

A Survey on the Usage of Eye-Tracking in Computer Programming

UNAIZAH OBAIDELLAH and MOHAMMED AL HAEK, University of Malaya
PETER C.-H. CHENG, University of Sussex

Traditional quantitative research methods of data collection in programming, such as questionnaires and interviews, are the most common approaches for researchers in this field. However, in recent years, eye-tracking has been on the rise as a new method of collecting evidence of visual attention and the cognitive process of programmers. Eye-tracking has been used by researchers in the field of programming to analyze and understand a variety of tasks such as comprehension and debugging. In this article, we will focus on reporting how experiments that used eye-trackers in programming research are conducted, and the information that can be collected from these experiments. In this mapping study, we identify and report on 63 studies, published between 1990 and June 2017, collected and gathered via manual search on digital libraries and databases related to computer science and computer engineering. Among the five main areas of research interest are program comprehension and debugging, which received an increased interest in recent years, non-code comprehension, collaborative programming, and requirements traceability research, which had the fewest number of publications due to possible limitations of the eye-tracking technology in this type of experiments. We find that most of the participants in these studies were students and faculty members from institutions of higher learning, and while they performed programming tasks on a range of programming languages and programming representations, we find Java language and Unified Modeling Language (UML) representation to be the most used materials. We also report on a range of eye-trackers and attention tracking tools that have been utilized, and find Tobii eye-trackers to be the most used devices by researchers.

CCS Concepts: • **Human-centered computing** → **Empirical studies in visualization**; *Visualization design and evaluation methods*; **Empirical studies in HCI**;

Additional Key Words and Phrases: Eye tracker, programming, empirical studies, HCI, debugging, comprehension, participants, test materials, eye-tracking metrics

ACM Reference format:

Unaizah Obaidellah, Mohammed Al Haek, and Peter C.-H. Cheng. 2018. A Survey on the Usage of Eye-Tracking in Computer Programming. *ACM Comput. Surv.* 51, 1, Article 5 (January 2018), 58 pages.
<https://doi.org/10.1145/3145904>

This work is supported by the University of Malaya and Ministry of Higher Education, Malaysia, under grants FP062-2014A and RP030C-14AET.

Authors' addresses: U. Obaidellah, Department of Artificial Intelligence, Faculty of Computer Science and Information Technology, University of Malaya, 50603, Kuala Lumpur, Malaysia; email: unaizah@um.edu.my; M. Al Haek, Department of Artificial Intelligence, Faculty of Computer Science and Information Technology, University of Malaya, 50603, Kuala Lumpur, Malaysia; email: alhaekmohammed@gmail.com; P. C.-H. Cheng, Department of Informatics, School of Engineering and Informatics, University of Sussex, Falmer, Brighton, BN1 9QH, United Kingdom; email: p.c.h.cheng@sussex.ac.uk. 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. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

© 2018 ACM 0360-0300/2018/01-ART5 \$15.00

<https://doi.org/10.1145/3145904>

1 INTRODUCTION

The lack of problem solving ability has been cited by many researchers as one of the main reasons students face difficulties in learning to write computer programs [29, 32, 33]. Based on a survey conducted on students and tutors, [43] associated the most difficult topics in learning programming to the lack of understanding of concepts and comprehension inability. In another attempt to understand the difficulties novice programmers encounter, [45] classified cognitive requirements of programming (i.e., cognitive demand, mental models, cognitive load), as one of the five main areas of difficulties in learning programming. Based on the Blooms Taxonomy [2] in relation to learning computer programming, [59] found that students have problems in all three domains of Blooms Taxonomy, with the majority of them having difficulties in the cognitive domain.

The relation between programming education and the cognitive process of students is, however, not a new finding. In 1990, [26] proposed the analysis of data collected from eye movement to investigate the cognitive process of programmers, based on theories that linked eye fixation to comprehension to understand how programmers read and understand algorithms. In their work, [26] used the Applied Science Laboratories (ASL) eye movement monitor to track the focused attention of the subjects. This work by [26] is one of the first attempts to study computer programming by tracking the eye movement of students. However, this approach was not adopted immediately by researchers in this field. In a review, [75] analyzed published research between 2005 and 2008 related to programming education, and found that most papers used quantitative research methods such as course assessments, questionnaires, and interviews. The traditional data collection methods adopted in most papers follows that of [75]. A limited number of reports used eye-tracking. In fact, eye-tracking was not in the top 12 data gathering techniques reported at that time. In 2013, [42] reviewed research using eye-trackers in studies related to learning, and showed the emerging use of eye-trackers in research related to learning.

Given the increasing use of eye-tracking in assessing the underlying cognitive processes of programming, this survey will provide a summary of existing work in this field. This is considered necessary due to the absence of a guideline or methodology specifically designed for this type of research. We decide to report this work as a mapping study to find linkages between the existing published studies, identify patterns of publication, and facilitate collection of literature retrieval in the field of computer programming using eye-tracking. In this mapping study, we will try to answer questions about eye-tracking experiments and its setup, and provide quantitative information from existing work, such as the type of materials used and the number of subjects recruited for the study. This work will also provide a wide overview of eye-tracking in programming, and try to identify areas suitable for conducting a Systematic Literature Review (SLR) with more appropriate primary studies. This survey can serve as a reference for researchers who are using eye-tracking to conduct a study related to learning and education, and it can provide detailed information for those who want to get started on research in similar fields using eye-trackers. This mapping study discusses an extensive list of research areas studied using eye-trackers in relation to programming, as to provide information to students and researchers who are planning on exploring this field. Furthermore, it can be helpful to eye-tracking manufacturers in terms of providing quantitative information about the eye-tracking devices used by researchers in this field.

In this research, we describe the details of papers collection in Section 2 and the results of the data collection in Section 2.4. Then, we report the mapping from the papers in Section 3, starting with the number of research that used eye-tracking, the classification of experiments, the different materials used, the participant selection and sample details, followed by types of eye-trackers used in each paper, along with the metrics and variables used in the studies. We present a discussion of

the mapping results in Section 4, then state any threats to the validity of our study in Section 5, and conclude our survey in Section 6.

Throughout this work, we will use the LaTeX generated BibTeX bibliography style *Alpha* [49] to reference the papers included in our study. The *Alpha* bibliography style makes use of the first initials of the last name of the authors along with the year of the publication. We decided on using this style since it helps in reducing the space of the citation, and makes reporting 63 papers easier to fit into a diagram or a table. For example, if we want to reference paper [26] published by Crosby and Stelovsky in 1990, we will use the *Alpha* bibliography style to refer to it as [CS90]. If two papers have the same initials, such as in the case of [14] and [15], both papers published by Bednarik and Tukiainen in 2007 would have the same *alpha* style of [BT07]. These papers will be distinguished by a letter following the year, and the order of the letters is based on the alphabetical order of the papers' title. Therefore, [14] will be [BT07a] and [15] will be [BT07b]. It is also worth noting that in all the tables and the figures, we present the papers of our mapping study sorted and organized by the year of the publication, not alphabetically.

1.1 Related Work

In the course of this research, we came across a recent SLR that reported on the use of eye-trackers in software engineering. In their research, [67] provided details of different experiments that used an eye-tracker to examine studies on software engineering. Upon starting our research, we noticed similarities between our work, and the work done by [67], as well as differences. Although most studies reported by [67] are included in our reporting, it is worthy to note that our work focuses more on the programming aspects of these studies. As [67] has provided detailed information on the calculations and formulas for the metrics of visual effort, we advise readers to refer to [67] for further details. We noticed some differences in the reporting between our survey and [67]. However, we acknowledge the level of details reported on the use of eye-trackers in software engineering studies. In their work, [67] provided detailed information on the background of eye-tracking technology and different setup of devices, along with visual descriptions of the technology. Again, we advise the reader to refer to [67] for detailed information on the metrics used for calculating and processing the collected data from the eye-trackers, as we will not be reporting these topics to avoid repetition. Instead, we will focus on the findings of the studies in relation to experimental setup, analyzing the materials, type of trackers, and the types of participants recruited in the studies.

2 METHODOLOGY

This survey is based on the guidelines suggested by [4] and [51] for mapping studies. In our work, the mapping study helps to determine the amount and scope of research that has been conducted in the area of computer programming using eye-trackers within a certain period of time. We estimate that the reported work guides research and practices in this area such as the selection of participants, types of test materials, types of and eye-tracker devices, types of information an eye-tracker can help gather, and the variables it can measure. In this survey, we try to answer the following research questions:

- RQ1: How many papers have used eye-tracking in research on computer programming?
- RQ2: What kinds of programming tasks and areas were explored using eye-tracking?
- RQ3: What types programming materials were used as stimulus in eye-tracking experiments on computer programming?
- RQ4: Who were the subjects of eye-tracking experiments in computer programming?
- RQ5: What eye-tracking technology and devices were used in computer programming experiments?

—RQ6: What eye-tracking metrics and variables are commonly reported in programming studies?

RQ1 will help identify the number of published works that used eye-trackers, thereby identifying how the technology has been adopted over the years, and provide evidence to its emergence in this field. RQ2 will help to identify and classify programming topics suitable for conducting a systematic literature review, as well as to provide a detailed classification of research attributes of the reported studies. RQ3, RQ4, RQ5, and RQ6 will answer questions related to the experimental setup, and provide detailed information about the sample of subjects and its size, type of materials and stimulus, eye-tracking device, and the metrics and variables that researchers need to consider prior to conducting a similar study.

2.1 Papers Collection and Selection

We focus on the collection of papers related to the use of eye-trackers in programming tasks, or programming-related topics such as algorithms and diagrams. We performed our search on electronic databases using multiple search queries (see below). Then, we went through the initial search results returned by each search engine, and selected papers that fit our aim by analyzing the abstract and keywords of each paper. Finally, we performed the Snowballing process. Each search query consisted of a variation of the words “*EyeTrack*,” followed by stimuli and/or a task. The details of the search queries are as follows:

- (1) Eye-tracking: In order to ensure that the search results are related to eye-tracking or eye-trackers, the first keyword in every search query was “eye track” or “eye-track.”
- (2) Stimuli: To make sure we find papers related to programming, we included terms such as “code” and “program,” as well as “UML” (Unified Modeling Language) and “pseudocode,” followed by a programming task.
- (3) Task: Since our main focus is to find papers related to programming, we included programming-related tasks into the search queries such as “comprehension,” “debugging,” “scan,” and “read.”

