

Personalized Code Readability Assessment: Are We There Yet?

Antonio Vitale, Emanuela Guglielmi, Rocco Oliveto, Simone Scalabrino



ICPC  **2025**

**The 33rd IEEE/ACM International
Conference on Program Comprehension**

Ottawa, ON, Canada – April 27 – 28, 2025

CODE READABILITY

```
for (int i=1;i<=100;i++) {  
    int a=((528>>i%15-1)&1)*4;  
    int b=(-2128340926>>(i%15)*2)&3)*4;  
    System.out.println("FizzBuzz".substring(a,b)+(a==b?i:""));  
}
```

```
for (int i=1;i<=100;i++) {  
    String fizzBuzz = "";  
    if (i % 3 == 0)  
        fizzBuzz += "Fizz";  
  
    if (i % 5 == 0)  
        fizzBuzz += "Buzz";  
  
    if (fizzBuzz.isEmpty())  
        fizzBuzz += i;  
  
    System.out.println(fizzBuzz);  
}
```


CODE READABILITY

THE SPECIAL ISSUE ON THE BEST 2008 BEST PAPERS

Learning a Metric for Code Readability

Raymond PL. Buse, Wesley R. Werner

(Invited Paper)

Abstract: In this paper, we explore the concept of code readability and investigate its relation to software quality. With data collected from 1023 source code reviews, we derive associations between a sample set of code style features and human ratings of readability. Using these features, we constructed an automated readability measure and show that it can be 80% effective, and better than a human on average, at predicting readability judgments. Furthermore, we show that this metric correlates strongly with three measures of software quality: code changes, automated defect reports, and developer messages. We measure these correlations on over 2.8 million lines of code, as well as qualitatively, over many reviews of actual projects. Finally, we discuss the implications of this study on programming language design and engineering practice. For example, our results suggest that correctness, or its correlates, are more important than we have been able to find in past judgments of readability.

Index Terms: software readability, program understanding, machine learning, software maintenance, code metrics, Feature22

1 INTRODUCTION

When a new task is to be undertaken, the readability of a program is defined in its maintainability, and is thus a key factor in overall software quality. Typically, maintainability will consume over 70% of the total lifecycle cost of a software product [4]. Aggravated claims that source code readability and documentation maintainability are both critical to the maintainability of a project [1]. Other researchers have noted that the use of reading code is the most time-consuming component of all maintenance activities [8], [32], [34]. Readability is so significant, in fact, that Elfrink and Manthey, after recognizing that many commercial programs were much more difficult to read than necessary, proposed adding a development phase in which the program is made more readable [8]. Knight and Meyer suggested that one phase of source code inspection should be a check of the source code for readability [23] to ensure maintainability, portability, and reusability of the code. Harrod proposed adding a dedicated readability and documentation group in the development team, observing that "without established and consistent guidelines for readability, individual reviewers may not be able to help much" [19].

His hypothesis that programmers have some intuitive notion of this concept, and that program features such as indentation (e.g., as in Python [9]), choice of identifier names [23], and comments [32] are important parts. Dijkstra, for example, claimed that the readability of a program depends largely on the quality of the comments in its sequencing control; he is conjectured that "correctness" encompasses program understanding and, therefore, readability [10].

It is important to note that readability is not the same

A Simpler Model of Software Readability

Daryl Poznett
University of California, Davis
Daryl Poznett
dpoznett@ucdavis.edu

Abram Hindle
University of California, Davis
Darin C. Casin
ahindle@ucdavis.edu

Prem Dewarshi
University of California, Davis
Daryl C. Casin
dewarshi@ucdavis.edu

ABSTRACT

Software readability is a property that influences how fast a group of people can read and understand. Since readability can affect maintainability, quality, etc., programmers are very concerned about the readability of code. If software readability checks could be built, they could be integrated into development toolchains, thus dramatically assisting developers about the readability level of the code. Unfortunately, readability is a subjective code property, and not amenable to direct automated measurement.

In a recently published study, Buse et al. asked 100 participants to rate code snippets for readability. Varying aspects of code readability scores of each snippet; they then built a fairly complex predictive model for these mean scores using a large, diverse set of directly measurable source code properties. We build on this work we present a simple, intuitive theory of readability, based on the use of code snippets, and show how this theory leads to a much simpler, yet also intuitively significant, model of the more readability scores than the previous 28 features. Our model uses well-known code metrics and features that are more directly measurable, and thus more likely to be used by programmers for the need to estimate source code readability. In particular, we present a theory of code readability that is more directly measurable, and thus more likely to be used by programmers for the need to estimate source code readability.

There is much previous work about readability [2], [5], [18], [21], [25]. One major line of studies of readability is the difficulty of experimentally inferring what is necessary for a subjective perception. However, of subjective perception, the difficulty of inferring what is necessary for a subjective perception is a critical to software readability. In particular, we present a theory of code readability that is more directly measurable, and thus more likely to be used by programmers for the need to estimate source code readability. In particular, we present a theory of code readability that is more directly measurable, and thus more likely to be used by programmers for the need to estimate source code readability.

Categories and Subject Descriptors

D.1.2 Software Systems Applications; D.1.2 Software Systems Applications; D.1.2 Software Systems Applications; D.1.2 Software Systems Applications

General Terms

Human Factors, Theory, Measurement

Keywords

Readability, Harrod, Entry, Revision

Permission to make digital or hard copies of part of this work for personal or classroom use is granted by the publisher for copies that are not made for profit or for commercial purposes and that are not made for profit or for commercial purposes. This permission is granted without fee or charge, except for the copying fee. This permission is granted without fee or charge, except for the copying fee. This permission is granted without fee or charge, except for the copying fee.

A General Software Readability Model

Jonathan Dunn
Department of Computer Science
University of Virginia
Charlottesville, Virginia
jcd2@virginia.edu

Abstract: Code readability is a property that influences how fast a group of people can read and understand. Since readability can affect maintainability, quality, etc., programmers are very concerned about the readability of code. If software readability checks could be built, they could be integrated into development toolchains, thus dramatically assisting developers about the readability level of the code. Unfortunately, readability is a subjective code property, and not amenable to direct automated measurement.

1. INTRODUCTION

Modern software developers spend more time maintaining and evolving existing software than writing new code [13], [23]. Software readability, a fundamental notion related to the comprehension of text, is critical to software readability: reading code is a necessary first step toward maintaining it.

Much research, both recent and longstanding, has argued that readability plays a large role in software maintenance. A well-known example is Knuth, who viewed readability as essential to the success of Linear Programming [4]. He argued that a program should be viewed as a single programming language, and that the programmer should be able to read it as if it were a single programming language. He argued that a program should be viewed as a single programming language, and that the programmer should be able to read it as if it were a single programming language.

