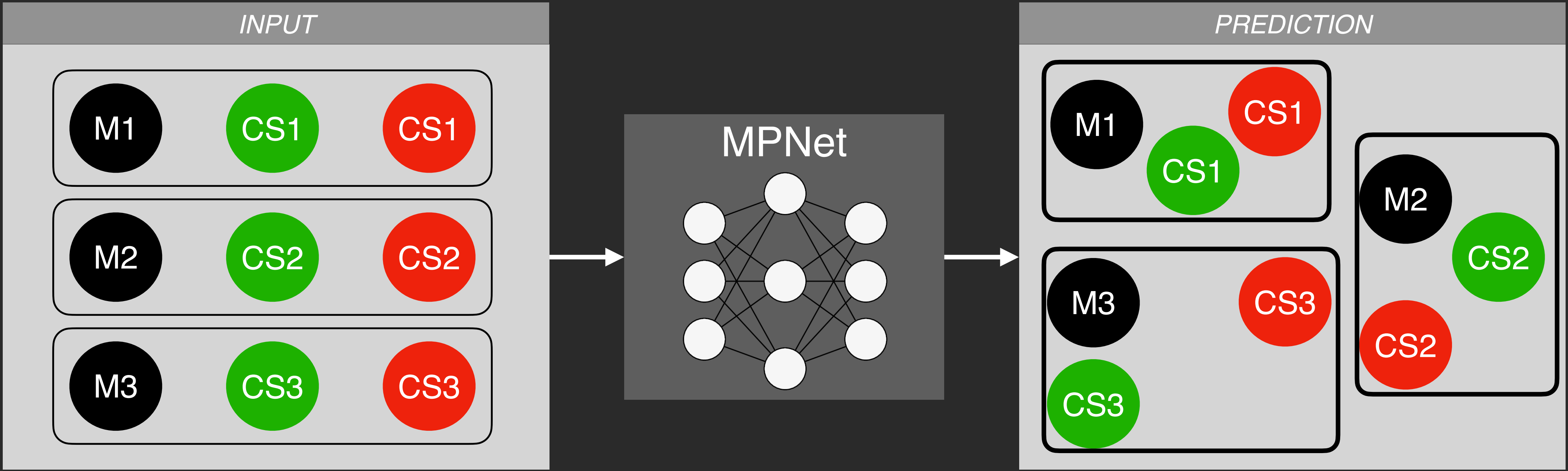
A Contrastive Learning-based Method



A Contrastive Learning-based Method



DL-Model Training



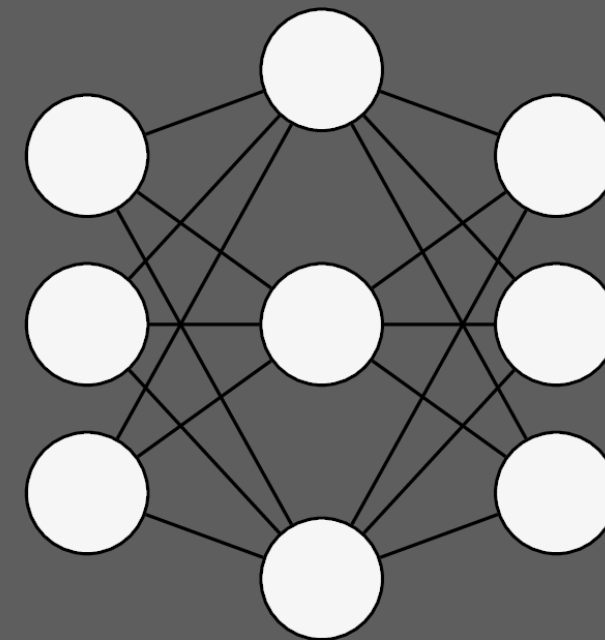
MPNet: Masked and permuted pre-training for language understanding

SIDE

INPUT

```
Create a connection to the consumer.  
public ConnectionConsumer  
  createConnectionConsumer  
  ...{  
  if(LOGGER.isTraceEnabled())  
  {  
    ActiveMQRALogger.LOGGER.  
      trace("Create  
            connectionConsumer");  
  }  
  else {  
    ...  
  }  
}
```

MPNet



PREDICTION

0.81



Evaluation of SIDE

Whether **SIDE** can be effectively used to capture new “**dimensions**” where state-of-the-art metrics fail



Evaluation of SIDE

Whether **SIDE** can be effectively used to capture new “**dimensions**” where state-of-the-art metrics fail

Whether **SIDE** allows for a more effective representations of developers’ evaluations regarding the **quality** of automatically generated summaries

Evaluation of SIDE

Reassessing Automatic Evaluation Metrics for Code Summarization Tasks

Devjeet Roy
devjeet.roy@wsu.edu
Washington State University
Pullman, WA, USA

Sarah Fakhoury
sarah.fakhoury@wsu.edu
Washington State University
Pullman, WA, USA

Venera Arnaudova
venera.arnaudova@wsu.edu
Washington State University
Pullman, WA, USA

ABSTRACT

In recent years, research in the domain of source code summarization has adopted data-driven techniques pioneered in machine translation (MT). Automatic evaluation metrics such as BLEU, METEOR, and ROUGE, are fundamental to the evaluation of MT systems and have been adopted as proxies of human evaluation in the code summarization domain. However, the extent to which automatic metrics agree with the gold standard of human evaluation has not been evaluated on code summarization tasks. Despite this, marginal improvements in metric scores are often used to discriminate between the performance of competing summarization models.

In this paper, we present a critical exploration of the applicability and interpretation of automatic metrics as evaluation techniques for code summarization tasks. We conduct an empirical study with 226 human annotators to assess the degree to which automatic metrics reflect human evaluation. Results indicate that metric improvements of less than 2 points do not guarantee systematic improvements in summarization quality, and are unreliable as proxies of human evaluation. When the difference between metric scores for two summarization approaches increases but remains within 5 points, some metrics such as METEOR and chrF become highly reliable proxies, whereas others, such as corpus BLEU, remain unreliable. Based on these findings, we make several recommendations for the use of automatic metrics to discriminate model performance in code summarization.

CCS CONCEPTS

- Software and its engineering → Software maintenance tools;
- General and reference → Metrics; Evaluation.

KEYWORDS

automatic evaluation metrics, code summarization, machine translation

ACM Reference Format:

Devjeet Roy, Sarah Fakhoury, and Venera Arnaudova. 2021. Reassessing Automatic Evaluation Metrics for Code Summarization Tasks. In *Proceedings of the 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '21)*, August 23–28, 2021, Athens, Greece. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3468264.3468588>



This work is licensed under a Creative Commons Attribution International 4.0 License. *ESEC/FSE '21, August 23–28, 2021, Athens, Greece*
© 2021 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-8562-6/21/08.
<https://doi.org/10.1145/3468264.3468588>

1 INTRODUCTION

Useful source code comments play a vital role in program comprehension and other software maintenance activities [45, 58]. However, proper documentation comes at a cost and producing well-written comments requires a substantial amount of effort on the part of software developers. This is one of the motivating factors behind why source code summarization is a rapidly growing research area—at least 18 papers proposing or evaluating automated summarization approaches are published in 2020 [1, 2, 10, 16, 18, 20, 22, 24, 27, 32, 47, 52, 55, 56, 60, 61, 63, 66].