[4] suggested a list of electronic sources to consider for finding studies relevant to software engineers, including: IEEEExplore, ACM Digital library, SCOPUS, CiteSeer, ScienceDirect, and Springer. However, in an experience report, [30] reported on multiple returned similarities or no unique results returned from some digital libraries, stating that

“after performing the searches we found that we could have saved ourselves some work as none of the publisher-specific databases except the IEEEExplore and ACM Digital Library returned any unique “hits”. That is, all articles returned by Kluwer Online, ScienceDirect, SpringerLink, and Wiley Inter Science Journal Finder were also returned by either ISI Web of Science or Compendex.” ([30], p. 229)

This point was later echoed by [38], who stated:

“This is similar to the point made by [30]...they could have saved time and effort for general searches by using ACM, IEEE, plus two indexing systems rather than searching multiple publishers digital libraries.” ([38], p. 2063)

In an updated guidelines for mapping studies, [52] cited the recommendations of [30] and [38] in searching digital libraries, and stated that using IEEE, ACM, and two indexing systems such as Inspec/Compendex and Scopus is sufficient.

Table 1. Search Strings Used and the Number of Papers Returned

| Database | Search string | Number of papers |
|-------------|--|------------------|
| ACM | ((comprehension OR understand OR debug OR debugging OR scan OR read) AND code OR program OR programming) AND eye-tracker OR eye-tracking) | 74 |
| ACM | ((uml OR diagram OR pseudocode OR flowchart) AND eye-tracking) | 138 |
| IEEEExplore | ((comprehension OR understand OR debug* OR scan OR read) AND code OR program OR programming OR) AND eye-track*) | 83 |
| IEEEExplore | ((uml OR diagram OR pseudocode OR flowchart) AND eye-track*) | 12 |
| SCOPUS | ((comprehension OR understand OR debug OR debugging OR scan OR read) AND code OR program OR programming) AND eye AND tracker OR eye AND tracking)) | 163 |
| SCOPUS | (eye AND tracker OR eye AND tracking) AND (uml OR "class diagram" OR pseudocode OR flowchart) | 17 |

We conducted our electronic database search in two phases. Phase one was an initial search using IEEEExplore, ACM Digital library, and SCOPUS. This initial search phase will validate our keywords search selection, and help edit the search query if no adequate results were returned. The libraries or databases used in phase one were ACM, IEEEExplore, and SCOPUS, in that order. The searching string used for the digital libraries is shown in Table 1, where each search string was modified accordingly. For instance, IEEEExplore ignores most punctuation, and when searching for “Eye-track,” it also looks for “Eye track” and “Eye_track.” IEEEExplore also makes use of the wild card (*) at the end of a word to search for words with different endings. The search strings for other digital libraries were modified as well.

Phase two took place after the returned results from phase one were analyzed. Phase two depended on phase one, and it included either editing the search queries, or expanding the search to other recommended electronic databases. As our search queries and keywords selection returned sufficient results, we expanded our search to other digital libraries. For the second phase, we explored ScienceDirect, Springer, Web of Science, and Citeseer. However, no new results were retrieved from these sources, as all returned papers were already collected from the initial searching phase. For example, some of the results returned by conducting the keywords search on Web of Science were [Bed12, BDL+13, Dcl13, ASGA15], which were already retrieved through Scopus. Similarly, [SUGGR15, JF15] were retrieved through IEEEExplore and [FBM+14, RLMM15] were retrieved from ACM. As no new results were retrieved from the second phase search, we relied on the snowballing process to find additional papers. The snowballing process depended on exploring the references, suggestions from related work, and recommendations from Mendeley [35]. The papers selection process carried out is shown in Figure 1, where the number on the left is the number of papers extracted from each of the following steps:

Table 2. Classification Aspects for Data Extraction and Reporting of the Study Papers

| Aspect | Data Extracted |
|-----------------------------|---|
| <i>Extent</i> | Author(s) names, Year of publication, Journal or Conference |
| <i>Topic</i> | From title, keywords, abstract, problem statement and objective |
| <i>Subjects</i> | Details of the sample for the study, their affiliation, number, gender, experience and grouping |
| <i>Task and Environment</i> | The programming materials used, Eye-tracker and its configuration, eye-tracking metrics and variable of the study |
| <i>Replication</i> | Form subjects, topics, authors and full analysis |

- (1) Use of eye-tracking in programming or programming-related context.
- (2) Having been published in a journal, conference, or proceeding reporting the results of an experiment using eye-tracker.
- (3) Papers not specifically on computer programming and program code, but a related topic such as software engineering, and using stimuli other than source code such as UML diagrams and flowcharts, as included in the context of programming-related topics.

We excluded papers based on the following:

- (1) Papers not reporting the use of eye-trackers.
- (2) Papers using eye-trackers in a context not related to computer programming.
- (3) Papers not published in English.
- (4) Papers that did not go through a referring process, such as posters, work sessions, lecture notes, and dissertations.
- (5) Papers re-reporting the results, or doing a re-analysis of a previously published experiment were studied, and the most comprehensive paper was selected.
- (6) Other materials such as books, technical papers, government reports, letters and editorials, possession papers, and papers with an abstract but no full text available.
- (7) Papers not involving an empirical study or only those that propose a proof of concept.

2.3 Classification Scheme

We adopted the classification scheme used by [76], as it was cited and recommended by [51]. In their work, [76] classified controlled experiments in software engineering based on the aspects of *extent*, *topic*, *subjects*, *task and environment*, and *replication*. Given that we consider our work as a mapping study that reports experiments related to computer programming, we see this criteria to fit our objective. We found the aspects listed in Table 2 to be more fit to the nature of the work we are reporting.

The aspect of *extent* helped classify publication frequency and venues such as journals and proceedings. Additional information from *extent* also assisted the authors in identifying and classifying repeated experiments. As for *topic*, we identify and classify the programming area or task that was the focus of each study. We attended to the problem statement and objective of each study to accurately classify papers into topics such as debugging, collaborative programming, and comprehension; however, keywords and complete abstracts were mostly sufficient in classifying the programming topic or area each paper addressed. Details of the subjects affiliation, number, grouping, and other information were classified mostly for the purpose of identifying the targeted

audience of eye-tracking research in programming studies. As for the aspect of *task and environment*, we were able to identify and classify the materials used in each experiment, the type of eye-tracking technology utilized, and details of the metrics and variables each study tried to evaluate and examine. The last aspect of *replication* ensures that papers which re-reported or re-published the results of a previous experiment were not included. This classification of unique and repeated experiments ensures more accurate statistical reporting of information from the collected papers.

For the first two aspects: *extent* and *topic*, the data collected from title, publisher, year, abstract, and keywords were in most cases sufficient enough to categorize papers. However, for the other aspects of *subjects*, *task and environment*, and *replication*, a full analysis of each paper was required by one of the first two authors to accurately classify each paper. While one author extracted and recorded information from a selected paper and classified it, the other author confirmed the classification.

2.4 Data Collection Results

Figure 1 shows a total of 63 papers were selected out of 487 returned by the search engines we used, in response to our keywords search queries. Of these 63 papers, 16 were removed after we performed a full analysis. Then, 15 more were added from the Snowballing process. 13 papers that reported the results from the same experiment in different publications were excluded. The details of the excluded and included papers from Figure 1 are as follows:

- (1) Removed papers: 16 papers were not included for the following reasons:
 - (a) Papers that do not contain an eye-tracking study:
 - [34] presented a working session of a conference with no experiment.
 - [77] contained a preliminary idea to study the factors driving the cognitive process during program comprehension, including vision, and measure their impact using eye-tracking.
 - [65] contained an analysis of the visualization techniques used in software maintenance and code comprehension.
 - [71] presented a case for the use of eye-trackers in software traceability and what eye-tracking can contribute to the task. It did not report an experiment that used eye-tracker.
 - [53] inspected the comprehension factors of a business process model.
 - [37] proposed an evaluation process and method based on eyes course, to measure programming difficulties among students.
 - [41] investigated the possibility of using eye-tracking data and navigation paths in order to identify and find source code dependencies.
 - [48] introduced a fixation shift algorithm in an attempt to correct the location of a fixation automatically, in cases of fixation drift or calibration inaccuracy.
 - [83] discussed the use of eye -tracking technology in programming and the relation between eye-tracking data and cognitive load in programming.
 - (b) Papers on the use of systems or tools:
 - [79] presented a Computer-Aided Empirical Software Engineering (CAESE) framework and introduces the Ginger2 environment, which was developed based on the CAESE framework, and has a variety of data collection and analysis tools including audio and video recorder, monitoring tools, and eye-tracking.
 - [82] designed and implemented the *DRESREM* system, a single-document review evaluation system, capable of measuring and recording eye movements of reviewers during a document review task.

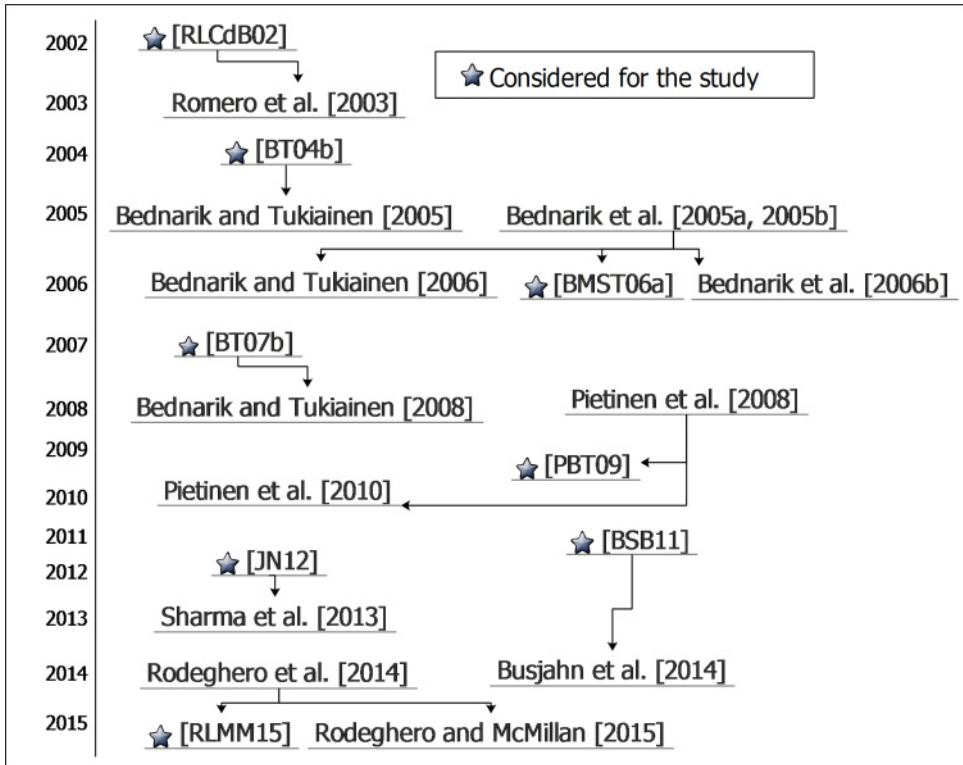


Fig. 2. Repeated experiments.

