# Improving Code Readability Models with Textual Features

Simone Scalabrino, Mario Linares-Vásquez, Denys Poshyvanyk, Rocco Oliveto



### **Software maintenance** accounts for



### of the costs of a project



```
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:""));
}
```

# What does it do?

# Not really clear...



# What about this one?

```
for (int i=1;i<=100;i++) {
   String fizzBuzz = "";
   if (i % 3 == 0)
      fizzBuzz += "Fizz";</pre>
```

if (i % 5 == 0)
 fizzBuzz += "Buzz";

```
if (fizzBuzz.isEmpty())
  fizzBuzz += i;
```

```
System.out.println(fizzBuzz);
}
```

Why is it easier to understand?

# Code readability

### THE SPECIAL ISSUE ON THE ISSTA WAS BEET BADEDO

### Learning a Metric for Code Readability

Raymond PL Ruse Westley Weimer

Abstract-In this paper, we explore the concept of code readability and investigate its relation to software quality. With data collected from 120 human annotators, we derive associations between a simple set of local code features and human notions of readability. Using those features, we construct an automated readability measure and show that it can be 80% effective, and better than a human Using these features, we construct an automated readability measure and show that it can be 80% effective, and before than a human on average, at greeding readability alguest. Furthermore, we show that this meric constraists stoogly with here measures of software quality: code changes, automated defect reports, and defect log messages. We measure these correlations on owe 2.2 million lines of code, as well as ingulicationally, ore may relates or beloted projects. Fund, we discuss the indication of this study on programming language design and engineering practice. For example, our data suggests that comments, in of themselves, are less important than simple tank lines to calculatories.

its sequencing control (e.g., he conjectured that goto

Index Terms-software readability, program understanding, machine learning, software maintenance, code metrics, FindBugs

### **1** INTRODUCTION

We define readebility as a human judgment of how program is related to its maintainability, and is thus a protoch to system design [9]. We define readebility and is thus a structure of the protoch to system design [9]. employed that notion to help motivate his top-down approach to system design [9]. We present a descriptive model of software readability factor in overall software quality. Typically, maintenance factor in overall software quality. Typically, maintenance will consume over 70% of the total lifecycle control lifecycle control and pe strategies automati-based on simple fostius bill and personal automati-simple fostius bill and personal autom code readability and documentation readability are both critical to the maintainability of a project [1]. Other westernal (widel) variable) roisons of software quality, researchers have noted that the act of reading code is the most time-consuming component of all maintenance advirtuis [8], [33]. [53]. Readability is so significant, of advirtuis age (13), [53]. Readability is so significant, of advirtuis age (13) and between the use of read-bility matrice in particel harvess. The Bierk-Kinesid acurumes [8], [15], [35], Readability is as significant, in *u* software readability is useful, consider the use of read-fact, that Elshoff and Marcotty, after recognizing that a billy metrics in natural languages. The Flesch-Kincald many commercial programs were much more difficult for the second sec to read than necessary, proposed adding a development to read than necessary, proposed adding a development phasie in which the program is made more readable [10]. The program is made more readable [10] and the additional build be a check of the source cost Knight and Myers suggested that one phase of soft ware inspectime hould be a check of the source cost of the sour for readability [22] to ensure maintainability, portability, quite useful in practice. Flesch-Kincaid, which has been in use for over 50 years, has not only been integrated into popular text editors including Microsoft Word, but and reusability of the code. Haneef proposed adding a dedicated readability and documentation group to the development team, observing that, "without established has also become a United States governmental standard. Agencies, including the Department of Defense, require and consistent guidelines for readability, individual re-viewers may not be able to help much" [16].

viewers may not be able to help much<sup>2</sup> [16]. We hypothesize that programmers have some intra-titive notion of this concept, and that program features such as indentation (e.g., as in Pytohn [40]), choice (DOD MIL-435788). Determining a some com-identifier names [34], and comments are likely to play a part. Dijktra, for example, claimed that the readability of a program depends largely upon the simplicity of a program depends largely upon the simplicity of