The earliest code summarization approaches are based on strong syntactic theories of comment structure, information retrieval, and textual templates. These approaches typically evaluate the quality of a generated summary using human annotators who rate summaries on metrics such as: content adequacy [36], conciseness [23, 37], and fluency [37, 38, 50].

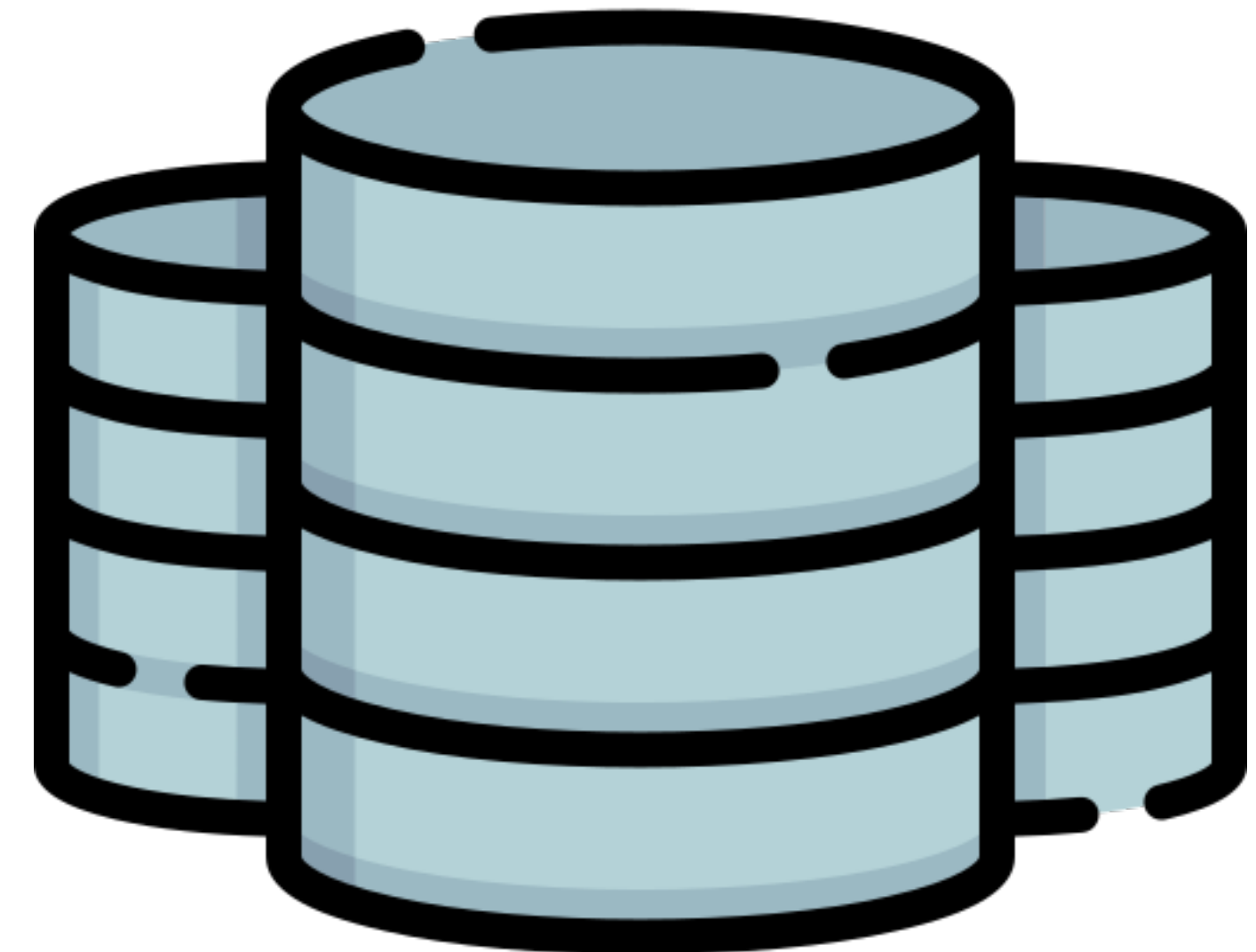
In the last 5 years, the solution space for code summarization has significantly shifted towards the widespread adoption of techniques and evaluation metrics from the Machine Translation (MT) domain. Prior to 2015, virtually all summarization approaches involved template or information retrieval based techniques. In 2015, the first paper using an MT approach was published at an SE conference [39], and, to the best of our knowledge, there are 35 papers thus far that propose an approach, an evaluation, or a critique of MT summarization approaches [2–5, 8–10, 12, 16–21, 23–25, 27, 28, 31, 35, 39, 47, 48, 53–57, 59–62, 65, 66]. Note that 7 of those are published in 2019 and 13 in 2020 alone.

One of the primary reasons for this shift is the idea that translating source code to its natural language equivalent holds many parallels with the concept of translating one natural language to another [23]. A side effect of this shift is a substantial change in the amount of data needed to evaluate code summarization approaches; training generative models requires a large amount of data. Naturally, the evaluation techniques used must also scale, and automated metrics provide an efficient way to evaluate the quality of generated summaries en masse.

Automated metrics designed to evaluate natural language translation approaches, such as BLEU [40], METEOR [6], and ROUGE [30], have been adopted by code summarization researchers where a generated summary is compared to a 'gold standard' or 'reference' summary. In the context of leading comment summaries, the reference summary is often the original accompanying comment written by a developer.

In the MT domain, BLEU has long been accepted as a de-facto best-practice metric. Recently, however, prominent machine translation researchers have raised concern over the use of BLEU [34, 42, 43], warning the MT community that the way BLEU is used

~5K



Evaluation of SIDE

Conciseness: Assesses the degree to which the summary contains unnecessary information

Reassessing Automatic Evaluation Metrics for Code Summarization Tasks

Devjeet Roy
devjeet.roy@wsu.edu
Washington State University
Pullman, WA, USA

Sarah Fakhoury
sarah.fakhoury@wsu.edu
Washington State University
Pullman, WA, USA

Venera Arnaoudova
venera.arnaoudova@wsu.edu
Washington State University
Pullman, WA, USA

ABSTRACT

In recent years, research in the domain of source code summarization has adopted data-driven techniques pioneered in machine translation (MT). Automatic evaluation metrics such as BLEU, METEOR, and ROUGE, are fundamental to the evaluation of MT systems and have been adopted as proxies of human evaluation in the code summarization domain. However, the extent to which automatic metrics agree with the gold standard of human evaluation has not been evaluated on code summarization tasks. Despite this, marginal improvements in metric scores are often used to discriminate between the performance of competing summarization models.