- [81] reported on an enhanced version of the *DRESREM* system to a multi-document review evaluation system named *DRESREM2*.
 - [17] reported on the research and development of a program animation system named Jeliot.
 - [40] examines the possibility of using eGlasses as an eye-tracking tool. It focused more on the tool, rather than programming.
- (c) Workshop with papers on eye-tracking data provided by the organizers:
- [5]: the results of this workshop were reported in [BSS+14].
 - [24] contained a workshop on analyzing eye-tracking data from novice programmers.
- (2) Snowballing papers: 15 papers were obtained through the Snowballing process, 11 of these were from the *Psychology of Programming Interest Group (PPIG)* [CSW02, RBCL03, NS04, BT04a, NS05, BMST05b, BMST06b, BT07a, Dub09, PBT09, Loh14], and four other papers were from different sources [RCdBL02, AC06, BMST06a, DcI13].
- (3) Repeated experiments: 13 papers that re-reported or did a re-analysis of 8 published experiments were not included in the study. These repeated experiments are shown in Figure 2 where the selected paper for a repeated experiment has a star symbol to highlight it:
- [63]: A quantitative analysis of the same experiment was reported in [RLCdB02]. In this new paper, [63] analyzed the data for two of the six most vocal participants from the experiment reported in [RLCdB02], due to their differing scores.

- [12]: Re-reports the experiment from [BT04b] to further explore the effect of RFV’s blurring condition on participants with different levels of expertise.
- [6], [7], [13], and [9]: All reported results from the same experiment, but the most comprehensive report was in [BMST06a].
- [16]: Extended the analysis from [BT07b] by dividing the data into a series of shorter intervals.
- [54] and [55]: While [54] focused more on the setup of the eye-tracking environment for paired programmers, [55] was a follow-up on the results achieved and reported previously in [PBT09].
- [21]: Used the same data from the feasibility study in [BSB11], in order to study attention distribution on code elements.
- [74]: Used the same experiment data from [JN12], but with different analysis and research question.
- [62] and [61]: A more comprehensive analysis of the data from the same experiment was reported in [RLMM15].

3 MAPPING

After an analysis of the selected papers, detailed information from each publication was extracted for reporting, starting with the year of publication, the general purpose of each paper, and ending with the types of variables and eye-tracking metrics performed on each study. This section will answer our research questions and present all the information available. Figure 3 shows the summary of all papers included in our study, with basic information about the participant and the type of eye-tracker used in each study, categorized into five groups (i.e., program comprehension, debugging, non-code comprehension, collaborative programming, and requirements traceability). Over the years, research on program comprehension has received regular and more recently, increased interest among researchers. In contrast, research on collaborative programming and requirement traceability did not receive the same attention, evident by the low number of publications and experiments on these topics.

3.1 [RQ1:] How Many Papers have Used Eye-tracking in Research on Computer Programming?

Figure 4 shows the number of papers published by year. The earliest paper that reported an experiment on eye-tracking in programming comprehension was published in 1990. From the figure, we can see that nearly 62% (39 papers) of the papers on the use of eye-tracking in programming have been published after 2012, with the highest number of publications per year reaching eight papers in 2015. This reflects the increasing popularity of using eye-tracking in programming studies in recent years.

This increased interest in using eye-trackers in programming research aligns with findings reported by [42], where they found an emerging use of eye-trackers in research related to learning.

Out of the 63 papers listed in Table 3, 16 (25%) were journal papers and 47 (75%) were from conferences or proceedings. Table 4 shows the names of journal papers, while Table 5 shows conferences and workshop proceedings. On the topic of eye-tracking in computer programming, more conference papers were published than journals. The highest number of journal papers was published by “Empirical Software Engineering” with 3 journal papers. As for conference papers, “ICPC” published 9, followed by “PPIG” with 8 and “ETRA” with 4 papers on eye-tracking in relation to programming, while the remaining venues published 1 or 2 papers at most. Data analyzed for RQ1 also suggests the potential journal and conference or workshop venues for readers to submit their work related to eye-tracking and computer programming.



















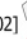








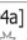











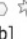
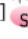










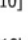


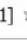



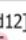













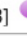





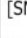
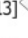


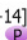

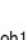

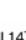


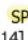


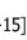



















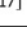





| | Program Comprehension | Debugging | Non-code Comprehension | Collaborative programming | |
|------|--|--|--|--|--|
| 1990 | [CS90]   | | | | Participants  Students  Students and Faculty  Students and professionals  Professionals  Not Available Eye-Tracker  Tobii  ISCAN  ASL  Mirametrix  RFV  SMI iViewX  FaceLAB  EMR-NC  EyeLink  Eye Tribe  GazePoint |
| 2002 | [CSW02]   | [RCdBL02]   [RLCdB02]  | | | |
| 2003 | | | | | |
| 2004 | [NS04]    | [BT04a]    [BT04b] | | [SB04]   | |
| 2005 | [NS05]   | | | | |
| 2006 | [BMST06a]  [NS06]  [UNMM06]  [AC06]  | | [Gu'06]  | | |
| 2007 | | [BT07a]   [BT07b] | [YKM07]   | | |
| 2008 | | | | | |
| 2009 | [Dub09]   | | [JGSH09]   |  [PBT09]  | |
| 2010 | [SM10a]   | |  [PG10]   [SM10b]  | | |
| 2011 | [BSB11]   | | | | |
| 2012 | [SSGA12]   | [Bed12]   [HN12]  [SFM12] | [SSvdP+12]   | [JN12]   | Requirements Traceability   ←   ←   ← |
| 2013 | [BDL+13]   [DcI13]  | [SJAP13]    [CL13]  [HLL+13] | [SMS+13]   [CTKT13]   | | |
| 2014 | [FBM+14]   [Loh14]   [BSS+14] | [TFSL14]   | [DLSL+14]   | | |
| 2015 | [ASB+15]   [RLMM15]   [BBB+15]  [MD15]   [JF15]   | | [SUGGR15]   | [MGRB15]   | |
| 2016 |  [BdP16]  [MDPVRDVI16]  |  [NHMG16]  [LWH+16]   [PLS+16] [GWMP16] | | | |
| 2017 | [PIS17]   [PSGJ17]  |  [BSL+17]  [MNH+17]  | | [DB17]   | |

Fig. 3. A visual summary of all the papers used in this mapping study.

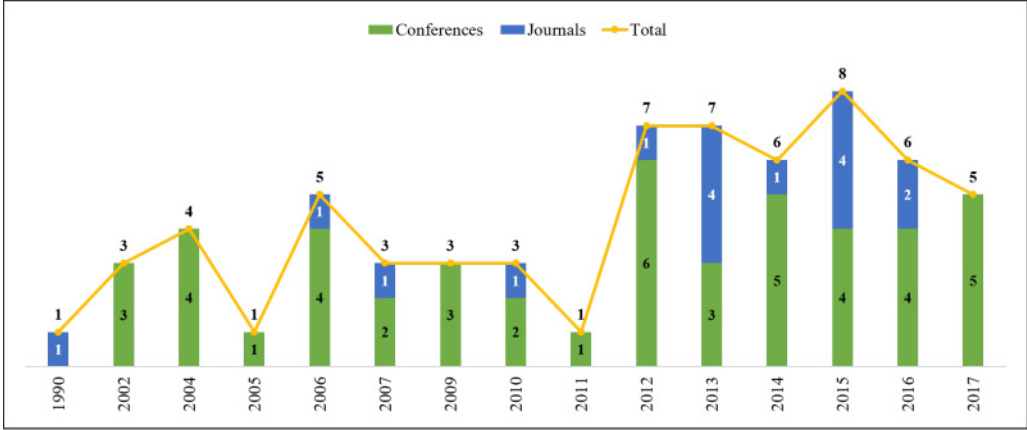


Fig. 4. Years of publication of the included papers in this mapping study.

Table 3. Selected Papers Classified into Groups Based on Tasks

| Tasks | Number of papers | List of papers |
|----------------------------|------------------|--|
| Program/Code comprehension | 26 | [CS90] [CSW02] [NS04] [NS05] [AC06] [BMST06a] [NS06] [UNMM06] [Dub09] [SM10a] [BSB11] [SSGA12] [BDL+13] [DcI13] [BSS+14] [FBM+14] [Loh14] [ASB+15] [BBB+15] [JF15] [MD15] [RLMM15] [BdP16] [MDPVRDVI16] [PIS17] [PSGJ17] |
| Debugging | 19 | [RCdBL02] [RLCdB02] [BT04a] [BT04b] [BT07a] [BT07b] [Bed12] [HN12] [SFM12] [CL13] [HLL+13] [SJAP13] [TFSL14] [GWMP16] [LWH+16] [NHMG16] [PLS+16] [BSL+17] [MNH+17] |
| Comprehension (non-code) | 10 | [Gu06] [YKM07] [JGSH09] [PG10] [SM10b] [SSVdP+12] [CTKT13] [SMS+13] [DSLS+14] [SUGGR15] |
| Collaborative | 5 | [SB04] [PBT09] [JN12] [MGRB15] [DB17] |
| Traceability | 3 | [ASGA12] [WSSK14] [ASGA15] |

3.2 [RQ2:] What Programming Tasks and Areas Were Explored Using Eye-tracking?

We looked into topics and areas of programming that have been studied by researchers using eye-trackers. The collected papers were categorized into five groups based on the type of task performed by participants. These groups are shown in Table 3. The categories shown in Table 3 were gathered through an analysis process, and they are similar to the categorization done by [67]. A pattern began to emerge, and we noticed that in relation to computer programming and using codes in eye-tracking experiments, researchers tended to perform one of these five tasks to collect eye-tracking data from participants.

Table 4. List of Journals and the Number of Papers Published

| Journal | Number of papers | List of papers |
|--|------------------|------------------------|
| Behavior Research Methods | 1 | [BT07b] |
| Communications in Computer and Information Science (CCIS) | 1 | [CL13] |
| Dyna Journal | 1 | [MGRB15] |
| Empirical Software Engineering | 3 | [PG10][BDL+13][ASGA15] |
| Journal of Systems and Software | 1 | [CTKT13] |
| IEEE Revista Iberoamericana de Tecnologías del Aprendizaje | 1 | [MDPVRDVI16] |
| IEEE Computer Journal | 1 | [CS90] |
| IEEE Transactions on Education | 1 | [LWH+16] |
| IEEE Transactions on Software Engineering | 1 | [RLMM15] |
| Interactive Learning Environments | 1 | [ASB+15] |
| International Journal of Human-Computer Interaction | 1 | [DcI13] |
| International Journal of Human-Computer Studies | 1 | [Bed12] |
| Science of Computer Programming | 1 | [DSLS+14] |
| Technology, Instruction, Cognition and Learning | 1 | [BMST06a] |

The highest number of papers focused on comprehension of code or program with 26 papers (41%). In this type of task, participants are asked to read a source code, and summarize it or answer questions related to the code for a variety of purposes such as finding reading patterns, or comparing the way an expert programmer examines a code compared to a novice.