### and Weinor are stilt de Department of Camputer Science at seni. We believe that similar mettres, tangeres operaneurs, ar de Surgin, Charlestenik, UA 22404. The anti-the senior senior for effectiveness, cam second-witter monitorie for the order for effectiveness, cam second-witter monitories for the order of the order. Linkning V Hysini, Chenttorowi, vo Lerse, Eusti, Hou, envirosi lovaropiantal This needado sua supportal in part hybe trang an effect in position of Arband Science That Science and Arbandian Cannis (CNOTHST and CNS 90021). And Markan Hyber 1999. A second seco

We describe the first general readability metric for source code

A Simpler Model of Software Readability

Daryl Posnett Abram Hindle Prem Devanbu University of California, Davis University of California, Davis Davis, CA Davis, CA Davis, CA Davis, CA devanbu@ucdavis.edu

ABSTRACT

General Terms

Keywords

### 1. INTRODUCTION

### A General Software Readability Model

as well as integ agree with each other and better than previous well. The falling integrate of the second s

 add is more likely to generalize.
 ceasaromary are retarroutly recent, first proposed by Buse et al. [13] and refined by Posnett et al. [14]. Such models are not coding standards (cf. [15]) but are based on combinations

 I. INTRODUCTION
 not coding standards (cf. [15]) but are based on combinations

Categories and Subject Descriptors

### Learning a Metric for Code Readability

### 1 INTRODUCTION

W E define readability as a human judgment of how We summe remainstrip as a human judgment of how unnecessarily complicates program understanding), and employed that notion to help motivate his top-down program is related to the maintainability, and is thus approach to system design [9].

### A Simpler Model of Software Readability

Daryl Posnett Abram Hindle Prem Devanbu University of California, Davis Davis, CA dpposnett@ucdavis.edu ah@softwareprocess.es devanbu@ucdavis.edu

ADS1KAC1 Softwar roadshifty is a property that influences how eas-ily a given piece of code can be road and understood. Since readshifty can affect maintainability, quality, etc., program-mers are very concerned labels the readshifty of code. If automatic readshifty, checkers to cold be built, heye could be integrated into developera tool cold be built, heye could be integrated into developerat tool-chains, and thus con-tinually inform developerat about the readshifty level of the Initially intern acevicepers about the reachanuity sevies of the code. Unfortunety, readability is a subjective code prop-erty, and not amenable to direct automated measurement. In a recently published study, Buse *et al.* asked 100 partici-pants to rate code snippets by readability, yielding arguably reliable mean 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 size and code entropy and show how this theory leads to a much sparser, yet sta-tistically significant, model of the mean readability scores produced in Buse's studies. Our model uses well-known size metrics and Halstead metrics, which are easily extracted us-ing a variety of tools. We argue that this approach pro-vides a more theoretically well-founded, practically usable. approach to readability measurement.

### Categories and Subject Descriptors

H.4 [Information Systems Applications]: Miscellaneous: D.2.8 [Software Engineering]: Metrics—complexity mea-sures performance measures

General Terms

ABSTRACT

Human Factors, Theory, Measurement

Keywords Readability, Halstead, Entropy, Replication

Permission to make digital or hard copies of all or part of this work for personal or classroom see is granted without fee provided that copies are to make of adiustrated for profits or conversional advantage and that copies republich, to part on servers or to redistribute to hits, requires prior specific permission and/or fee. ASSF '11, 21, AMX-2017, Waidak, Homoha, USA (Copyright 2011 ACM/97-14050, 0571-10165, -51006,

1. INTRODUCTION AT RODUCTION Readability of code is of central concern for developers [1, 18, 19, 24]. Code that is readable is often considered more maintainable: code that is more readable today is presumed to remain easier to read, comprehend, and maintain at a later data.