In this paper, we present a critical exploration of the applicability and interpretation of automatic metrics as evaluation techniques for code summarization tasks. We conduct an empirical study with 226 human annotators to assess the degree to which automatic metrics reflect human evaluation. Results indicate that metric improvements of less than 2 points do not guarantee systematic improvements in summarization quality, and are unreliable as proxies of human evaluation. When the difference between metric scores for two summarization approaches increases but remains within 5 points, some metrics such as METEOR and chrF become highly reliable proxies, whereas others, such as corpus BLEU, remain unreliable. Based on these findings, we make several recommendations for the use of automatic metrics to discriminate model performance in code summarization.

CCS CONCEPTS

- Software and its engineering → Software maintenance tools;
- General and reference → Metrics; Evaluation.

KEYWORDS

automatic evaluation metrics, code summarization, machine translation

ACM Reference Format:

Devjeet Roy, Sarah Fakhoury, and Venera Arnaoudova. 2021. Reassessing Automatic Evaluation Metrics for Code Summarization Tasks. In *Proceedings of the 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '21)*, August 23–28, 2021, Athens, Greece. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3468264.3468588>



This work is licensed under a Creative Commons Attribution International 4.0 License. *ESEC/FSE '21, August 23–28, 2021, Athens, Greece*
© 2021 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-8562-6/21/08.
<https://doi.org/10.1145/3468264.3468588>

1 INTRODUCTION

Useful source code comments play a vital role in program comprehension and other software maintenance activities [45, 58]. However, proper documentation comes at a cost and producing well-written comments requires a substantial amount of effort on the part of software developers. This is one of the motivating factors behind why source code summarization is a rapidly growing research area—at least 18 papers proposing or evaluating automated summarization approaches are published in 2020 [1, 2, 10, 16, 18, 20, 22, 24, 27, 32, 47, 52, 55, 56, 60, 61, 63, 66].

The earliest code summarization approaches are based on strong syntactic theories of comment structure, information retrieval, and textual templates. These approaches typically evaluate the quality of a generated summary using human annotators who rate summaries on metrics such as: content adequacy [36], conciseness [23, 37], and fluency [37, 38, 50].

In the last 5 years, the solution space for code summarization has significantly shifted towards the widespread adoption of techniques and evaluation metrics from the Machine Translation (MT) domain. Prior to 2015, virtually all summarization approaches involved template or information retrieval based techniques. In 2015, the first paper using an MT approach was published at an SE conference [39], and, to the best of our knowledge, there are 35 papers thus far that propose an approach, an evaluation, or a critique of MT summarization approaches [2–5, 8–10, 12, 16–21, 23–25, 27, 28, 31, 35, 39, 47, 48, 53–57, 59–62, 65, 66]. Note that 7 of those are published in 2019 and 13 in 2020 alone.

One of the primary reasons for this shift is the idea that translating source code to its natural language equivalent holds many parallels with the concept of translating one natural language to another [23]. A side effect of this shift is a substantial change in the amount of data needed to evaluate code summarization approaches; training generative models requires a large amount of data. Naturally, the evaluation techniques used must also scale, and automated metrics provide an efficient way to evaluate the quality of generated summaries en masse.

Automated metrics designed to evaluate natural language translation approaches, such as BLEU [40], METEOR [6], and ROUGE [30], have been adopted by code summarization researchers where a generated summary is compared to a 'gold standard' or 'reference' summary. In the context of leading comment summaries, the reference summary is often the original accompanying comment written by a developer.

In the MT domain, BLEU has long been accepted as a de-facto best-practice metric. Recently, however, prominent machine translation researchers have raised concern over the use of BLEU [34, 42, 43], warning the MT community that the way BLEU is used

Evaluation of SIDE

Conciseness: Assesses the degree to which the summary contains unnecessary information

Fluency: Evaluates the “smoothness” rate in the generated summary

Reassessing Automatic Evaluation Metrics for Code Summarization Tasks

Devjeet Roy
devjeet.roy@wsu.edu
Washington State University
Pullman, WA, USA

Sarah Fakhoury
sarah.fakhoury@wsu.edu
Washington State University
Pullman, WA, USA

Venera Arnaoudova
venera.arnaoudova@wsu.edu
Washington State University
Pullman, WA, USA

ABSTRACT

In recent years, research in the domain of source code summarization has adopted data-driven techniques pioneered in machine translation (MT). Automatic evaluation metrics such as BLEU, METEOR, and ROUGE, are fundamental to the evaluation of MT systems and have been adopted as proxies of human evaluation in the code summarization domain. However, the extent to which automatic metrics agree with the gold standard of human evaluation has not been evaluated on code summarization tasks. Despite this, marginal improvements in metric scores are often used to discriminate between the performance of competing summarization models.

In this paper, we present a critical exploration of the applicability and interpretation of automatic metrics as evaluation techniques for code summarization tasks. We conduct an empirical study with 226 human annotators to assess the degree to which automatic metrics reflect human evaluation. Results indicate that metric improvements of less than 2 points do not guarantee systematic improvements in summarization quality, and are unreliable as proxies of human evaluation. When the difference between metric scores for two summarization approaches increases but remains within 5 points, some metrics such as METEOR and chrF become highly reliable proxies, whereas others, such as corpus BLEU, remain unreliable. Based on these findings, we make several recommendations for the use of automatic metrics to discriminate model performance in code summarization.

CCS CONCEPTS

- Software and its engineering → Software maintenance tools;
- General and reference → Metrics; Evaluation.

KEYWORDS

automatic evaluation metrics, code summarization, machine translation

ACM Reference Format:

Devjeet Roy, Sarah Fakhoury, and Venera Arnaoudova. 2021. Reassessing Automatic Evaluation Metrics for Code Summarization Tasks. In *Proceedings of the 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '21)*, August 23–28, 2021, Athens, Greece. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3468264.3468588>



This work is licensed under a Creative Commons Attribution International 4.0 License. *ESEC/FSE '21, August 23–28, 2021, Athens, Greece*
© 2021 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-8562-6/21/08.
<https://doi.org/10.1145/3468264.3468588>

1 INTRODUCTION

Useful source code comments play a vital role in program comprehension and other software maintenance activities [45, 58]. However, proper documentation comes at a cost and producing well-written comments requires a substantial amount of effort on the part of software developers. This is one of the motivating factors behind why source code summarization is a rapidly growing research area—at least 18 papers proposing or evaluating automated summarization approaches are published in 2020 [1, 2, 10, 16, 18, 20, 22, 24, 27, 32, 47, 52, 55, 56, 60, 61, 63, 66].

The earliest code summarization approaches are based on strong syntactic theories of comment structure, information retrieval, and textual templates. These approaches typically evaluate the quality of a generated summary using human annotators who rate summaries on metrics such as: content adequacy [36], conciseness [23, 37], and fluency [37, 38, 50].

In the last 5 years, the solution space for code summarization has significantly shifted towards the widespread adoption of techniques and evaluation metrics from the Machine Translation (MT) domain. Prior to 2015, virtually all summarization approaches involved template or information retrieval based techniques. In 2015, the first paper using an MT approach was published at an SE conference [39], and, to the best of our knowledge, there are 35 papers thus far that propose an approach, an evaluation, or a critique of MT summarization approaches [2–5, 8–10, 12, 16–21, 23–25, 27, 28, 31, 35, 39, 47, 48, 53–57, 59–62, 65, 66]. Note that 7 of those are published in 2019 and 13 in 2020 alone.