As for debugging task with 19 papers (30%), participants were asked to find defect(s) in a source code, or perform a debugging process from a given program. The majority of the papers on debugging tasks investigated the visual strategies of programmers during debugging, and find the relation between different representations and debugging performance. As for comprehension task with 10 papers, it refers to papers examining resources other than codes, such as UML diagram or flowchart, or focus on a task related to software engineering and not specifically a source code. Most of the papers from this task used UML diagrams. In five papers, collaborative programming research focused on the visual attention of pair of programmers or more, to evaluate an emerging trend of collaborative programming. We also found three papers that used eye-tracking to evaluate links traceability in software or programs.

Figure 3 shows that the early studies to use eye-trackers in programming research mostly focused on code comprehension and debugging, while few works studied non-code representations and collaborative programming in the early years. The earliest work on collaborative programming done by [SB04] did not study the visual attention of pair programmers simultaneously, since the eye-tracking technology was not suitable for that type of research, but used the recorded eye gaze of one programmer as a cue of another programmer attempting to solve the same task. Additional setup to the eye-tracking environment was required in order to simultaneously recode the viewing habits and collaboration of pair programmers, which was done by [PBT09] who presented the details of the hardware setup for the eye-tracking environment in [54]. Although the use of eye-trackers in collaborative programming and tractability studies showed an increase in recent years,

Table 5. List of Conference and the Number of Papers Published

| Conference/Workshop | Number of papers | List of papers |
|--|------------------|--|
| CHI Conference on Human Factors in Computing Systems | 1 | [DB17] |
| Computer Software and Applications Conference (COMPSAC) | 1 | [PLS+16] |
| Computer Vision, Pattern Recognition, Image Processing and Graphics (NCVPRIPG) | 1 | [MD15] |
| Conference of the Center for Advanced Studies on Collaborative Research | 1 | [Gu06] |
| Conference of the South African Institute of Computer Scientists and Information Technologists | 1 | [BdP16] |
| Conference on Computer Supported Cooperative Work | 1 | [JN12] |
| Conference on Interaction Design and Children | 1 | [PSGJ17] |
| Diagrammatic Representation and Inference | 1 | [RCdBL02] |
| Global Engineering Education Conference | 1 | [NHMG16] |
| Hawaii International Conference on System Sciences | 1 | [AC06] |
| International Conference on Augmented Cognition | 1 | [PIS] |
| International Conference on Computer Supported Education (CSEDU) | 1 | [HLL+13] |
| International Conference on Multimodal Interfaces (ICMI) | 1 | [SB04] |
| International Conference on Program Comprehension (ICPC) | 9 | [YKM07] [SM10a] [SSVdP+12] [SSGA12] [SMS+13] [WSSK14] [BBB+15] [JF15] [MNH+17] |
| International Conference on Software Engineering (ICSE) | 2 | [FBM+14] [BSL+17] |
| International Conference on Software Maintenance (ICSM) | 2 | [SM10b] [ASGA12] |
| International Symposium on Empirical Software Engineering and Measurement (ESEM) | 2 | [JGSH09] [GWMP16] |
| International Working Conference on Source Code Analysis and Manipulation (SCAM) | 1 | [SUGGR15] |
| International Workshop on Computing Education Research (ICER) | 2 | [NS06] [BSS+14] |
| Koli Calling International Conference on Computing Education Research | 1 | [BSB11] |
| Nordic Conference on Human-computer Interaction | 1 | [BT04b] |
| Symposia on Human Centric Computing Languages and Environments | 1 | [RLCdB02] |
| Symposium on Eye Tracking Research and Applications (ETRA) | 4 | [UNMM06] [SFM12] [HN12] [TFSL14] |
| Working Conference on Software Visualization (VISOFT) | 1 | [SJAP13] |
| Workshop of Psychology of Programming Interest Group (PPIG) | 8 | [NS04] [NS05] [Loh14] [CSW02] [BT04a] [BT07a] [PBT09] [Dub09] |

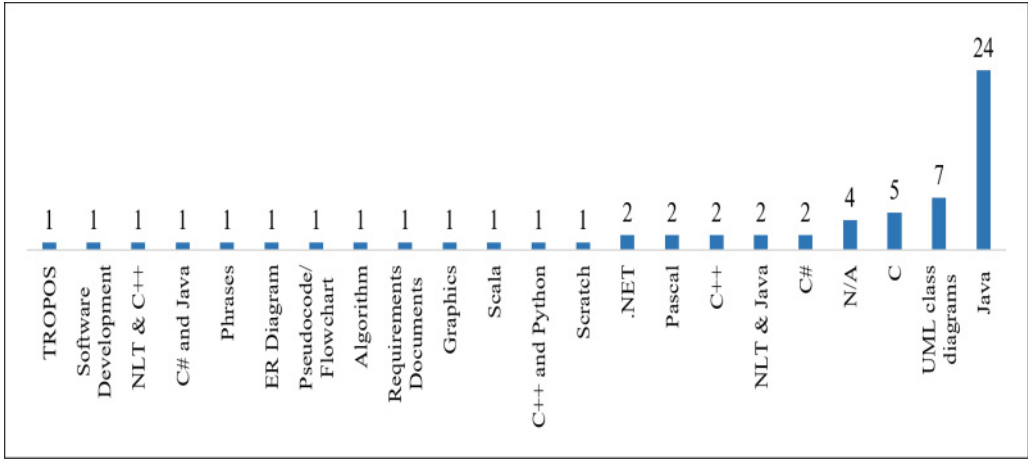


Fig. 5. Programming languages and programming representations used by papers in this study.

the number of publications in these areas is still small compared to comprehension and debugging studies. Since the early years of eye-tracking research in programming, code comprehension has been a regular interest of researchers, and saw an increase in the number of publications in recent years, as more than half the papers on code comprehension (15 out of 26) have been published in the past 5 years. The same trend of increased interest can be observed for debugging research, as nearly 70% (13 out of 19) of the published work was produced in the last 5 years.

3.3 [RQ3:] What Programming Materials were Used as Stimulus in Eye-tracking Experiments on Computer Programming?

To answer RQ3, we will look into the programming materials used by researchers. The selection of participants and materials can correlate in some cases, where both can have a great influence on the outcome of the study. Most of the papers in this study are related to code comprehension and debugging, hence the selection of the programming language and the participants' familiarity and skills in the selected language are the major factors to be considered in similar studies. While a variety of programming languages were used by researchers in this area, Java programming language stands out as the favorite choice, with students as the major participants in these studies. In answering this research question, we will list details of the type of materials and the participants considered by researchers in an eye-tracking experiment related to programming.

Materials refer to all types of stimuli used by researchers during the eye-tracking experiment. In empirical studies involving the collection of data from eye gaze and eye fixation to examine the visualization patterns or viewing habit, participants are subjected to a task related to programming. Often, stimuli for the task are prepared using a code or other types of programming data presentation, written in a programming language or presented in a format such as graphs or algorithms. The selection of materials can be in direct relation to the topic(s) being examined or selected based on popularity or participants' preferences.

While performing full analysis of the selected papers for our study, we kept track and collected information about the materials used during the experiment, and the associated programming language used. Table 9 (Appendix A) presents the type of materials used for each experiment and the details provided by researchers. It is shown in Figure 5 that the majority of the source code used during the experiment were written in Java programming language. While 24 papers

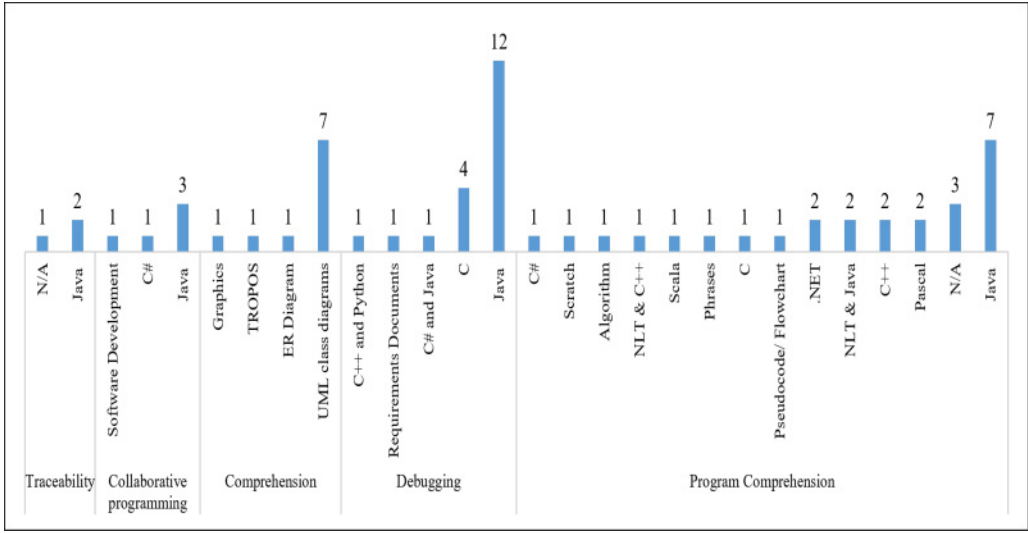


Fig. 6. Materials used by each paper grouped by tasks.

(38%) used Java alone as a main stimulus, 3 other papers used Java alongside other stimuli such as [CL13] which used Java with C# and [BSB11] and [BBB+15] which used Java with Natural Language Text (NLT). The experiment done by [BBB+15] to compare Java and NLT reading was later replicated in [PIS17], but replacing Java with C++. Another paper that used multiple materials is [TFSL14], which compared the languages of C++ and Python to assess their impact on students' comprehension. Figure 6 shows that out of the 24 papers that used Java, 9 papers (nearly 37%) were on programming comprehension, while 12 out of 19 debugging papers (63%) used Java. As for comprehension task, 7 out of the 10 papers on this task used the UML class diagram.

For a programming language in eye-tracking experiments, no language stands out as much as Java, which seems to be researchers' favorite type of language used in stimuli for eye-tracking studies related to programming. Some of the experiments listed the following reasons for choosing Java codes:

- [RLCdB02] All participants enrolled in an introductory course in Java, and the debugging environment used in the experiment was a Java software.
- [SB04] The first author was a Java programmer.
- [BT04b, BT04a] The software development environment used in the experiment was for Java debugging.
- [BMST06a] Used the Jeliot 3 visualization tool that automatically visualizes execution of Java programs.
- [AC06] Java is the primary language used to teach programming at the University of Hawaii.
- [BSB11] Because of its wide use and representativeness.
- [ASGA12, ASGA15] Java programming language was well known to the participants, and it contains a variety of different Source Code Entities.
- [HN12] Subjects had a Java programming experience of minimum 6 months.
- [JN12] Used a custom Java programming editor based on Eclipse [46] to present the code.
- [SJAP13] The software visualization tool examined, SeeIT 3D [44], is implemented as an Eclipse plugin.
- [RLMM15] The participants were professional Java programmers.