maintainside: code that is more readable today is pressmed to remain easies to read, comprehend, and maintain at a lab. To its. The relationship between readablity and un-programmers have of the difficulty of code, as they try to understand it. The relationship between readablity and un-erstanding is analogous to syntactic and semantic analysis, readablity is the syntactic appet while understandablity is barrier to understanding that the programmer fields the need to evercome before working with a body of code: the more translable it is the locer the barrier. The combening of the main state of the syntactic appet with studies of readablity is is the difficulty of experimentally residential experimen-tial syntamic studies and the syntactic studies of the synta subjective perception. Measures of subjective per-eption are body of experimentally independent studies of working multiple human rates, and careful attaitial anal-volving multiple human rates, and careful attaitial anal-volving multiple human rates, and careful attaitial and the area. Hey conducted a fully happensed study, akking readablity of code snippets. These scores were validated and aggregated to tyde time-consuming study was a set of code-snippets, 2, accound wide by unan a study-ied of the scettarious and time-consuming study was a study code-snippets, 2, accound wide by unan a shapettre readablity of code-snippets, 2, accound wide by unan a shapettre readablity of code-snippets, 2, accound wide by unan a shapettre readablity of code-snippets, 2, accound wide by unan a shapettre readablity of code-snippets, 2, accound wide by unan a shapettre readablity of code-snippets, 2, accound wide by unan a shapettre readablity of code-snippets, 2, accound wide by unan a shapettre readablity of code-snippets, 2, accound wide by unan a shapettre readablity of code-snippets, 2, accound wide by unan a shapettre readablity of code-snippets, 2, accound wide by unan a shapettre readablity of code-snippets, 2, accound wide by unan a shapettre readablity of code-snippets, S, accompanied by mean subjective readabil-ity scores: O(s),  $s \in S$ . Buse *et al.* then gathered direct. my scores: O(a),  $s \in S$ . Indie et al. Inen gathered affect, automatically-derived, token-level measures,  $\mathcal{M}(a)$ , of the code anippets. They then built a logistic regression model to predict the subjective, laboriously gathered scores, O(a), using the collection of automatically gathered direct metrics,  $\mathcal{M}(a)$ . The  $\mathcal{M}(a)$  were essentially token-level measures, but even so, were able to predict O(s) to some degree. This was an important contribution since it opens up the possibility of automatic tools that can provide readability score as feedback to developers for every commit they make: this

as resonance to developers for every commit usey make; this continuous feedback might potentially improve readability, and thus maintainability, of code over time. In this paper we improve upon Buse *et al.*'s model, yield-ing a model that is *simpler*, *better performing and theoreti*cally well-founded in both classical software engineering and basic information theory. In particular we argue that funda-mental aspects of readability have been measured since the 1970s by Halstead's software science metrics [21]. We sho that the Halstead metrics can be used to improve upon the

### A General Software Readability Model

as well as they agree with each other and better than pervisor work. We identify universal features of resultability and languages works with the dominant of software, formall metrics for a schernal notion of defect density. We address multiple program in order of magnitude more participants than previsors work, all gargesting our models in more Back by generating.

 Jel is more likely to generalize.
 e-consumity are treativity recent, first proposed by Buse e el. [13] and refined by Posnett et al. [14]. Such models are not coding standards (cf. [15]) but are based on combinations

### Learning a Metric for Code Readability

### 1 INTRODUCTION

We summe remainstry as a human judgment of how unnecessarily complicates program understanding), and employed that notion to help motivate his top-down approach to system design [9].

W is define readability as a human judgment of how

### A Simpler Model of Software Readability

Daryl Posnett Abram Hindle Prem Devanbu University of California, Davis University of California, Davis Davis, CA Davis, CA Davis, CA Davis, CA devanbu@ucdavis.edu

ABSTRACT

Keywords

Categories and Subject Descriptors

### 1. INTRODUCTION

### A General Software Readability Model

Jonathan Dorn Department of Computer Science University of Virginia Charlottesville, Virginia iad5iu@vireinia.edu

Abstract—We present a generalizable formal model of soft-ware readiability based on a human study of 5000 participants. are commoly used in commercial software and policies. Menotype of the software and policies with the software software and policies. Such as the code singlets. By contrast, we approach code as read on screens. For example, Flexi-Kincaid by humans and propose to analyze varia, spatial and linguistic is integrated into popular editors such as the to manufacture projects to almyter totals, spatial and highested provide the standard projects to almyter totals, spatial and highested highest provide the standard with the standard standard with the US Department on these notions and show that it agrees with human judgments as well as the syster with each other and better than provide the standard st a volta i they agree with cach other and better than previous over, We identify universal features of readbility and languages or experimence-specific ones. Our metric also correlations with any external notion of defect density. We analysis and any extension of the domain of software, formal metrics to translability are well-established in particular domains such as byperices (12) correlations and the second second second second second an order of magnitude more participants than previous work, and and effect density. We appendix the previous work, and an order of magnitude more participants than previous work, and and effect density we are realized as and the second second second second second and effect density. We appendix the previous work, and and effect density. We appendix the previous work, and and effect density we are realized by the second second second and (13) and refined by Poorent et al. [14]. Second consisting such as operative counts of an order of magnitude on experiments were appendix to the second second and (13) and refined by Poorent et al. [14]. Second consisting such as operative counts of an order of magnitude on experiments were appendix to the second second