One of the primary reasons for this shift is the idea that translating source code to its natural language equivalent holds many parallels with the concept of translating one natural language to another [23]. A side effect of this shift is a substantial change in the amount of data needed to evaluate code summarization approaches; training generative models requires a large amount of data. Naturally, the evaluation techniques used must also scale, and automated metrics provide an efficient way to evaluate the quality of generated summaries en masse.

Automated metrics designed to evaluate natural language translation approaches, such as BLEU [40], METEOR [6], and ROUGE [30], have been adopted by code summarization researchers where a generated summary is compared to a ‘gold standard’ or ‘reference’ summary. In the context of leading comment summaries, the reference summary is often the original accompanying comment written by a developer.

In the MT domain, BLEU has long been accepted as a de-facto best-practice metric. Recently, however, prominent machine translation researchers have raised concern over the use of BLEU [34, 42, 43], warning the MT community that the way BLEU is used

FSE'2021

Evaluation of SIDE

Reassessing Automatic Evaluation Metrics for Code Summarization Tasks

Devjeet Roy
devjeet.roy@wsu.edu
Washington State University
Pullman, WA, USA

Sarah Fakhoury
sarah.fakhoury@wsu.edu
Washington State University
Pullman, WA, USA

Venera Arnaoudova
venera.arnaoudova@wsu.edu
Washington State University
Pullman, WA, USA

ABSTRACT

In recent years, research in the domain of source code summarization has adopted data-driven techniques pioneered in machine translation (MT). Automatic evaluation metrics such as BLEU, METEOR, and ROUGE, are fundamental to the evaluation of MT systems and have been adopted as proxies of human evaluation in the code summarization domain. However, the extent to which automatic metrics agree with the gold standard of human evaluation has not been evaluated on code summarization tasks. Despite this, marginal improvements in metric scores are often used to discriminate between the performance of competing summarization models.

In this paper, we present a critical exploration of the applicability and interpretation of automatic metrics as evaluation techniques for code summarization tasks. We conduct an empirical study with 226 human annotators to assess the degree to which automatic metrics reflect human evaluation. Results indicate that metric improvements of less than 2 points do not guarantee systematic improvements in summarization quality, and are unreliable as proxies of human evaluation. When the difference between metric scores for two summarization approaches increases but remains within 5 points, some metrics such as METEOR and chrF become highly reliable proxies, whereas others, such as corpus BLEU, remain unreliable. Based on these findings, we make several recommendations for the use of automatic metrics to discriminate model performance in code summarization.

CCS CONCEPTS

- Software and its engineering → Software maintenance tools;
- General and reference → Metrics; Evaluation.

KEYWORDS

automatic evaluation metrics, code summarization, machine translation

ACM Reference Format:

Devjeet Roy, Sarah Fakhoury, and Venera Arnaoudova. 2021. Reassessing Automatic Evaluation Metrics for Code Summarization Tasks. In *Proceedings of the 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '21)*, August 23–28, 2021, Athens, Greece. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3468264.3468588>



This work is licensed under a Creative Commons Attribution International 4.0 License. *ESEC/FSE '21, August 23–28, 2021, Athens, Greece*
© 2021 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-8562-6/21/08.
<https://doi.org/10.1145/3468264.3468588>

1 INTRODUCTION

Useful source code comments play a vital role in program comprehension and other software maintenance activities [45, 58]. However, proper documentation comes at a cost and producing well-written comments requires a substantial amount of effort on the part of software developers. This is one of the motivating factors behind why source code summarization is a rapidly growing research area—at least 18 papers proposing or evaluating automated summarization approaches are published in 2020 [1, 2, 10, 16, 18, 20, 22, 24, 27, 32, 47, 52, 55, 56, 60, 61, 63, 66].

The earliest code summarization approaches are based on strong syntactic theories of comment structure, information retrieval, and textual templates. These approaches typically evaluate the quality of a generated summary using human annotators who rate summaries on metrics such as: content adequacy [36], conciseness [23, 37], and fluency [37, 38, 50].

In the last 5 years, the solution space for code summarization has significantly shifted towards the widespread adoption of techniques and evaluation metrics from the Machine Translation (MT) domain. Prior to 2015, virtually all summarization approaches involved template or information retrieval based techniques. In 2015, the first paper using an MT approach was published at an SE conference [39], and, to the best of our knowledge, there are 35 papers thus far that propose an approach, an evaluation, or a critique of MT summarization approaches [2–5, 8–10, 12, 16–21, 23–25, 27, 28, 31, 35, 39, 47, 48, 53–57, 59–62, 65, 66]. Note that 7 of those are published in 2019 and 13 in 2020 alone.

One of the primary reasons for this shift is the idea that translating source code to its natural language equivalent holds many parallels with the concept of translating one natural language to another [23]. A side effect of this shift is a substantial change in the amount of data needed to evaluate code summarization approaches; training generative models requires a large amount of data. Naturally, the evaluation techniques used must also scale, and automated metrics provide an efficient way to evaluate the quality of generated summaries en masse.

Automated metrics designed to evaluate natural language translation approaches, such as BLEU [40], METEOR [6], and ROUGE [30], have been adopted by code summarization researchers where a generated summary is compared to a ‘gold standard’ or ‘reference’ summary. In the context of leading comment summaries, the reference summary is often the original accompanying comment written by a developer.

In the MT domain, BLEU has long been accepted as a de-facto best-practice metric. Recently, however, prominent machine translation researchers have raised concern over the use of BLEU [34, 42, 43], warning the MT community that the way BLEU is used

Conciseness: Assesses the degree to which the summary contains unnecessary information

Fluency: Evaluates the “smoothness” rate in the generated summary

Content Adequacy: Assess the extent to which the summary lacks information needed to understand the code

FSE'2021

Evaluation of SIDE

Reassessing Automatic Evaluation Metrics for Code Summarization Tasks

Devjeet Roy
devjeet.roy@wsu.edu
Washington State University
Pullman, WA, USA

Sarah Fakhoury
sarah.fakhoury@wsu.edu
Washington State University
Pullman, WA, USA

Venera Arnaudova
venera.arnaudova@wsu.edu
Washington State University
Pullman, WA, USA

ABSTRACT

In recent years, research in the domain of source code summarization has adopted data-driven techniques pioneered in machine translation (MT). Automatic evaluation metrics such as BLEU, METEOR, and ROUGE, are fundamental to the evaluation of MT systems and have been adopted as proxies of human evaluation in the code summarization domain. However, the extent to which automatic metrics agree with the gold standard of human evaluation has not been evaluated on code summarization tasks. Despite this, marginal improvements in metric scores are often used to discriminate between the performance of competing summarization models.