Knuth and Myers argued that the effectiveness of syntax highlighting suggests that visual cues promote, particularly maintainability and readability and should thus be a first-class phase of software development [13]. Based et al. showed that inspections paid by making techniques are better at reducing phase of software development [13]. Based et al. showed that inspections paid by making techniques are better at reducing phase of software development [13]. Based et al. showed that inspections paid by making techniques are better at reducing phase of software development [13].

Based et al. showed that inspections paid by making techniques are better at reducing phase of software development [13]. Based et al. showed that inspections paid by making techniques are better at reducing phase of software development [13]. Based et al. showed that inspections paid by making techniques are better at reducing phase of software development [13]. Based et al. showed that inspections paid by making techniques are better at reducing phase of software development [13].

Improving Code Readability Models with Textual Features

Simone Scabbini*, Mario Linarez-Vázquez†, Denis Poshynyak† and Rocco Oliviero†
*University of Molise, Pesche (IS), Italy
†The College of William and Mary, Williamsburg, Virginia, USA

Abstract: Code readability is one of the most frequent activities in software maintenance. Before implementing changes, it is necessary to fully understand source code often written by other developers. Thus, readability is a critical factor in software maintenance. In this paper, we propose a new model of code readability that incorporates textual features. We show that this model is more effective than previous models in predicting the quality of source code reviews.

Under the hypothesis that source code readability should be captured using both syntactic and textual code features, in this paper we propose a new model of code readability that incorporates textual features. We show that this model is more effective than previous models in predicting the quality of source code reviews. Under the hypothesis that source code readability should be captured using both syntactic and textual code features, in this paper we propose a new model of code readability that incorporates textual features. We show that this model is more effective than previous models in predicting the quality of source code reviews.

Under the hypothesis that source code readability should be captured using both syntactic and textual code features, in this paper we propose a new model of code readability that incorporates textual features. We show that this model is more effective than previous models in predicting the quality of source code reviews. Under the hypothesis that source code readability should be captured using both syntactic and textual code features, in this paper we propose a new model of code readability that incorporates textual features. We show that this model is more effective than previous models in predicting the quality of source code reviews.

Under the hypothesis that source code readability should be captured using both syntactic and textual code features, in this paper we propose a new model of code readability that incorporates textual features. We show that this model is more effective than previous models in predicting the quality of source code reviews. Under the hypothesis that source code readability should be captured using both syntactic and textual code features, in this paper we propose a new model of code readability that incorporates textual features. We show that this model is more effective than previous models in predicting the quality of source code reviews.

Under the hypothesis that source code readability should be captured using both syntactic and textual code features, in this paper we propose a new model of code readability that incorporates textual features. We show that this model is more effective than previous models in predicting the quality of source code reviews. Under the hypothesis that source code readability should be captured using both syntactic and textual code features, in this paper we propose a new model of code readability that incorporates textual features. We show that this model is more effective than previous models in predicting the quality of source code reviews.

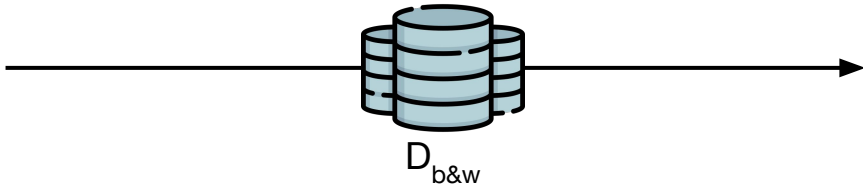
Under the hypothesis that source code readability should be captured using both syntactic and textual code features, in this paper we propose a new model of code readability that incorporates textual features. We show that this model is more effective than previous models in predicting the quality of source code reviews. Under the hypothesis that source code readability should be captured using both syntactic and textual code features, in this paper we propose a new model of code readability that incorporates textual features. We show that this model is more effective than previous models in predicting the quality of source code reviews.

2009

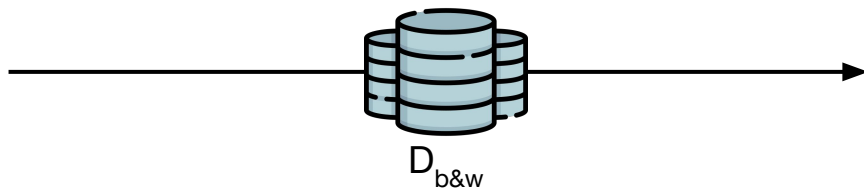
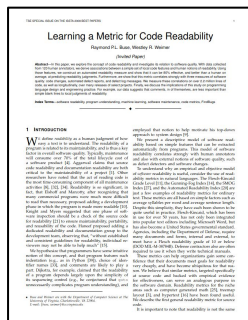
2011

2012

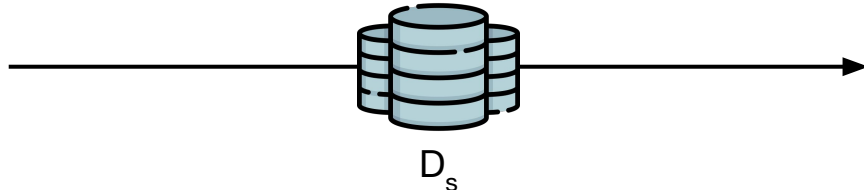
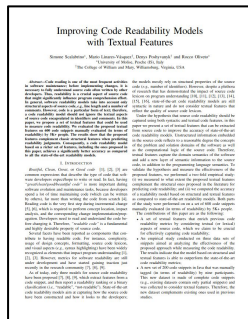
2016



CODE READABILITY



Snippets: 100
Developers: 120
Granularity: Snippet
Assessments: Complete



Snippets: 200
Developers: 9
Granularity: Method
Assessments: Complete

CODE READABILITY

Developers Assessments

Snippets

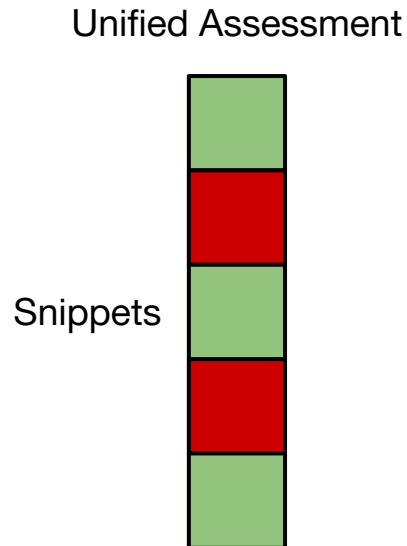
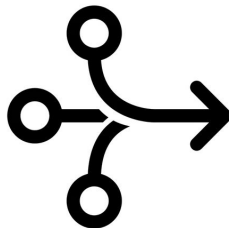
3	4	2	4	5
1	3	2	1	1
4	1	5	4	5
5	1	1	1	3
4	3	2	5	5