 I. INTRODUCTION
 for coung standards (c), [15] out are reased on communities
 of surface-level synactic features such as operator counts or
 Modern software developers spend more time maintaining
 line lengths, aim to agree with human judgments, and have
 and evolving existing software than writing new code [1], [2],
 been found to correlate with external notions of software
 [3]. Software readability, a fundamental notion related to the quality [16]. Such software readability nodels do not attempt (2). Solvate remaining include to the comprehension of text, is critical to software maintenance reading code is a necessary first step toward maintaining it. Much research, both recent and established, has argued that on readability as a controllable accidental complexity [18]. known example is Knuth, who viewed readability as essential readability previous readability metrics do not adequately gento his notion of Literate Programming [4]. He argued that a eralize. They are based on small (typically 7-line) snippets of

to his notion of Literate Programming [4]. He argued that a emfize. They are based on small (typically 7-line) snippets of program should be viewed as "a project of literature, addressed code from a snighe programming language, are tited to shallow to humon beings" and that a readable program is "more robat, and more costable, and more costable and more solutions and more possible, ling more casily antamineds". Hancef transmissioner and the snight programming language, are tited to shallow or innove possible, ling more casily annuanced". Hancef transmissioner and the snight programming language, are tited to shallow of a development group dedicated to readability and relatively small number of statems [13], [14]. We propose a documentation: "windue tablibles and consistent guidage that a source state is readability inducid that a description of the abbe to below remaining lightweight and apprecises of systax highlighting sug-check for readability induces a provide that a source that is supported that a source constrained in the snipport of the snipport o erance ucevapament passe anime as improving reasoning was inconnects, we imis propose the init incorporation or geomet-proposed by Elshöff and Marcoty, who observed that many rice, pattern-based and linguistic aspects and features into an commercial programs were unnecessarily difficult to read [8]. automated readability metric. For example, code in which the More recently, a 2012 survey of over 100 developers and "=" operators in a sequence of assignment statements "line managers at Microsoft by Buse and Zimmermann found that up" verically on the screen may be viewed as more readable, weight or statistically out the scientific and the science of t non-software natural language. Metrics such as the Automated readability.

### Learning a Metric for Code Readability Raymond PL Ruse Westley Weimer

A Simpler Model of Software Readability

Abram Hindle ersity of California, Davis Davis, CA 2softwareprocess.es is edu ah@soft

1. INTRODUCTION

### Structural features

es and Subject Descriptors

ABSTRACT

Keywords

### A General Software Readability Model

### Learning a Metric for Code Readability

### A Simpler Model of Software Readability

Abram Hindle Prem Devanb ersity of Galifornia, Davis Davis, CA Davis, CA essies devanbi

1. INTRODUCTION

### Structural features

Visual features

### A General Software Readability Model

### Learning a Metric for Code Readability

### A Simpler Model of Software Readability

Abram Hindle ersity of California, Davis Unive Davis, CA

Prem Devanb

1. INTRODUCTION

### Structural features

Visual features

A General Software Readability Model

Two datasets

Something is missing...

# Code is text!

DIVINA MEDIA Gustavo Dore



# New features

# Comments readability

nealize

Ara

actos

ennalmente:

A Periodo 1860

# Comments and **identifiers** consistency

# Identifier terms in dictionary

# Narrow meaning identifiers



# Number of meanings

# Textual coherence



# Case study

# **200** Java snippets

# annotators

# New dataset

Do textual features complement the others proposend in the literature?

# **Overlap** metrics

### Textual Features vs Buse's



### Textual Features vs Posnett's



### Textual Features vs Dorn's



### Textual Features vs Dorn's



What is the accuracy of a readability model based on structural and textual features?

### Dataset by Buse and Weimer



## Dataset by Dorn



### New dataset



### New dataset



### In summary...













# Code readability for defect prediction

# Thanks.