In this paper, we present a critical exploration of the applicability and interpretation of automatic metrics as evaluation techniques for code summarization tasks. We conduct an empirical study with 226 human annotators to assess the degree to which automatic metrics reflect human evaluation. Results indicate that metric improvements of less than 2 points do not guarantee systematic improvements in summarization quality, and are unreliable as proxies of human evaluation. When the difference between metric scores for two summarization approaches increases but remains within 5 points, some metrics such as METEOR and chrF become highly reliable proxies, whereas others, such as corpus BLEU, remain unreliable. Based on these findings, we make several recommendations for the use of automatic metrics to discriminate model performance in code summarization.

CCS CONCEPTS

- Software and its engineering → Software maintenance tools;
- General and reference → Metrics; Evaluation.

KEYWORDS

automatic evaluation metrics, code summarization, machine translation

ACM Reference Format:

Devjeet Roy, Sarah Fakhoury, and Venera Arnaudova. 2021. Reassessing Automatic Evaluation Metrics for Code Summarization Tasks. In *Proceedings of the 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '21)*, August 23–28, 2021, Athens, Greece. ACM, New York, NY, USA, 12 pages. <https://doi.org/10.1145/3468264.3468588>



This work is licensed under a Creative Commons Attribution International 4.0 License. *ESEC/FSE '21, August 23–28, 2021, Athens, Greece*
© 2021 Copyright held by the owner/author(s).
ACM ISBN 978-1-4503-8562-6/21/08.
<https://doi.org/10.1145/3468264.3468588>

1 INTRODUCTION

Useful source code comments play a vital role in program comprehension and other software maintenance activities [45, 58]. However, proper documentation comes at a cost and producing well-written comments requires a substantial amount of effort on the part of software developers. This is one of the motivating factors behind why source code summarization is a rapidly growing research area—at least 18 papers proposing or evaluating automated summarization approaches are published in 2020 [1, 2, 10, 16, 18, 20, 22, 24, 27, 32, 47, 52, 55, 56, 60, 61, 63, 66].

The earliest code summarization approaches are based on strong syntactic theories of comment structure, information retrieval, and textual templates. These approaches typically evaluate the quality of a generated summary using human annotators who rate summaries on metrics such as: content adequacy [36], conciseness [23, 37], and fluency [37, 38, 50].

In the last 5 years, the solution space for code summarization has significantly shifted towards the widespread adoption of techniques and evaluation metrics from the Machine Translation (MT) domain. Prior to 2015, virtually all summarization approaches involved template or information retrieval based techniques. In 2015, the first paper using an MT approach was published at an SE conference [39], and, to the best of our knowledge, there are 35 papers thus far that propose an approach, an evaluation, or a critique of MT summarization approaches [2–5, 8–10, 12, 16–21, 23–25, 27, 28, 31, 35, 39, 47, 48, 53–57, 59–62, 65, 66]. Note that 7 of those are published in 2019 and 13 in 2020 alone.

One of the primary reasons for this shift is the idea that translating source code to its natural language equivalent holds many parallels with the concept of translating one natural language to another [23]. A side effect of this shift is a substantial change in the amount of data needed to evaluate code summarization approaches; training generative models requires a large amount of data. Naturally, the evaluation techniques used must also scale, and automated metrics provide an efficient way to evaluate the quality of generated summaries en masse.

Automated metrics designed to evaluate natural language translation approaches, such as BLEU [40], METEOR [6], and ROUGE [30], have been adopted by code summarization researchers where a generated summary is compared to a ‘gold standard’ or ‘reference’ summary. In the context of leading comment summaries, the reference summary is often the original accompanying comment written by a developer.

In the MT domain, BLEU has long been accepted as a de-facto best-practice metric. Recently, however, prominent machine translation researchers have raised concern over the use of BLEU [34, 42, 43], warning the MT community that the way BLEU is used

Conciseness: Assesses the degree to which the summary contains unnecessary information

Fluency: Evaluates the “smoothness” rate in the generated summary

Content Adequacy: Assess the extent to which the summary lacks information needed to understand the code

{0...5}

FSE'2021



Dependent Variables

Conciseness

Fluency

Content



Independent Variables

SIDE

Word/Overlap-based

Embedding-based



Independent Variables

SIDE

Word/Overlap-based

Embedding-based

13

1 insert an arbitrary number of days
add n days to this date

P
G

BLEU SCORE

ROUGE SCORE

BERT SCORE

.....

2 add a new tab with a canonical ...
add a new tab with a default ...

P
G

BLEU SCORE

ROUGE SCORE

BERT SCORE

.....

3 create a new result set supporting ...
implement a new result set ...

P
G

BLEU SCORE

ROUGE SCORE

BERT SCORE

.....

4 constructs an instance from...
constructs an instance from...

P
G

BLEU SCORE

ROUGE SCORE

BERT SCORE

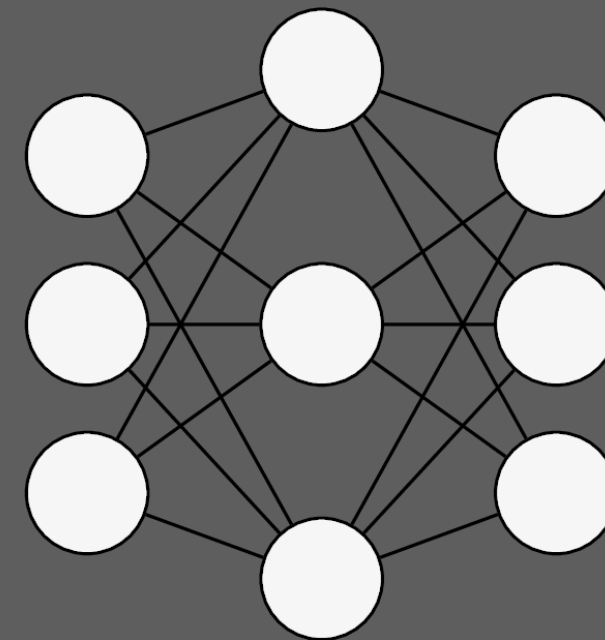
.....

SIDE

INPUT

```
Create a connection to the consumer.  
public ConnectionConsumer  
  createConnectionConsumer  
  ...{  
  if(LOGGER.isTraceEnabled())  
  {  
    ActiveMQRALogger.LOGGER.  
      trace("Create  
            connectionConsumer");  
  }  
  else {  
    ...  
  }  
}
```

MPNet



PREDICTION

0.81

MPNet: Masked and permuted pre-training for language understanding

~5K

1

Whether **SIDE** can be effectively used to capture new “**dimensions**” where state-of-the-art metrics fails

1

PRINCIPAL COMPONENT ANALYSIS (PCA)

Whether **SIDE** can be effectively used to capture new “**dimensions**” where state-of-the-art metrics fails

1

PRINCIPAL COMPONENT ANALYSIS (PCA)

PC1

55%

BERTScore
RougeScore

Whether **SIDE** can be effectively used to capture new “**dimensions**” where state-of-the-art metrics fails

1

PRINCIPAL COMPONENT ANALYSIS (PCA)