CODE READABILITY

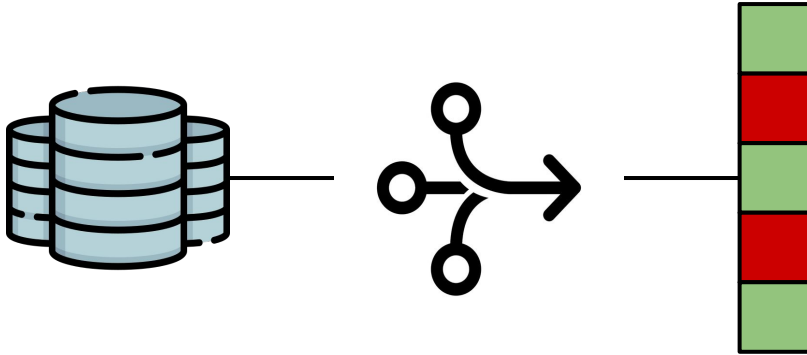
Developers Assessments

3	4	2	4	5
1	3	2	1	1
4	1	5	4	5
5	1	1	1	3
4	3	2	5	5

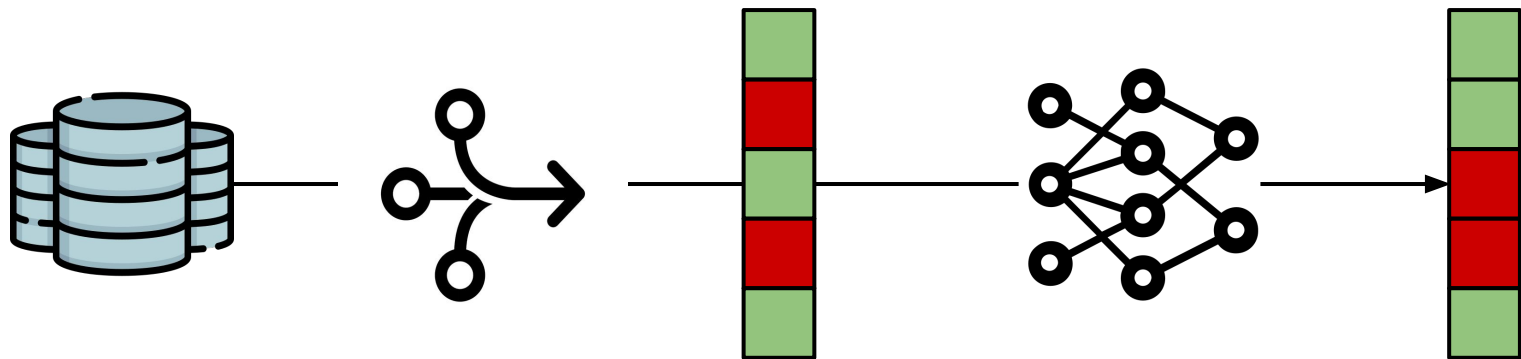
Snippets



CODE READABILITY



CODE READABILITY



CODE READABILITY

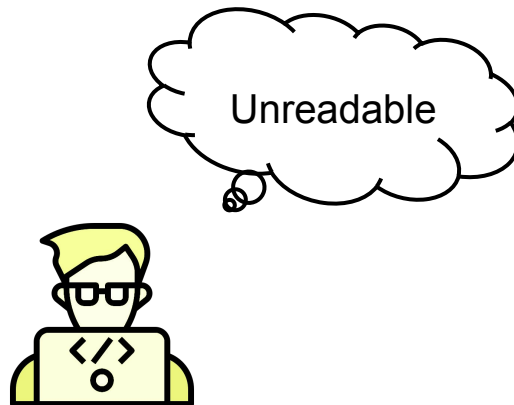
Readability is indeed subjective



CODE READABILITY

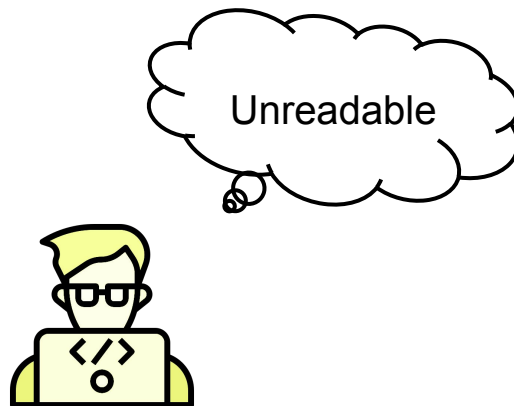
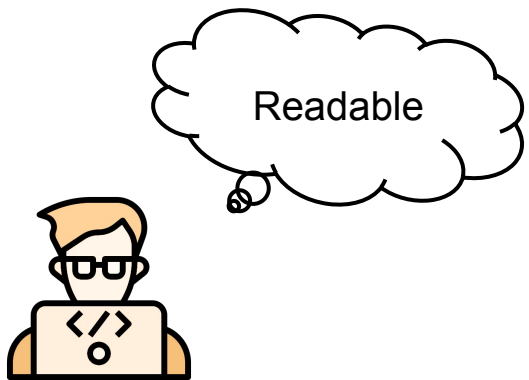
```
xsp = jj_scanpos;  
if (jj_scan_token (100) ) {  
  jj_scanpos = xsp;  
  if (jj_scan_token (101) ) return true;  
}
```

CODE READABILITY



```
xsp = jj_scanpos;  
if (jj_scan_token (100) ) {  
  jj_scanpos = xsp;  
  if (jj_scan_token (101) ) return true;  
}
```


CODE READABILITY



```
xsp = jj_scanpos;  
if (jj_scan_token (100) ) {  
  jj_scanpos = xsp;  
  if (jj_scan_token (101) ) return true;  
}
```

CODE READABILITY

To what extent is it possible to **assess code readability** as **subjectively** perceived by developers?