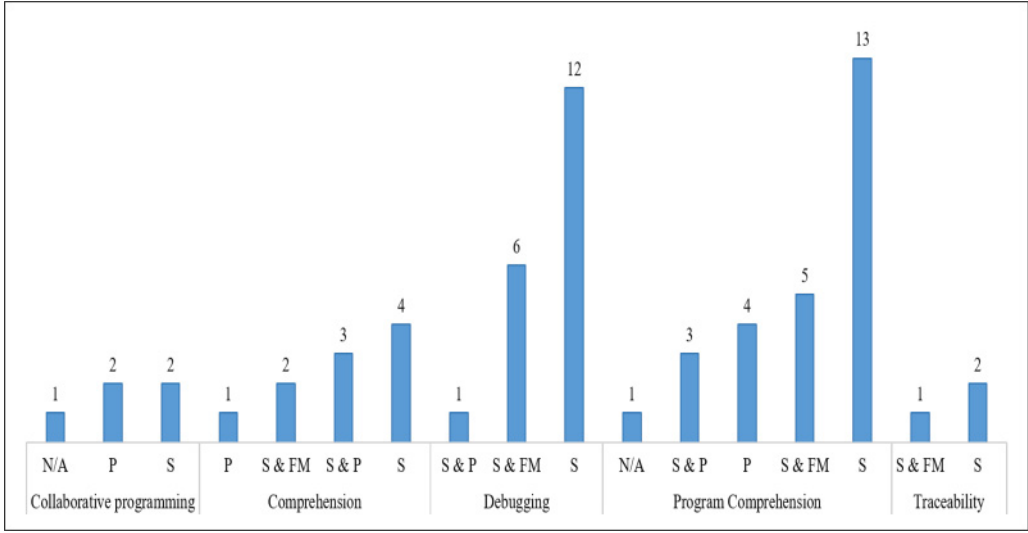


Fig. 7. Summary of the participants' affiliation grouped by task (S: Students, S&FM: Students and faculty members, P: Professionals, S&P: Students and professionals, N/A: Not available).

- [MGRB15] Used COLLECE system [19] for collaborative programming, which is compatible with and accommodates Java and C programs.
- [BBB + 15] Novice participants attended a Java beginners course, and the professionals were Java programmers.

3.4 [RQ4:] Who were the Subjects of Eye-tracking Experiments in Computer Programming?

The selection of both the materials and the subjects may be correlated in some cases. As discussed earlier in Section 3.3, some of the experiments used specific materials that required the subjects to be familiar with it. For instance, [ASGA12, ASGA15] used Java because it is well known to the participants, while [RLMM15] used Java as a material for the experiment, therefore the participants were professional Java programmers. The selection of the participants in the experiment is in some cases related directly to the aim of the study, and can have a significant impact on the outcome and the findings.

Since most of the experiments reported in relation to programming were conducted in an institute for higher learning, we see a pattern of using one or a compilation of three types of participants used as a sample in each study. As shown in Figure 7 and Table 10, each study either used Students (S), Faculty Members (FM), or professional programmers (P), or a compilation of the three to fit the purpose of the study. From Figure 7, it can be seen that more than 52% (33 papers) of the studies used students alone, while more than 22% (14 papers) used a compilation of students and faculty members (S&FM), which means that almost 75% of the studies used participants that were learning, teaching, or have a direct relation to computer programming courses. Figure 7 also provides a detailed look into the status of the participants in each task. From this figure, it can be seen that *students* were the subject of study in every task examined by researchers. We refer to the experiments with no available details about its participants with (N/A).

Figure 7 shows that seven experiments were conducted with professional programmers. For example, [SB04] examined whether the eye gaze of professional programmers in debugging a

Table 6. Details of the Sample Size for Each Task

| Tasks | Number of papers | Sum of participants | Minimum | Maximum | Average sample size | Standard Deviation |
|----------------------------|------------------|---------------------|---------|---------|---------------------|--------------------|
| Program/Code comprehension | 26 | 470 | 2 | 44 | 18.08 | 9.90 |
| Debugging | 19 | 432 | 5 | 56 | 22.74 | 13.39 |
| Comprehension (non-code) | 10 | 183 | 4 | 28 | 18.30 | 7.35 |
| Collaborative | 5 | 106 | 2 | 82 | 21.20 | 34.22 |
| Traceability | 3 | 48 | 8 | 26 | 16 | 9.17 |

program can provide hints for another programmer in finding bugs. While [DcI13] investigated reasons why the software engineering industry did not widely adopt visualization tools, with the help of professional software engineers, and [FBM+14] tried to classify the difficulty software developers experience while they work on their programming tasks, others studied the behavior of professional programmers during various programming tasks [CTKT13, BSS+14, RLMM15, DB17]. As for studies that used the compilation of students and professional programmers (S&P), three out of the six studies examined the differences between novice and experts [SSVdP+12, SUGGR15, BBB+15], while others did not state a direct reason [RCdBL02, DSLS+14, Loh14].

Sample size for each experiment can be seen in Figure 8. The largest sample size was [JN12] with 82 participants divided into pairs of programmers, followed by [PSL+17] with 56, [RLCdB02] with 49, [PSGJ17] with 44, and [TFSL14, BdP16, and LWH+16] with 38 participants in each. The average number for participant in all the studies was 19.6 with STDEV 13.5, and Table 6 shows the average sample size for each task, calculated by dividing the sum of participants by the number of papers. It also shows the minimum and maximum number of participants used, and the standard deviation of the sample size. It is important to address that the average value might be misleading due to outliers in the sample size. For instance, the average number of participants in collaborative programming studies is 21; however, Figure 8 shows that one study used 82 participants [JN12], while the others used 10 or less. A possible better alternative to average sample size might be frequency of the samples used. Nine studies used a sample of less than 10 participants [CL13, WSSK14, RCdBL02, UNMM06, CTKT13, PBT09, BSS+14], while the most common sample size was 15 participants, which was used in seven studies [AC06, SMI10b, SM10a, BSB11, SFM12, BDL+13, FBM+14, MD15, and DB17]. Other noticeable frequencies in the sample size were 20 participants which mostly occurred in recent studies [SJAP13, BBB+15, JF15, SUGGR15, PLS+16, MNH+17], 24 in four studies [JGSH09, PG10, SSGA12, NHMG16], and 38 participants in three studies [TFSL14, LWH+16, BdP16].

Table 10 (Appendix A) lists the details of the sample of participants in each paper. Table 10 (Appendix A) shows the status of the participants, the sample size considered in the experiment, the way the sample was split or divided in the experiment, the details on the split sample size, and the gender of the participant if mentioned.

In some studies, the sample considered for the experiment was not necessarily the same selection of sample considered for the eye-tracking data analysis for various reasons. Some samples were discarded due to technical issues with the eye-tracking device [BMST06a], data corruption [BT07a, SFM12, DSLS+14], participants withdrawing from the experiment [BT07b, FBM+14], or they were unfit for the experiment [ASGA12, ASB+15].

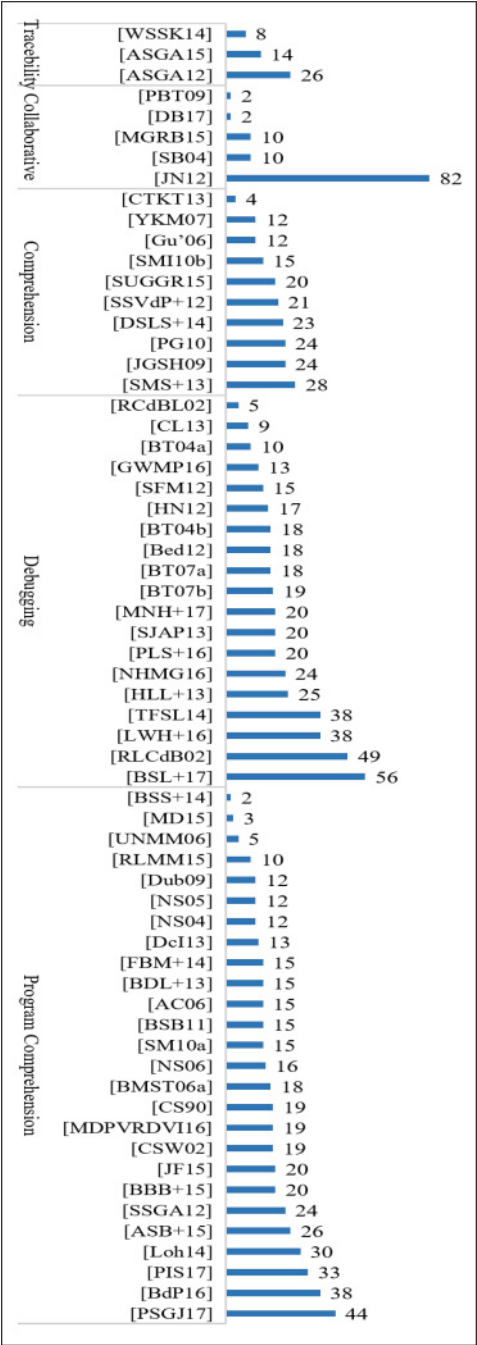


Fig. 8. Details of the participants in all papers grouped by task.