PC1

55%

BERTScore
RougeScore

PC2

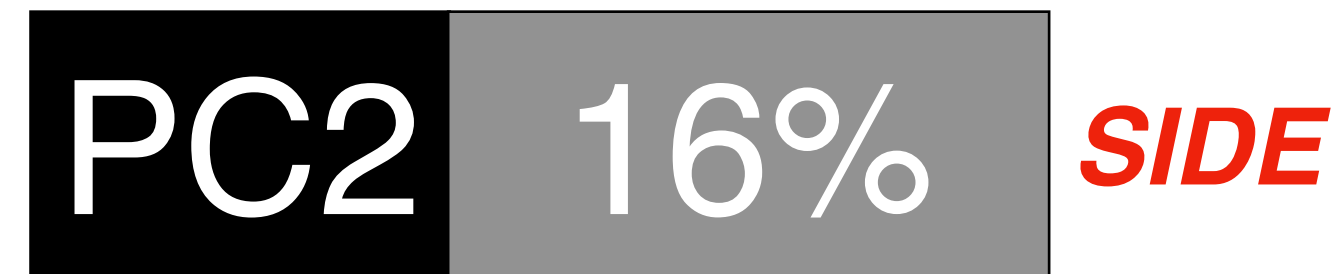
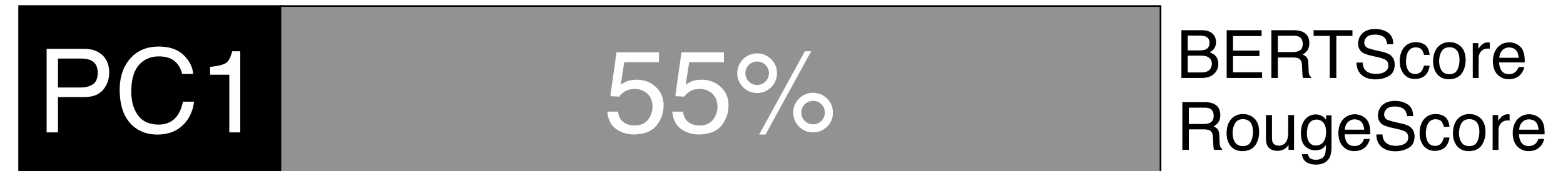
16%

SIDE

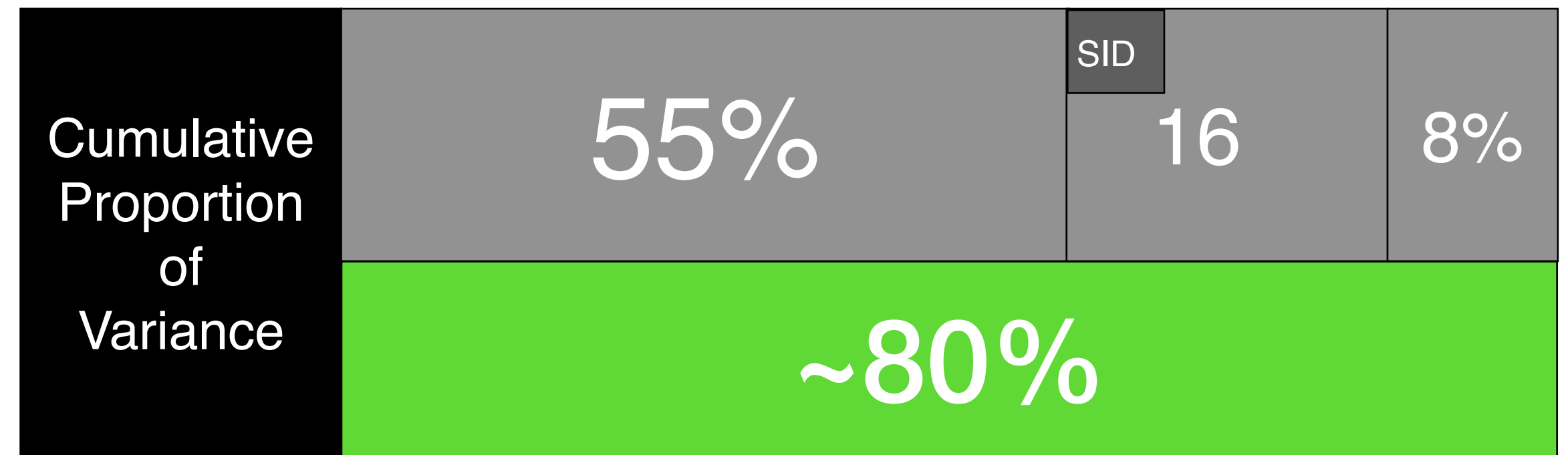
Whether **SIDE** can be effectively used to capture new “**dimensions**” where state-of-the-art metrics fails

1

PRINCIPAL COMPONENT ANALYSIS (PCA)



Whether **SIDE** can be effectively used to capture new “**dimensions**” where state-of-the-art metrics fails



2

ORDERED LOGISTIC REGRESSION MOD

Whether **SIDE** allows for a more effective representations of developers' evaluations regarding the **quality** of automatically generated summaries

2

ORDERED LOGISTIC REGRESSION MOD

Conciseness

38%

SIDE

24%

C_COEFF

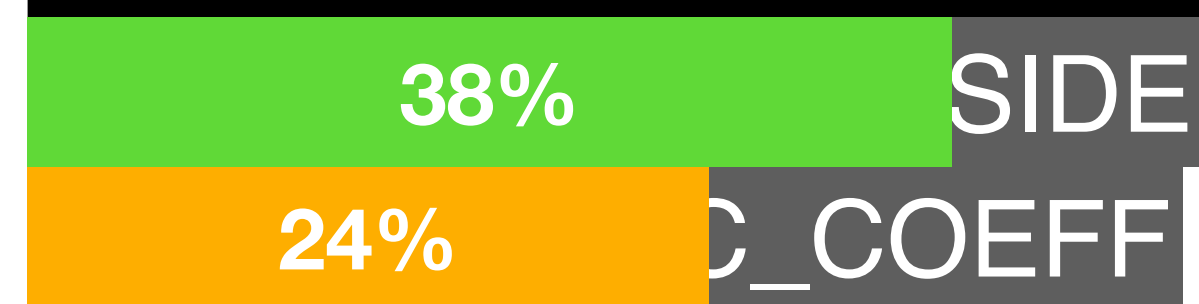
Whether **SIDE** allows for a more effective representations of developers' evaluations regarding the **quality** of automatically generated summaries

2

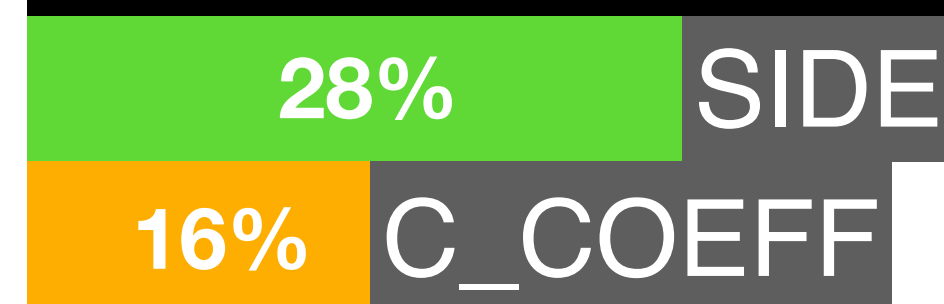
ORDERED LOGISTIC REGRESSION MODEL

Whether **SIDE** allows for a more effective representations of developers' evaluations regarding the **quality** of automatically generated summaries