CODE READABILITY

Developers Assessments

Snippets

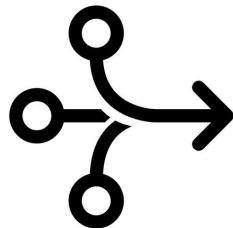
3	3	1	4	5
3	1	3	3	4
2	3	3	2	3
1	3	1	3	2
4	4	5	5	4

CODE READABILITY

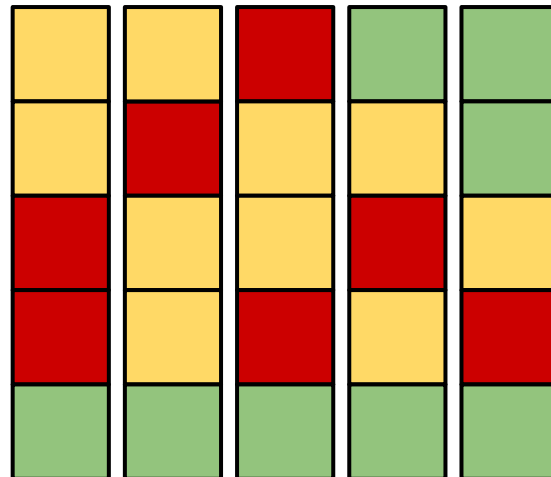
Developers Assessments

Snippets

3	3	1	4	5
3	1	3	3	4
2	3	3	2	3
1	3	1	3	2
4	4	5	5	4



Developers Assessments



Unreadable - {1, 2}



Neutral - {3}



Readable - {4, 5}

RQ1

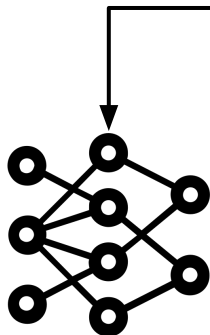
How effective are **generalist models** in predicting developers' subjective assessments of code readability