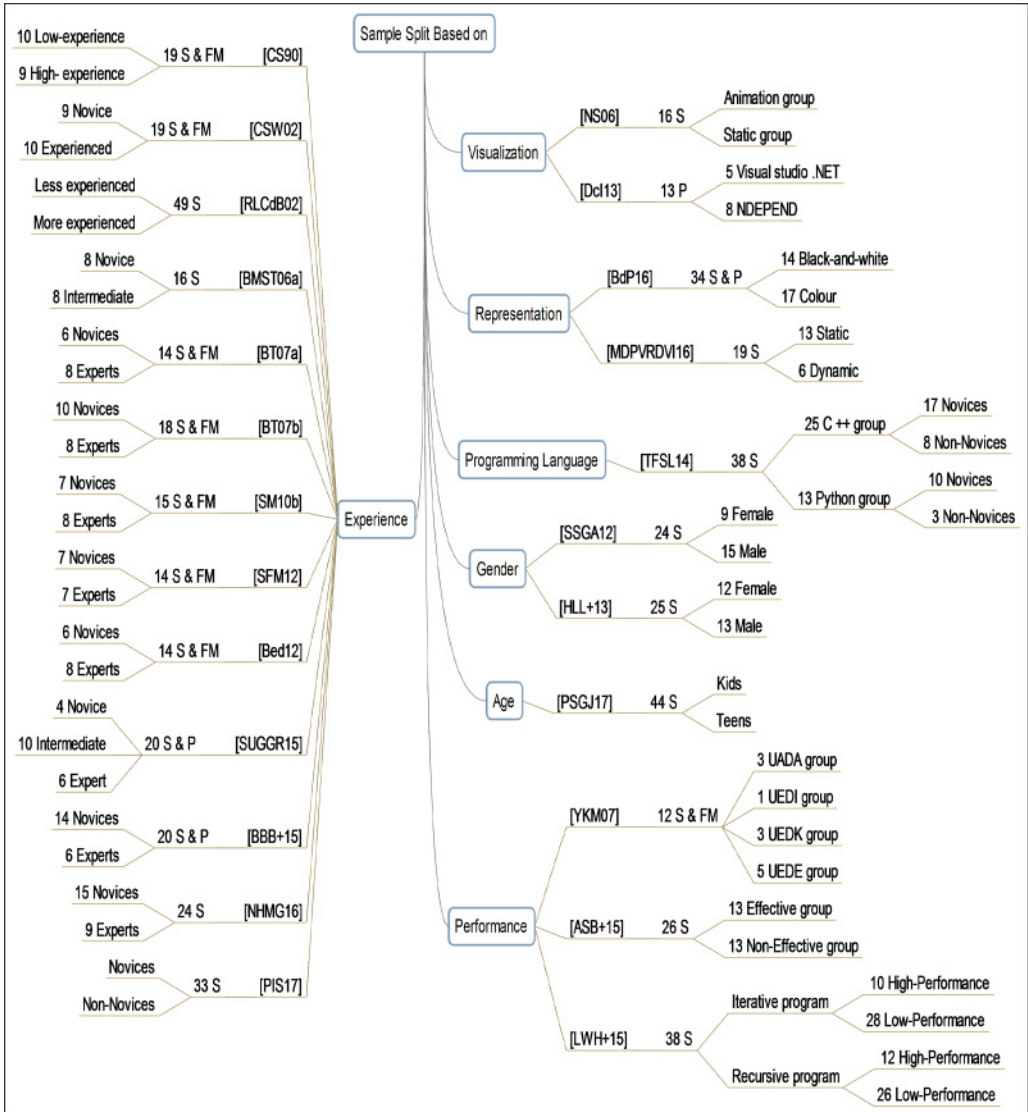


Fig. 9. Sample analysis is done based on a single split.

In the last column of Table 10, the gender of the participants is listed if it was mentioned by the researcher. Although in most cases the gender of the participants did not affect the analysis of the results, some of the work tried to provide as much information on the selected subjects as possible. Some work lists the experience of the participants, classes they have taken, and the department from which they were recruited. All the reported information regarding the participants can help provide some insights into the setup of the experiments, but there were few studies in which details of the participants were important to the study and the analysis of the results, especially in the cases where the sample was divided into two or more groups. Although columns 4 (Split) and 5 (Split details) in Table 10 list the way in which the sample was divided, the details of the sample division and split can be seen in detail in Figure 9.

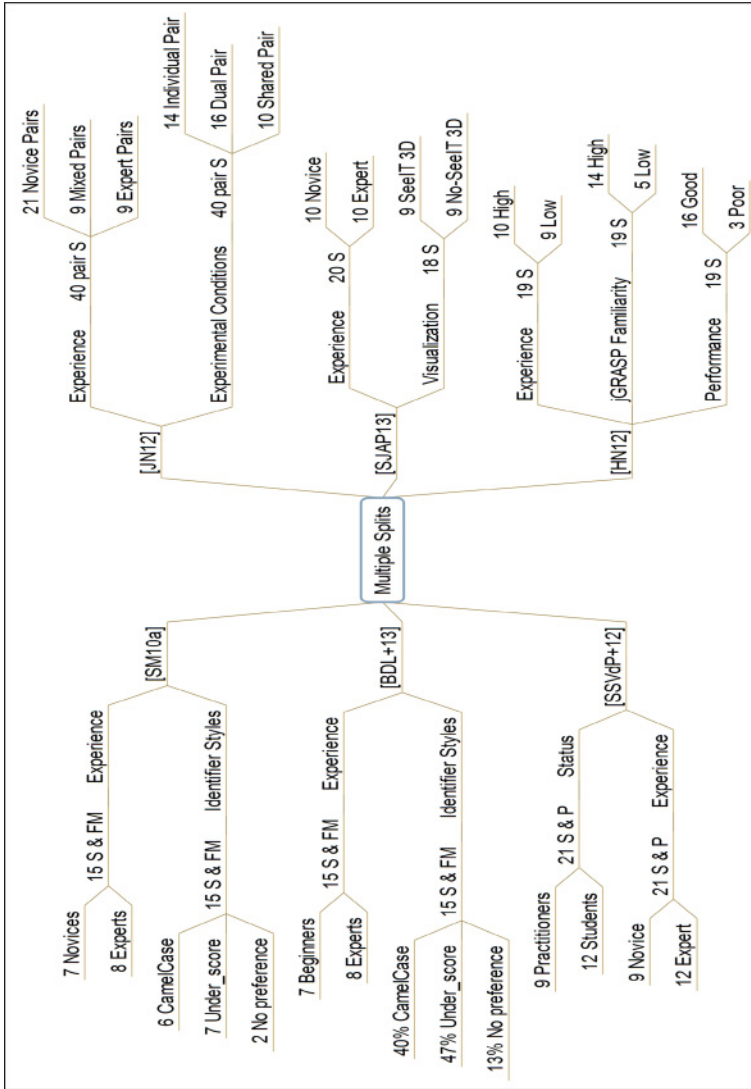


Fig. 10. Samples were split into multiple groups.

Sample Split. We use the term sample split to refer to the way a sample of participants or the collected data from participants were divided for purposes related to the experimental setup. Since the selected studies focused on computer programming, the participants had various knowledge and skills in one or more programming languages and programming representations.

While most of the studies analyzed the collected data from the sample of participants as a whole, 30 studies (48%) divided the selected sample into one or more groups for various reasons. Figure 9 shows the details of the papers that had a single sample split from the studies, while Figure 10 shows the studies that divide the participants based on two or more splits. Out of the 30 reported experiments, 24 did a single split of the sample, while the remaining 6 studies did two or more splits.

From the 24 single split experiments, 13 divided the sample based on experience, while the remaining 11 divided the sample based on other reasons such as visualization, representation,

identifier style, gender, age, programming language, or performance. Details of the single sample split of the participants in Figure 9 are as follows:

(1) Samples divided based on experience:

- [CS90] In order to examine the differences in reading algorithms, 10 students from the second semester of a computer science course were considered as the low-experience group, while eight graduates and one PhD faculty member were the high-experience group.
- [CSW02] To determine how programmers from different experience levels understand a typical simple program, nine students with one semester of programming experience were considered novices, while the experienced group consisted of computer science faculty members and students from advanced undergraduate and graduate classes.
- [RLCdB02] To find out if excellent debugging skills were associated with specific representation patterns, the sample was divided into two groups of less and more experienced students.
- [BMST06a] To improve program visualization systems, it is important to study how the behavior of the novices differs from intermediates in relation to program animation. Subjects who had less than 2 years of experience in programming were considered to be novices, while those with 2 years or more were intermediates.
- [BT07a] To investigate the visual strategies of programmers during debugging, six participants were in the novice group, while eight were in the experienced group, based on the reported period of experience by the participants.
- [BT07b] To investigate the differences in the allocation of visual attention between expert and novice programmers, subjects with an average of 8.13 months of Java experience were the novice group (10 programmers), while the remaining eight subjects who had an average of 16.25 months of Java experience formed the expert group.
- [SM10b] To examine the impact of design patterns in different layouts on the comprehension of UML class diagrams. Seven second-year students formed the novice group, while six graduates along with two faculty members were grouped as experts on UML diagrams and design patterns.
- [SFM12] To study the impact of scan time on debugging performance, the novice group had seven students in their second year of undergraduate study, while the expert group was formed from six graduates and two faculty members.
- [Bed12] In order to study the visual attention strategies of programmers during debugging, students and faculty members were divided into two distinct levels of experience consisting of eight experts and six novices.
- [SUGGR15] In order to find out what elements of the visualizations are most important for comprehension, participants were divided based on their performance, into 4 novices, 10 intermediates, and 6 experts.
- [BBB + 15] In order to identify the reading patterns of source-codes and its linearity (top to bottom, left to right), eye-tracking data was collected from 14 novice students from a beginners course in Java, and 6 professional software engineers.
- [NHMG16] to examine the difference in debugging behavior based on experience, the sample of 24 students was divided into 15 novice and 9 considered to be advanced.
- [PIS17] In order to examine the differences in reading patterns between NLT and C++, the sample consisted of 33 novice and non-novice students.

- (2) Samples divided based on gender:
 - [SSGA12] In order to study the relationship between source-code reading and identifier style in regard to gender, 9 female and 15 male students were the subject of the experiments.
 - [HLL + 13] To address and investigate the gender differences in program debugging, 12 female and 13 male students from the Department of Computer Science were recruited.
- (3) Samples divided based on code representation:
 - [BdP16] In order to investigate the effect syntax highlighting may have on the reading behavior of programmers, a sample of 34 students was divided in half into a black-and-white group and a color group, based on the syntax highlights. Then the data was compared with 4 experts.
 - [MDPVRDVI16] To assess the usefulness of the GreedEx representations, two experiments were performed, where the first one had 13 students work with a static representation, while the second experiment had 6 students work with a dynamic representation.
- (4) Samples divided based on performance:
 - [YKM07] In order to evaluate the comprehension of participants with varying knowledge of UML diagrams and their designs, the sample was divided based on performance into four groups: • UADA: UML and Design Agnostic had 3 subjects; • UEDI: UML Expert but Design Inexperienced had one subject; • UEDK: UML Expert and Design Knowledgeable had 3 subjects; and • UEDE: UML and Design Expert had 5 subjects.
 - [ASB + 15] To track and study the algorithmic problem solving techniques among students, 13 students who solved the algorithmic problem were selected to the effective group, and 13 who solved the problem incorrectly were selected to the non-effective group.
 - [LWH + 16] To investigate the behavior of students in debugging, each participant's performance was evaluated based on tasks related to debugging two programs provided in either an iterative structure (10 high-performance, 28 low-performance) or a recursive structure (12 high-performance, 26 low-performance).
- (5) Samples divided based on visualization:
 - [NS06] In order to study and compare the differences between animated and static visualization during program comprehension, the participants were grouped into either an animation or a static group, based on the score of a pre-test.
 - [Dcl13] In order to evaluate different software visualization tools, 5 professional programmers were assigned to the visual studio group, while 8 were in the NDEPEND group.
- (6) Sample divided based on programming language:
 - [TFSL14] To examine if program comprehension of students is influenced by the programming language selected, 25 subjects were assigned to the C++ group (17 novices, 8 non-novices) or Python group (10 novices, 3 non-novices) based on a background questionnaire.
- (7) Sample divided based on age:
 - [PSGJ17] To understand the differences in coding activities between kids (ages 8–12) and teens (ages 13–17), the eye movements of 44 students were recorded while working on Scratch tool [60].