Conciseness



Fluency

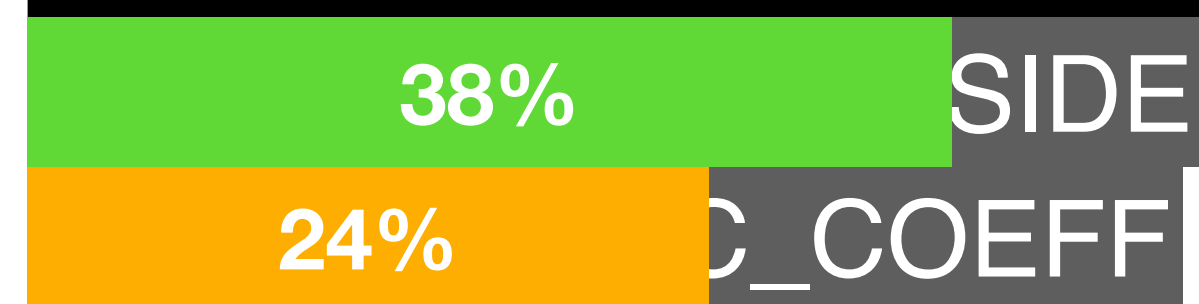


2

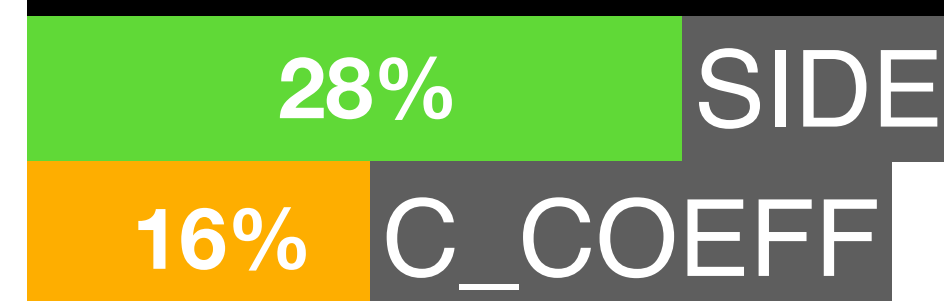
ORDERED LOGISTIC REGRESSION MODEL

Whether **SIDE** allows for a more effective representations of developers' evaluations regarding the **quality** of automatically generated summaries

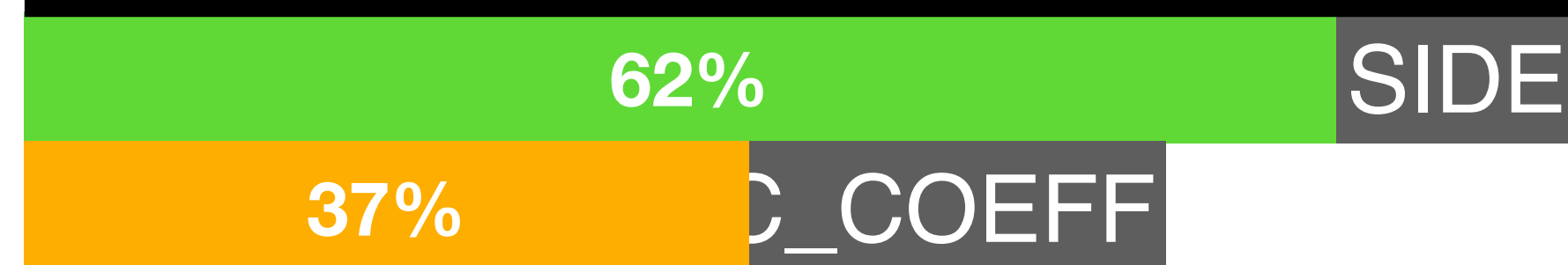
Conciseness



Fluency



Content Adequacy



METHOD

```
public Element asElement() {  
    return this.component.createCopy();  
}
```

REFERENCE SUMMARY

not yet documented

GENERATED SUMMARY

returns a copy of this component as an element

SCORES

BLEU-1	BertScore-R	SentenceBERT_CS	InferNet_CS	ROUGE-1-P	Rouge-4-R	Rouge-W-R	SIDE	DA	Adequacy	Concise	Fluency
0.34	0.00	0.08	0.43	0.00	0.00	0.00	0.91	88	4	4	5

Note: Red arrows in the original image point to InferNet_CS (0.43) labeled 'SOTA', SIDE (0.91) labeled 'OUR', and Concise (4) labeled 'HUMAN'.



State-of-the-art metrics are not **sufficient** to describe **humans' assessment** of automatically generated code summaries