RQ1

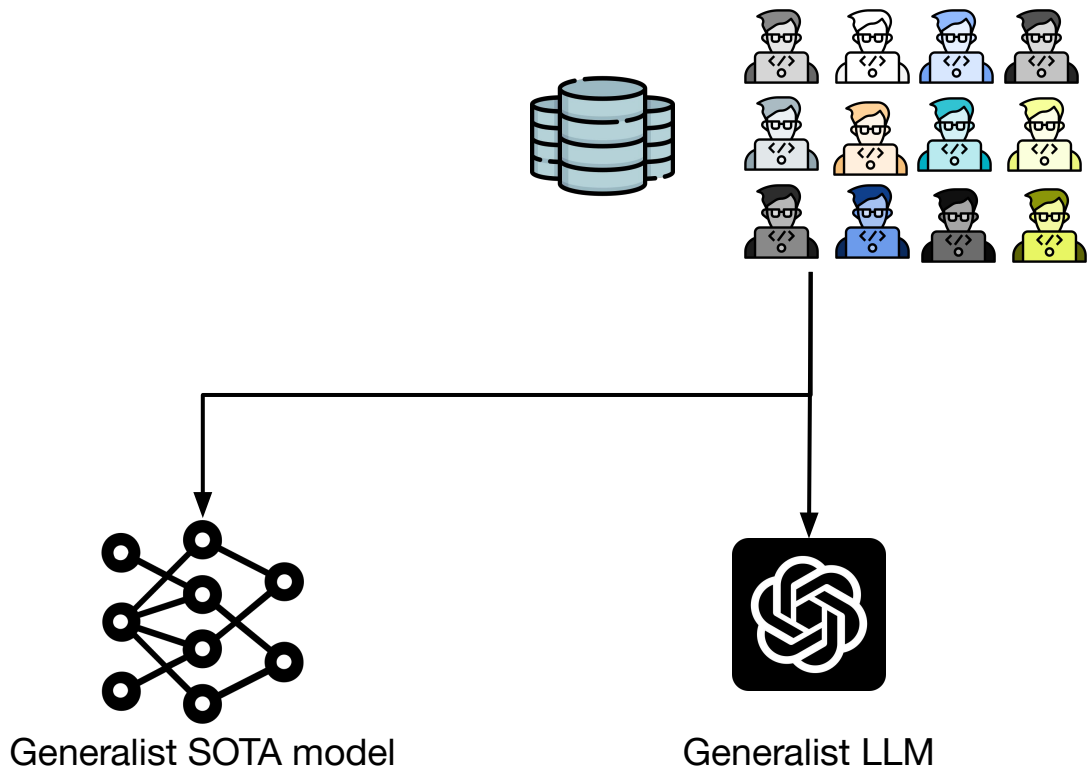


RQ1

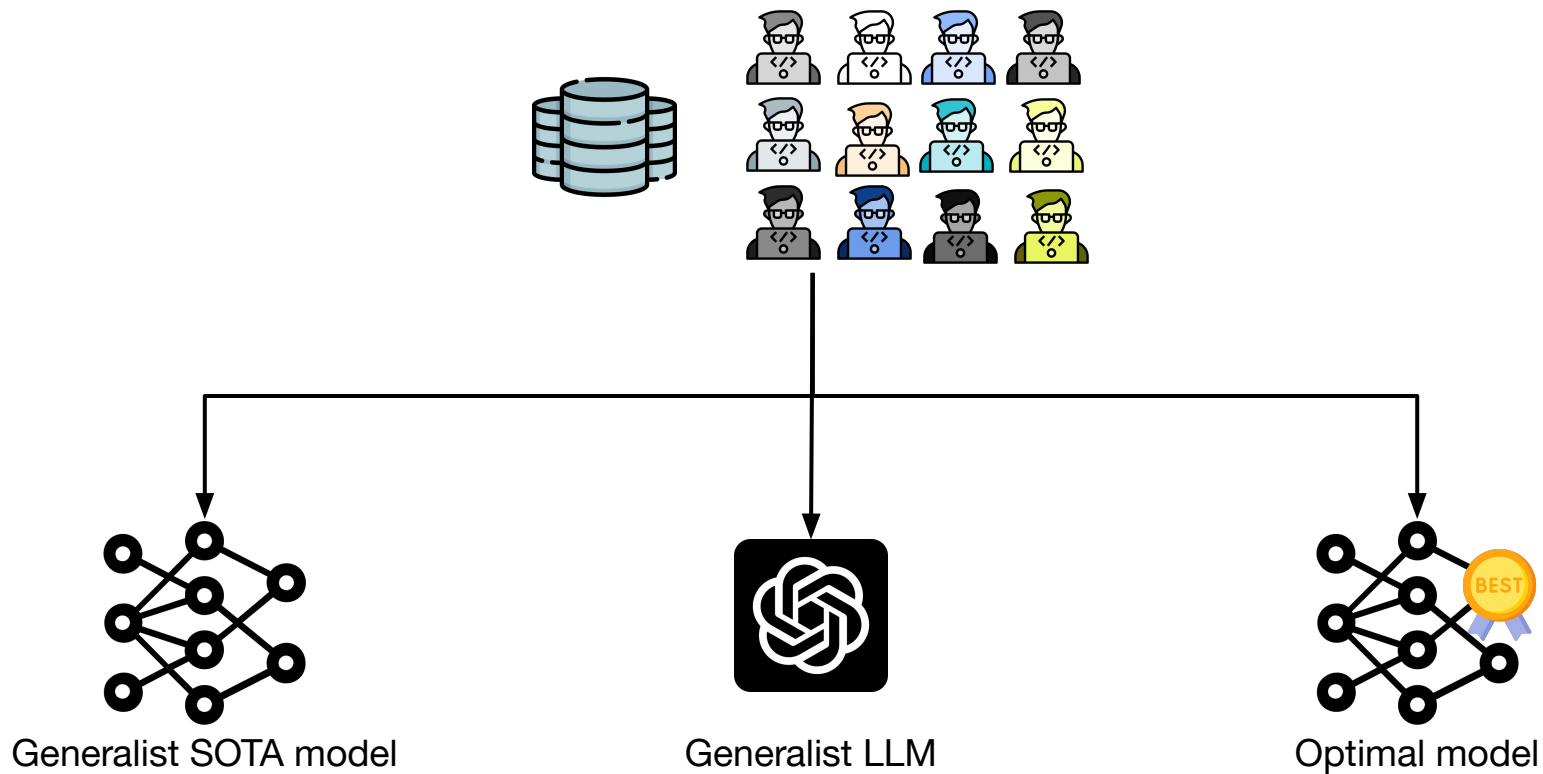


Generalist SOTA model

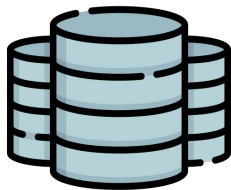
RQ1



RQ1

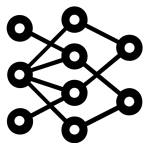


RQ1



D_s

F1-Score



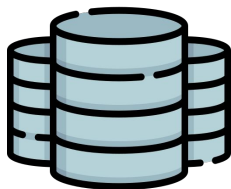
0.21



0.14

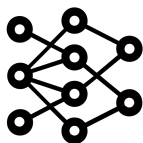


0.38



$D_{b\&w}$

F1-Score



0.40

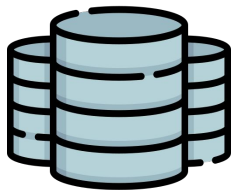


0.20



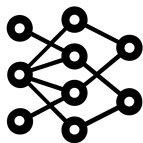
0.45

RQ1



D_s

F1-Score



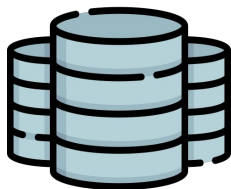
0.21



0.14

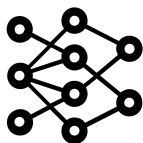


0.38



$D_{b\&w}$

F1-Score



0.40

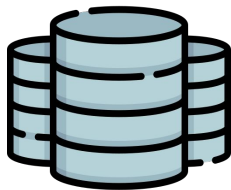


0.20

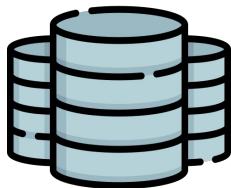
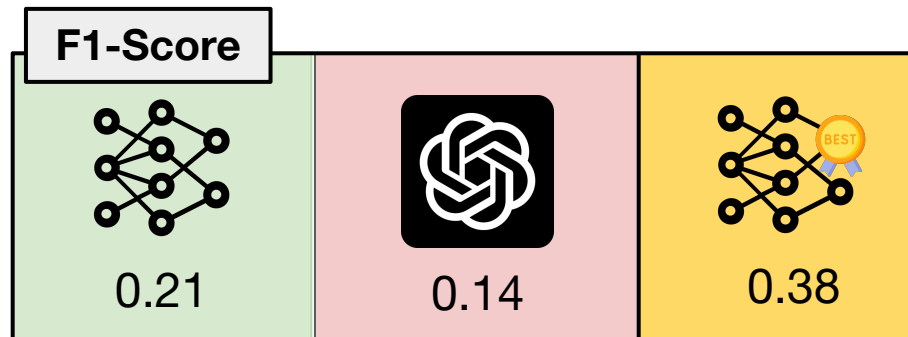


0.45

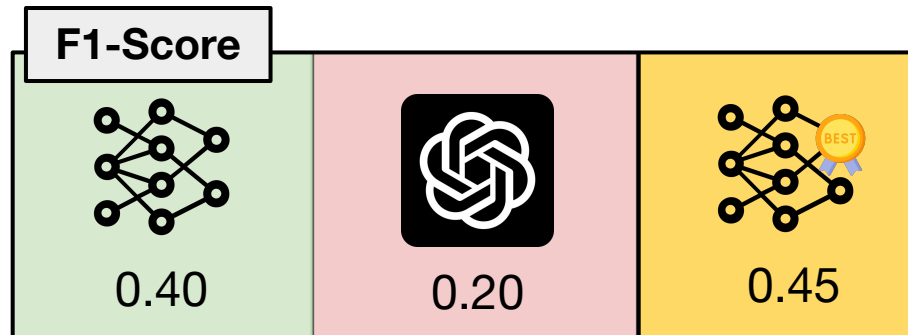
RQ1



D_s



$D_{b\&w}$



RQ1



Generalist models are not effective for personalized code readability assessments.

RQ1



Generalist models are not effective for personalized code readability assessments.



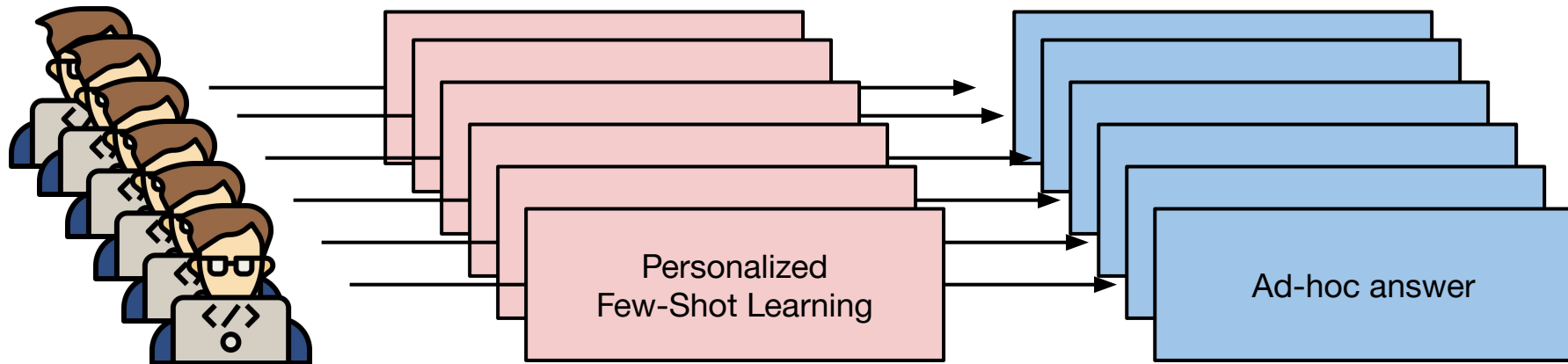
Generalist LLM achieves lower performance than a significantly smaller, feature-based model tailored for readability predictions.