Six experiments reported in this survey divided the selected sample of participants into two or more groups. Details of the multiple sample split of the participants in Figure 10 are as follows:

- (1) [SM10a] In order to examine how students from different backgrounds perform with different identifier styles, the collected data was analyzed based on the following:
 - Experience: 7 novices and 8 experts.
 - Identifier style: 6 participants stated they prefer camel-case style, 9 prefer underscore, and 2 had no preference.
- (2) [BDL + 13] Similar to [SM10a], participants were divided into groups based on the following:
 - Experience: 7 beginners and 8 experts.
 - Identifier style: 40% of the participants stated they prefer camel-case style, 47% prefer underscore, and 13% had no preference.
- (3) [SSVdP + 12] In order to evaluate the impact expertise and experience has on the time and accuracy of UML class diagrams, students and professional programmers were evaluated based on the following:
 - Status: 9 practitioners and 12 students.
 - Experience: 9 novices and 12 experts.
- (4) [JN12] 82 students were recruited for a collaborative programming experiment, and divided into pair programmers based on the following:
 - Experience: 40 pair of 21 novices, 9 experts, and 9 mixed expertise.
 - Experimental conditions: 40 pair where 14 were individual pair, 16 dual pair and 10 shared pair.
- (5) [SJAP13] In order to assess the effect of SeeIT 3D, students were divided based on experience and visualization conditions as follows:
 - Experience: 10 novice and 10 experts.
 - Visualization: 9 with SeeIT 3D, and 9 with No-SeeIT 3D tools.
- (6) [HN12] 19 students were evaluated and given a score based on their experience, performance, and familiarity with the Integrated Development Environment (IDE) used:
 - Experience: 10 high and 9 low.
 - jGRASP familiarity: 14 high and 5 low.
 - Performance: 16 good and 3 poor.

3.5 [RQ5:] What Eye-tracking Technology and Devices were Used in Computer Programming Experiments?

Before we discuss the devices used by researchers to track participants' attention while working on a task related to programming, we will provide a brief overview of the eye-tracking technology. For more details on the eye-tracking devices used and their measures and abilities, we refer the reader to a recent review that examined the use of eye-tracking in learning [42], and [67] for details on the metrics used by researchers for data collected from experiments related to software engineering. The use of an eye-tracker in an experiment involves the selection of the devices, as well as the types of data eye-trackers can collect to calculate the variables that help the researcher evaluate their hypothesis for the experiment. In this part, we will discuss the different types of eye-trackers used in experiments related to programming, then in the next part we will examine the common eye-tracking metrics these devices used to collect data before looking into the variables evaluated based on these eye-tracking metrics.

Eye-tracking Devices. From Figure 11 and Table 11 (Appendix A), the details of the devices used are listed for each experiment. In relation to programming, the literature we reviewed used one of two methods to track the attention of programmers: Restricted Focus Viewer (RFV) or eye-tracker.

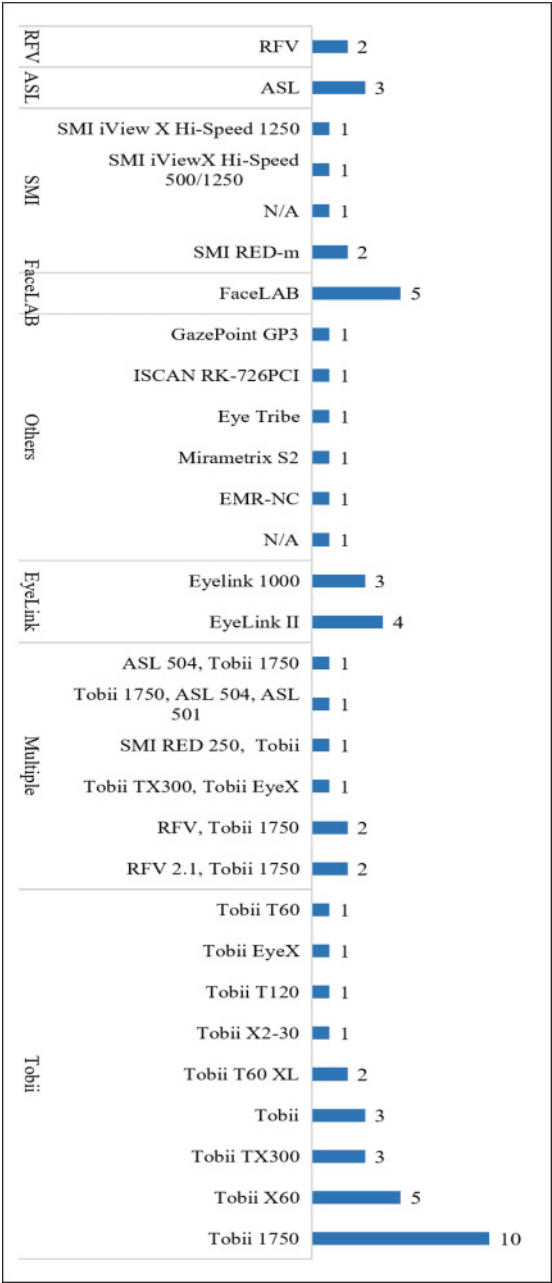


Fig. 11. The version of eye-tracker used, grouped by manufacturer.

The RFV was developed as an alternative to eye-trackers [18]. RFV allows subjects to see a limited focused region of a material at a time, and the subject uses a mouse to move the focused area, while the RFV tracks and records these moves of the mouse and the timestamps for analysis [10, 11, 18]. Figure 12 shows that the RFV alone was used in two studies, while it was used in four studies

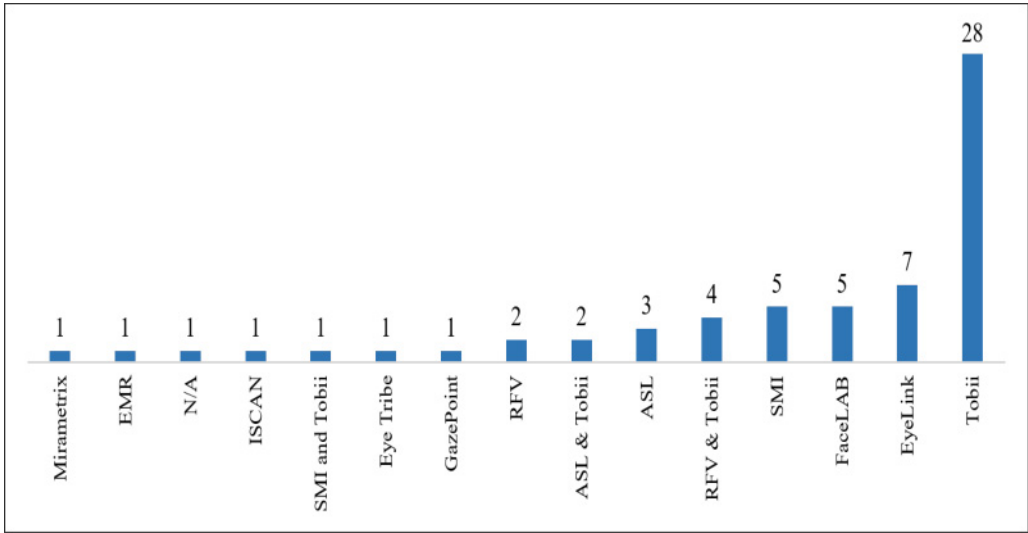


Fig. 12. The version of eye-tracker used grouped by manufacturer.

along with a Tobii eye-tracker. In [RCdBL02, RLCdB02], RFV was used to track the visual attention of programmers during a debugging task, focusing mainly on switching behavior between different representations of the code. Later, these experiments were replicated by [BT04a, BT04b], with the addition of a Tobii eye-tracker, in order to compare the performance and accuracy of the two technologies in measuring visual attention, and verify the effect of RFV on the behavior of participants. The validity of the RFV was later questioned again, and examined in [BT07a, BT07b] but this time with an RFV 2.1 with a Tobii eye-tracker. Figure 11 and Table 11 show that RFV was not widely adapted and used in experiments related to computer programming, and its validity was questioned and examined multiple times.

The reported studies indicated two types of trackers: intrusive tracker and non-intrusive tracker. Intrusive devices are head-mounted onto the subjects' head with a headband, and require the subject to be situated in front of the screen to keep a relatively steady position. Examples of these intrusive devices are the *EyeLink II* tracker from SR Research (<http://www.eyelinkinfo.com/eyelinkII.html>), which was used by [Gu06, JGSH09, PG10, SSVdP+12], the *ISCAN RK-726PCI* (<http://www.iscaninc.com>) pupil/corneal reflection eye-tracker used by [SB04], and the head-mounted *ASL 501* from Applied Science Laboratories (<http://host.web-print-design.com/asl/>), which was compared to two non-intrusive devices by [NS04]. Few researchers stated a clear reason for not using this type of device, such as [PBT09] who in their attempt to build a system for tracking the visual attention of a pair of programmers stated:

“As our goal was to combine both programmers eye-tracking data with the single display screen, we found that the field of view of the head mounted camera was too large, and coordination of its output with the screen was too coarse.”

Also [JF15] stated:

“There is an obvious advantage to using a remote eye-tracker over a head mounted device, especially when considering intrusiveness and how natural is the experiment environment.”

As for non-intrusive devices, these devices were the most frequently used. It allows the visual attention to be measured without interfering with the subjects thought process, and without restricting their head movement, except during the calibration process. Examples of non-intrusive eye-trackers are the FaceLAB from Seeing Machine (<https://www.seeingmachines.com>) used by [ASGA12, SSGA12, DSLS+14, ASGA15, SMS+13] and Tobii 1750.

We found a variety of devices used in the experiments, with one manufacturer standing out with the most popular device for tracking the visual attention of participants in relation to programming, which is Tobii. The first eye-tracker used in a study in relation to computer programming was the ASL eye-tracker used by [CS90] to examine the way programmers read algorithms. It was used again by the same first author in [CSW02], and by [AC06]. Two different distributions of ASL (ASL 501, ASL 504) were evaluated and compared with Tobii by [NS04], especially focusing on the head-mounted version (ASL 501). ASL was also part of a system developed and built by [PBT09] for tracking the eye movement of pair programmers during collaborative programming. The least used devices in the experiments are (1) ISCAN: A lightweight head-mounted eye-tracker used by [SB04] to demonstrate that debugging performance can improve by viewing another person's eye gaze; (2) NAC Incorporated Non-Contact Eye Mark Recorder (EMR-NC) (<http://www.nacinc.com>): Used by [UNMM06] to characterize the performance of individuals in reviewing source code; (3) Mirametrix S2 eye-tracker (<http://www.mirametrix.com>): Used by [TFSL14] to compare C++ and Python; (4) Eye Tribe (<http://www.theeyetribe.com>): Used by [JF15] to verify and quantify the effect of repetitions of code pattern comprehension; and (5) GazePoint GP3 (<https://www.gazept.com>): Used by [BSL+17] to examine if developers pay attention to error messages from the IDE compiler. On the other hand, the most widely used device by researchers in programming research is Tobii. Tobii was reported in more than 55% of the gathered papers (35 papers). In 7 of these, Tobii was used along with another device such as ASL, RFV, and SMI, while it was the main device for eye-tracking in 28 papers (44%). The versions of Tobii eye-trackers used in the 35 experiments are shown in Figure 13, where the papers are organized by the year of publication. Figure 13 shows that version 1750 was the type of eye-tracker most widely used in a total of 16 studies (25% of all the papers reported in this study) between the years 2004 and 2013. However, Tobii 1750 is now considered an outdated model prior to the introduction of T60 and X60 which have been replaced by the newer and mobile X2 version. Until recent years when it was replaced by newer systems, Tobii 1750 was leading the way in tracking the eye movement of participants in a programming task.