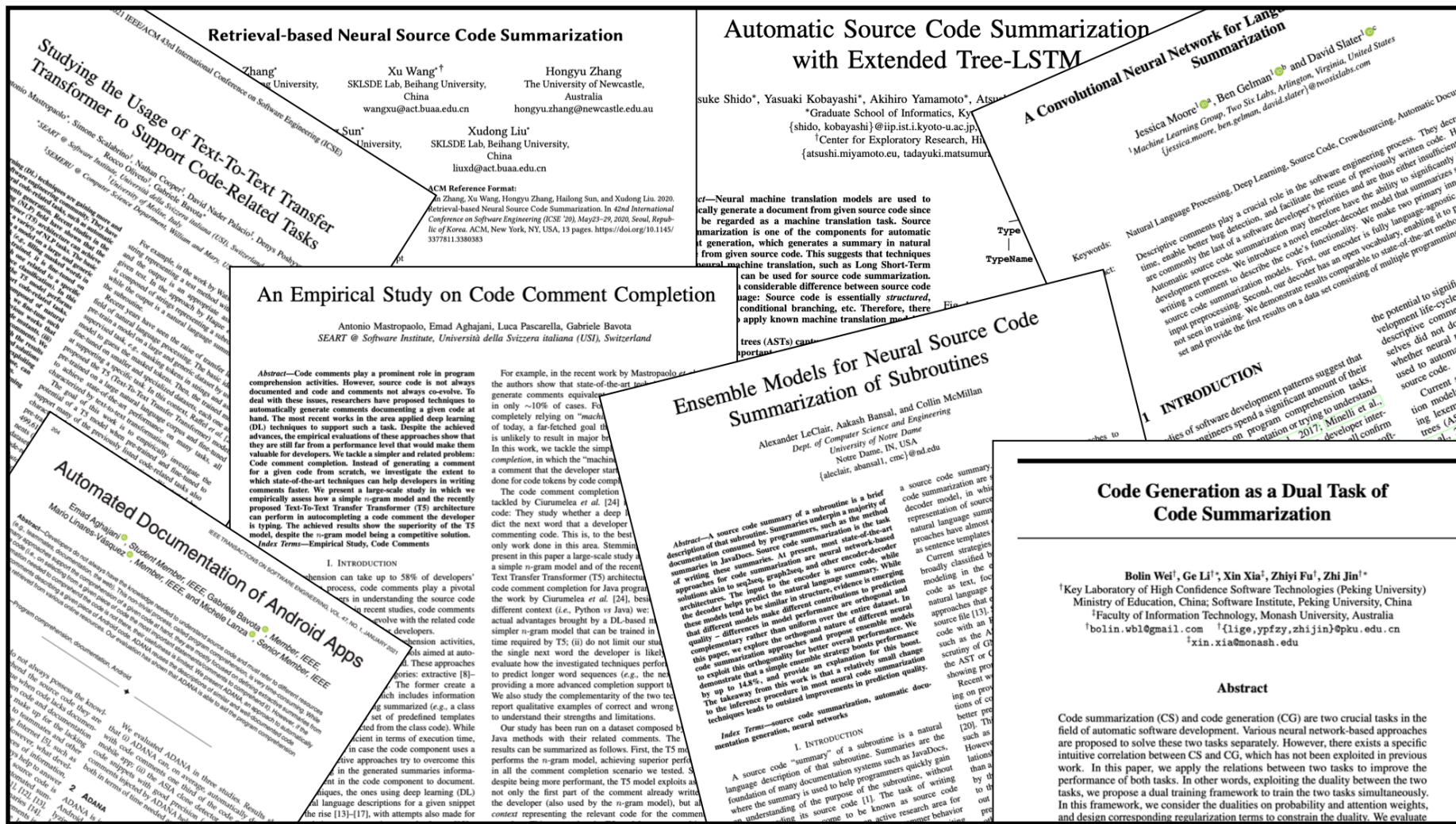


State-of-the-art metrics are not **sufficient** to describe **humans' assessment** of automatically generated code summaries



A more comprehensive evaluation **framework** which takes into account more traditional metrics and ***SIDE-like*** metrics are needed to capture more **dimensions**





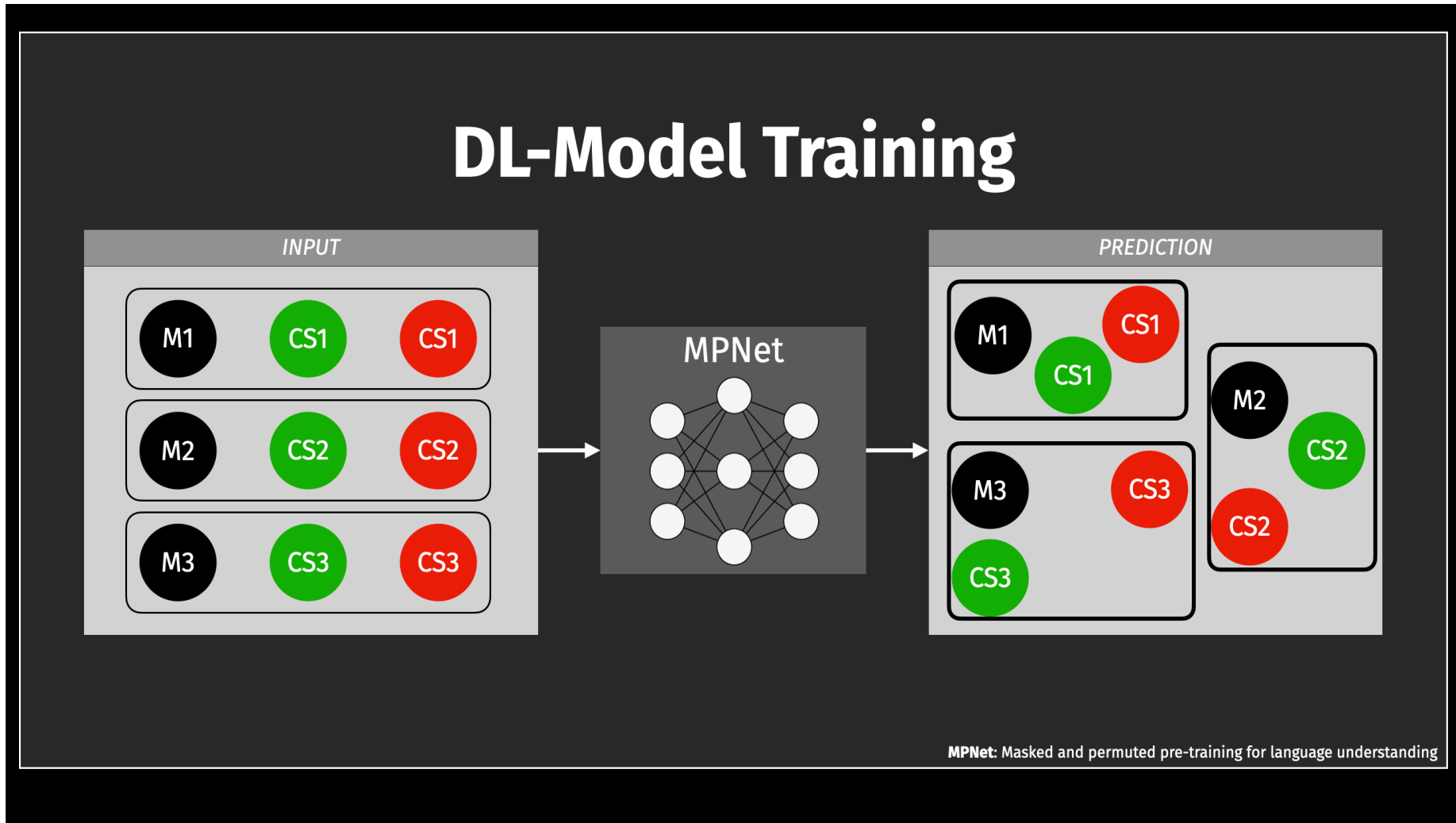
Reads the contents of this source as a string. PR

Get the textual information from this source and represent it as a string. GT

```
public String read() throws IOException {
    Closer closer = Closer.create();
    try {
        Reader reader = closer.register(openStream());
        return CharStreams.toString(reader);
    } catch (Throwable e) {
        throw closer.rethrow(e);
    } finally { closer.close(); }
}
```

SEMANTICALLY EQUIVALENT CODE SUMMARIES

BLEU SCORE: 0.21



2

Whether **SIDE** allows for a more effective representations of developers' evaluations regarding the **quality** of automatically generated summaries

ORDERED LOGISTIC REGRESSION MODEL

Conciseness	
38%	SIDE
24%	C_COEFF

Fluency	
28%	SIDE
16%	C_COEFF

Content Adequacy	
62%	SIDE
37%	C_COEFF

<https://github.com/side-metric/summarization>



Software Institute
Seminars @USI
2023-2024



Stay in touch!

✉ antonio.mastropaolo@usi.ch

🌐 antoniomastropaolo.com

🐦 AntonioMastro2