RQ2

How **effective** is an LLM-based **personalized** code readability assessment model in predicting developers subjective assessments of code readability?



RQ2



RQ2

MODEL INPUT

You are an expert and personal code readability labeler.

Your role is to assign a code readability label based on the already known preferences of the developer. The labels you can assign are: *Unreadable*, *Neutral*, and *Readable*.

Below you find examples of the already known developer preferences.

```
0-shot  
1-shot  
2-shot
```

Now, assign the developer code readability label for the following one:

```
target-snippet
```

MODEL

Readable

Neutral

Unreadable

RQ2

MODEL INPUT

You are an expert and personal code readability labeler.

Your role is to assign a code readability label based on the already known preferences of the developer. The labels you can assign are: *Unreadable*, *Neutral*, and *Readable*.

Below you find examples of the already known developer preferences.

```
0-shot  
1-shot  
2-shot
```

Now, assign the developer code readability label for the following one:

```
target-snippet
```

MODEL

Readable

Neutral

Unreadable

RQ2

MODEL INPUT

You are an expert and personal code readability labeler.

Your role is to assign a code readability label based on the already known preferences of the developer. The labels you can assign are: *Unreadable*, *Neutral*, and *Readable*.

Below you find examples of the already known developer preferences.

```
0-shot  
1-shot  
2-shot
```

Now, assign the developer code readability label for the following one:

```
target-snippet
```

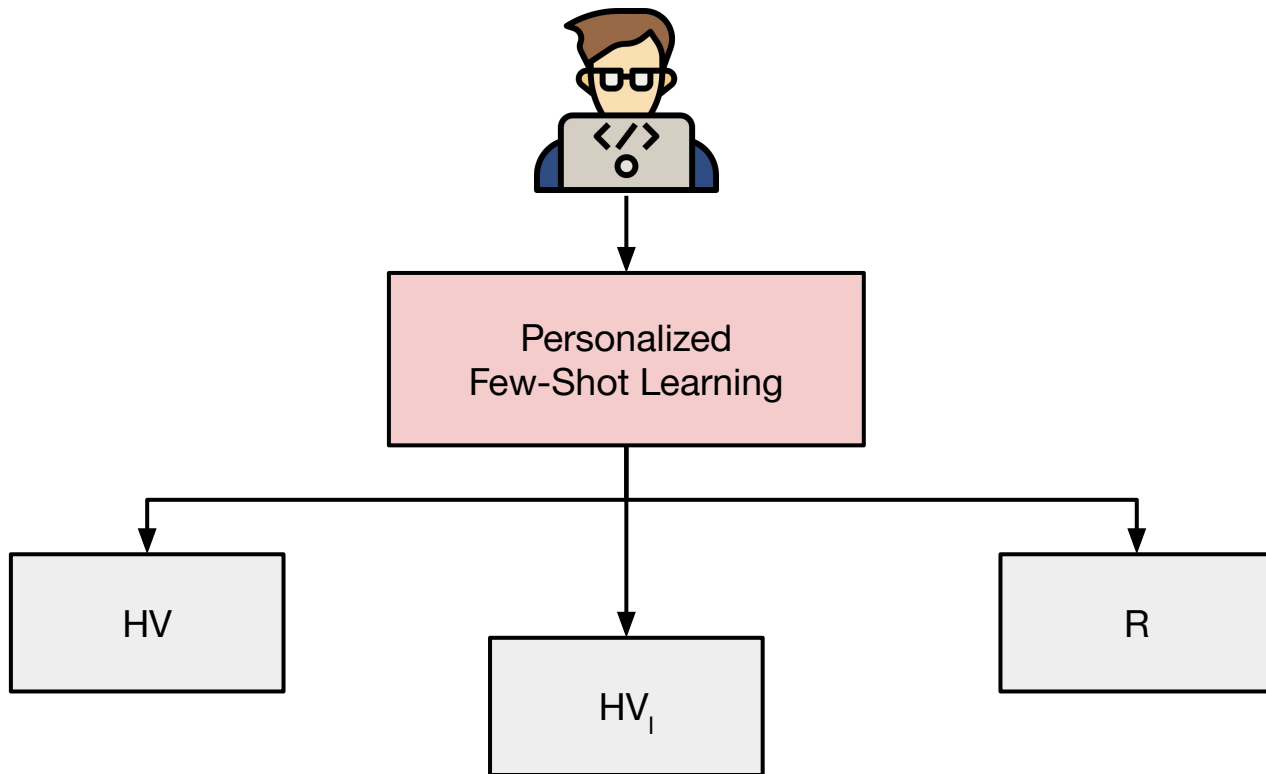
MODEL

Readable

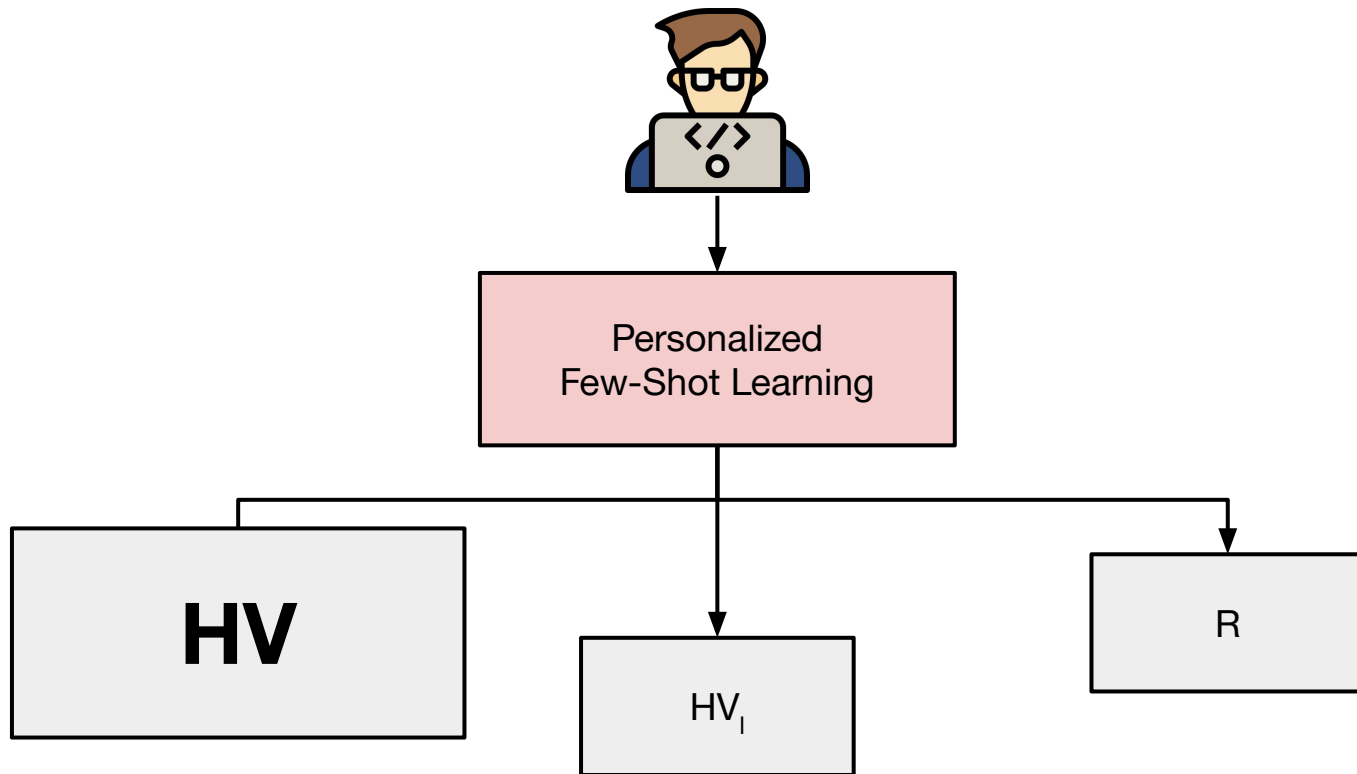
Neutral

Unreadable

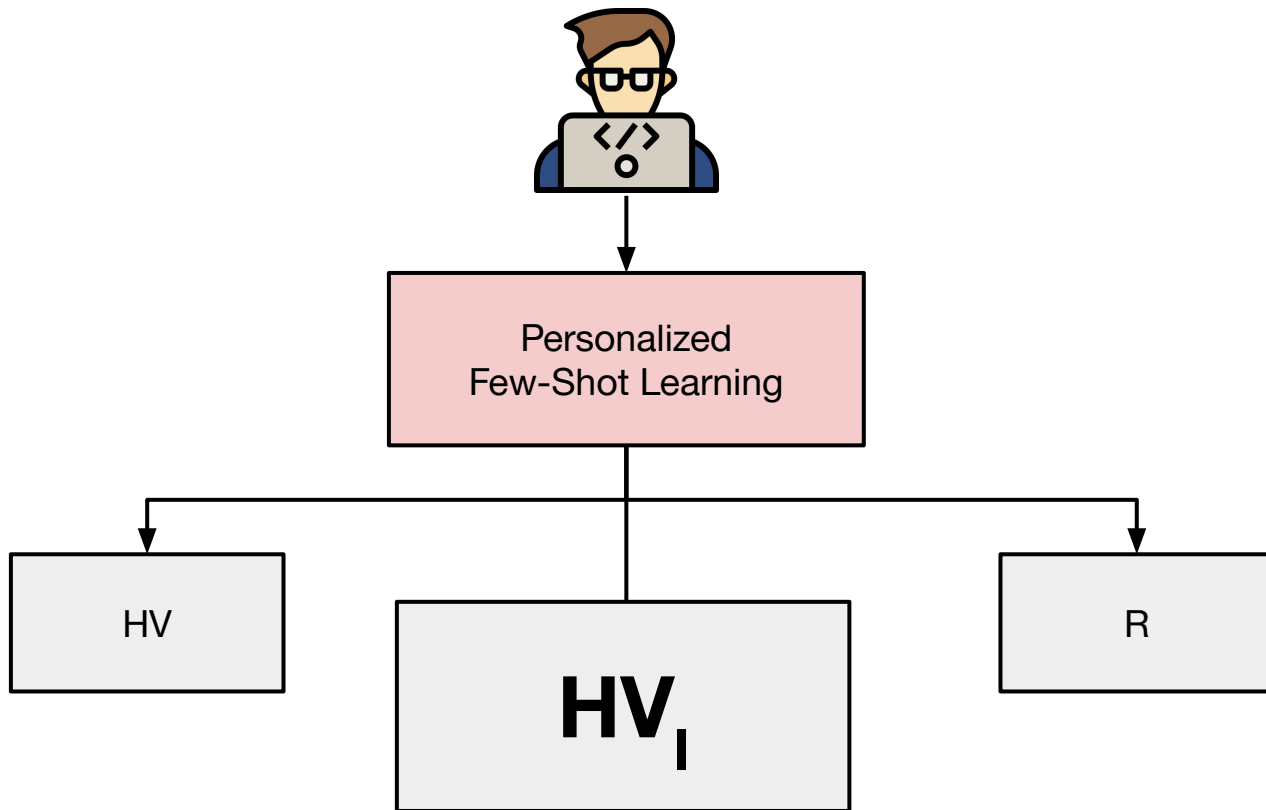
RQ2



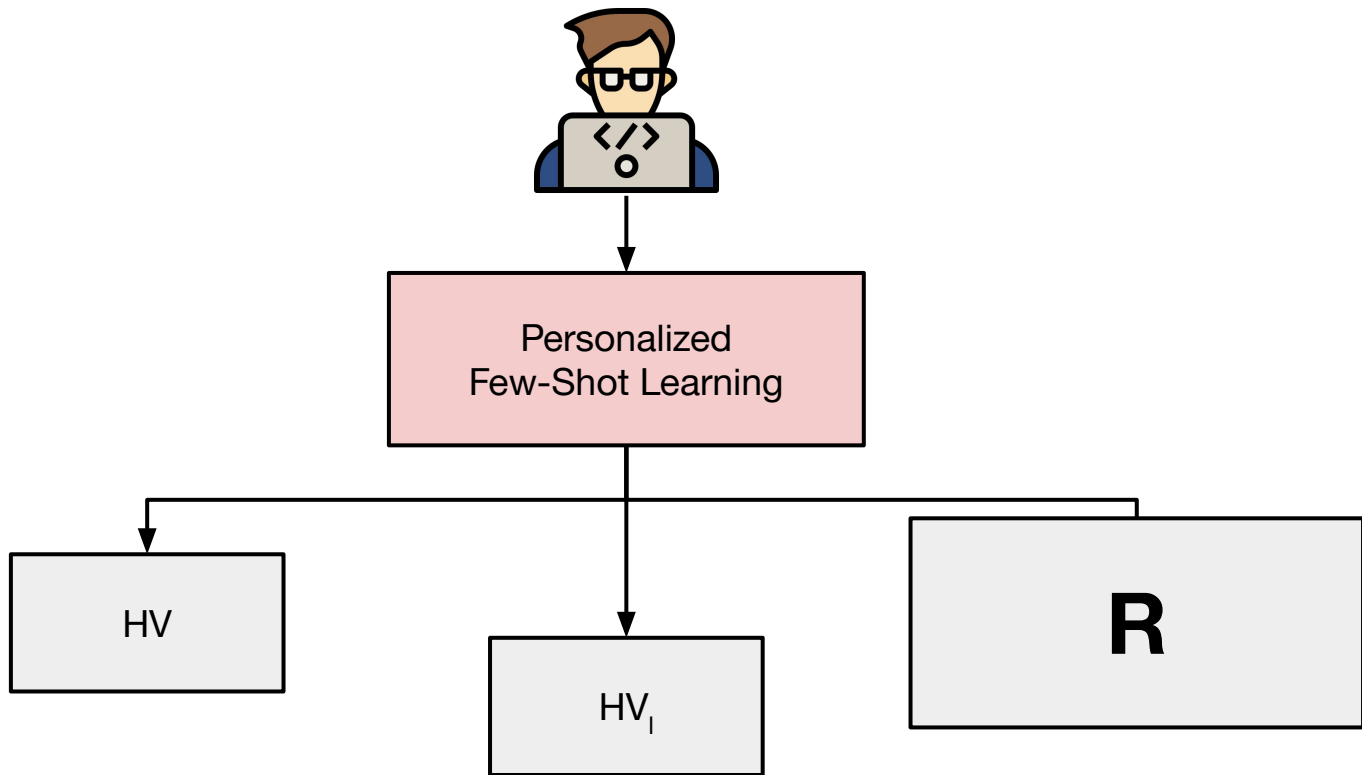
RQ2



RQ2



RQ2

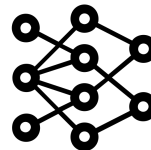
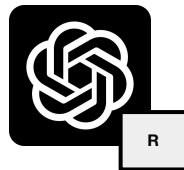
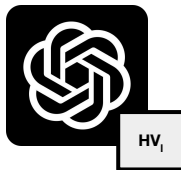


RQ2

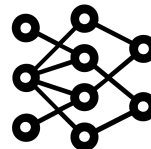
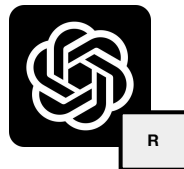
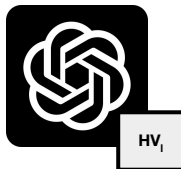
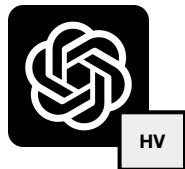


D_s

F1-Score



F1-Score



$D_{b\&w}$

RQ2



D_s

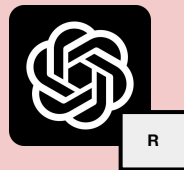
F1-Score



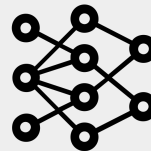
0.21



0.19



0.19

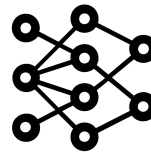
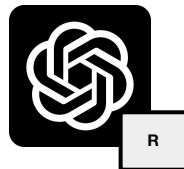
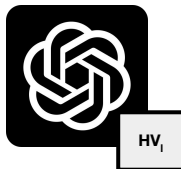
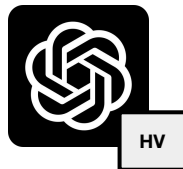


0.21



$D_{b\&w}$

F1-Score



RQ2



D_s

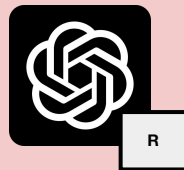
F1-Score



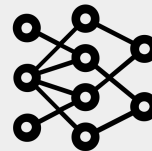
0.21



0.19



0.19

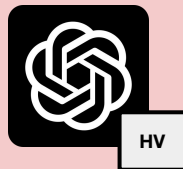


0.21



$D_{b\&w}$

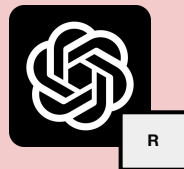
F1-Score



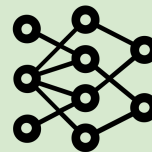
0.23



0.20



0.23



0.47

RQ2



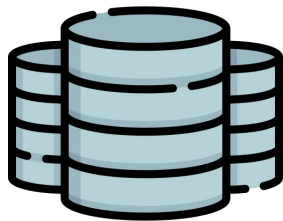
Personalized LLM-based code readability models are **less effective** than a state-of-the-art generalist model.

RQ3

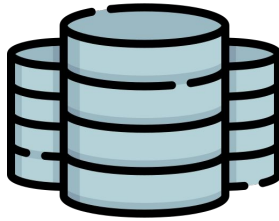
To what extent are developers' code readability assessments **consistent**



RQ3



$D_{b\&w}$



D_s



384 samples



384 samples

RQ3



Developers code readability assessments show moderate consistency, with **32% inconsistency in $D_{b\&w}$** and **23% in D_s** .

OUTCOMES



The definition of
**personalized code
readability prediction
models** is worth future
investigation.

OUTCOMES



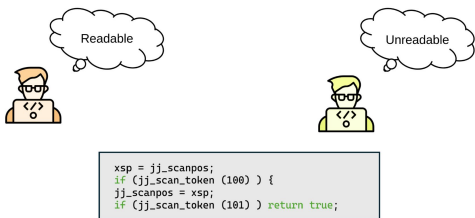
The definition of **personalized code readability prediction models** is worth future investigation.



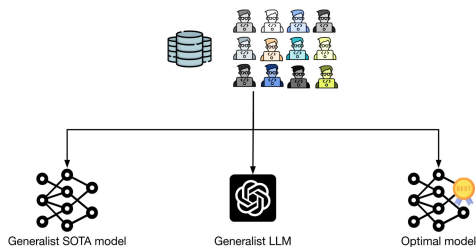
Future work should aim at using more principled procedures for defining code readability datasets.

SUMMARY

CODE READABILITY



RQ1

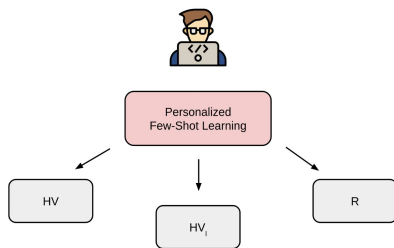


RQ1

Generalist models are not effective for personalized code readability assessments. ⚠️

Generalist LLM achieves lower performance than a significantly smaller, feature-based model tailored for readability predictions. ⚠️

RQ2



RQ2

❌ Personalized LLM-based code readability models are **less effective** than a state-of-the-art generalist model.

RQ3

⚠️ Developers code readability assessments show moderate consistency, with **32% inconsistency in $D_{b&w}$** and **23% in D_s** .