3.6 [RQ6:] What Eye-tracking Metrics and Variables are Commonly Reported in Programming Studies?

The main reason for using an eye-tracking device in an experiment is to be able to gather information on the visual attention of subjects while they examine a stimuli or material on a computer screen. The visual attention gives some indication of the underlying cognitive processes occurring at the time a participant performs a given task. In this part, we will map the common eye-tracking metrics and variables used by researchers in programming studies.

3.6.1 Eye-tracking Metrics. According to [42], eye movements are generally made from a series of *fixations* and *saccades*. While *fixations* are a state of relatively stable eye movement ranging from 100 to 500ms, *saccades* occur between two consecutive fixations in the form of rapid eye movement which in a typical reading task usually last from 30 to about 50ms [57, 58].

A range of eye-tracking metrics have been utilized and used in the studies, and we will focus on reporting the most common eye-tracking metrics in Table 7. Most of the experiments on computer

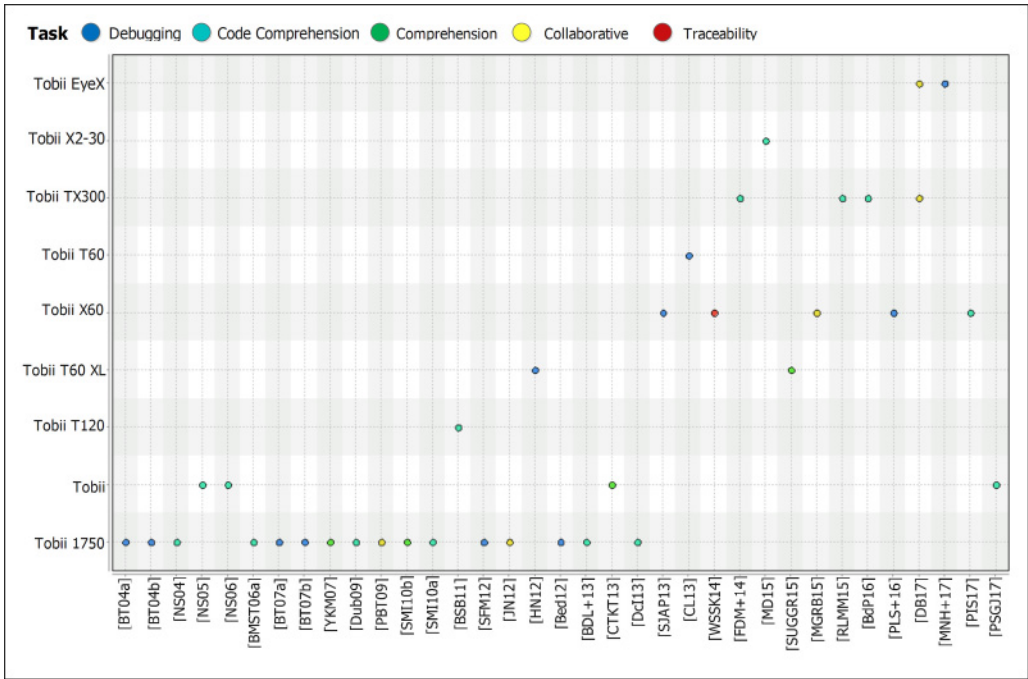


Fig. 13. Version of Tobii devices used in each experiment, organized by year of publication.

Table 7. List of Most Common Eye-tracking Metrics Used in the Papers

| Based on | Measurements | Number of Papers | Papers |
|-----------------------|---------------------------|------------------|--|
| Number of Fixations | Fixation Count | 18 | [BT04b] [BMST06a] [BT07b] [UNMM06] [SM10a] [SM10b] [HN12] [SFM12] [SSGA12] [BDL+13] [SJAP13] [TFSL14] [ASB+15] [JF15] [BdP16] [GWMP16] [MDPVRDVI16] [MNH+17] |
| | Fixation Rate | 5 | [SM10a] [SM10b] [SSGA12] [BDL+13] [TFSL14] |
| | Spatial Density | 4 | [Gu06] [BT07a] [SSVdP+12] [BSS+14] |
| | Convex Hull Area | 4 | [SSGA12] [SSVdP+12] [SMS+13] [DSLS+14] |
| Duration of Fixations | Fixation Time | 18 | [CS90] [CSW02] [RLCdB02] [BT04b] [BMST06a] [BT07b] [PG10] [BSB11] [BDL+13] [JF15] [ASGA15] [SUGGR15] [BdP16] [GWMP16] [NHMG16] [PLS+16] [MNH+17] [PSGJ17] |
| | Average Fixation Duration | 11 | [CS90] [CSW02] [SM10a] [SM10b] [PG10] [SSVdP+12] [BDL+13] [SMS+13] [ASGA15] [GWMP16] [MDPVRDVI16] |
| Attention Switching | Attention Switching | 12 | [RLCdB02] [RCdBL02] [BT04b] [BMST06a] [BT07a] [BT07b] [BSB11] [Bed12] [HN12] [PLS+16] [MNH+17] [PSGJ17] |
| Scan-paths | Scan-path | 11 | [SM10b] [PG10] [PBT09] [DcI13] [SJAP13] [BSS+14] [MGRB15] [ASB+15] [SUGGR15] [BdP16] [MDPVRDVI16] |

programming used eye-trackers to collect the number and duration of the subjects' fixations. We categorized the eye-tracking metrics used in the studies into measures calculated using the number of fixations, measures calculated using fixation duration; attention measures and scan behavior measures. This information, along with other metrics that eye-tracking devices are able to collect, can be used in many formulas that will help researchers evaluate the participants level of cognitive processing for the research questions studied. Details of the eye-tracking metrics listed in Table 7 are as follows:

(1) **Measures related to the number of fixations:**

- *Fixation Count*: Defined as the total number of eye fixations on a specific area, and can be interpreted as an indication to the importance of an area [3, 8, 36]. A higher fixation count on a specific area indicates more interest in that part, and suggests more effort spent in solving the task [73].
- *Fixation Rate*: Calculated by dividing the number of fixations on an Area of Interest (AoI) by the total number of fixations, and it can be an indicator of the importance of a specific part [31, 69, 72, 73]. The AoI on which the fixation rate was calculated by researchers is related to the task or the objective of the research. For example, the AoI for [TFSL14] was buggy lines of code.
- *Spatial Density*: Proposed by [31], it is calculated by dividing the screen into a grid, then finding the proportion of cells with at least one fixation to the total number of cells in the grid [23, 28]. More detailed information about this measure was presented by [DSLS+14] as they presented the Taupe visualization tool.
- *Convex Hull Area*: This was introduced and used by [31] to evaluate user interfaces quality, and it represents the smallest convex [31, 66, 78] or polygon [28] set of fixations which includes all fixations. Closer fixations are indicated by a smaller value, which in turn indicates that the participants spend less effort to find usable elements [SSVdP+12] or areas [SMS+13].

(2) **Measures related to the fixations duration:**

- *Fixation Time*: Total fixation time is calculated as the sum of all fixation durations, and it measures the time spent on an AoI or the stimuli by the participants [1, 13]. It is calculated and presented either as the total time in seconds of all fixations or as a percentage of the total time [26]. This measure does not only help calculate the average fixation duration, but it can also determine the visual effort [69, 72, 73]. More visual effort from the participants to solve the task is indicated by a higher fixation count, duration, and fixation rate [73].
- *Average Fixation Duration (AFD)*: The AFD refers to the portion of fixation time to total time spent on a particular AoI or the stimuli [31], and it can be calculated for an AoI or the stimuli using the fixation time and number of fixations. Longer AFD means that more time is needed by the participants to analyze and build a mental model of the task [31, 66].

(3) **Attention switching**: Refers to the number of switches between two or more AoIs [12]. It was mainly used to find the switching behavior of participants with different representations [BMST06a, BT07a, Bed12], or the effect of RFV on programmers' behavior [BT04b, BT07b]. Attention switching can be interpreted as a measure of an existing relation between two or more AoIs, that is, the more frequent the switches between AoIs, the more related they are.

(4) **Scan-paths**: A path formed by *saccades* between *fixations*, in the form of a directed sequence [23, 64, 70, 73]. It provides the order and directionality in which elements of the

material were examined by the subject. For example, what parts were immediately visited and which area directed the subject to where they are at a specific point. This type of information can be used to try and find reading patterns [56], and help organizing visual space and designing visual layouts [71].

These are the most commonly used eye-tracking metrics to calculate the data collected from the eye-tracker. However, there are some variations and more complex eye-tracking metrics that can be used, depending on the type of information the researcher is aiming to find. Along with these metrics, eye-trackers can also provide some techniques to visualize and view details of the participants' behavior during the task performance. For example, the heat-map technique is used to visualize the intensity of fixations using a color spectrum. It overlays on top of the stimulus to show the viewing patterns of the subjects [22, 36, 73]. Jbara and Feitelson [36] provides detailed information on the use of heat maps, along with examples and the participants responses to their heat maps. Also, some examples of heat maps can be found in [BSB11, DSLs+14, ASB+15, JF15, SUGGR15, MNH+17]. Other eye-tracking data can be presented in the form of a gaze plot, which can help in visualizing the scan-path by displaying the eye gaze data for each stimulus in a static viewpoint [73]. Some of the experiments used this technique, and examples can be seen in [YKM07, SM10a, SM10b, BDL+13, MGRB15, SUGGR15].

3.6.2 Variables. Based on the collected data from an eye-tracking experiment, along with the metrics these devices calculate, researchers can evaluate the variables corresponding to their hypothesis and thus make sense of the data and derive a conclusion for their work. Table 12 (Appendix A) lists all the available variables from the literature.

Dependent Variables. Dependent variables are selected based on the hypotheses for the experiment [56, 80], and are directly related to the analysis and interpretation of the collected data. Table 8 shows the most common dependent variables from the study papers, which are as follows:

- (1) *Time*: The most commonly measured and used variable in studies of eye-tracking in computer programming. While the term *Time* in Table 8 is general and refers to all the studies that used a variation of time measurements, Table 12 (Appendix A) shows more details on what time was measured and used to evaluate the findings. The most common types of time measured in the study papers were the following:
 - *Completion Time, Required Time, Speed or Efficiency*: The time needed by a participant to perform and complete a task [SSGA12, SSVdP+12], or the time spent on performing each task [JF15, SM10b, SMS+13, TFSL14, GWMP16, NHMG16, DB17].
 - *Fixation Time*: The time participants spent focusing on the AoI [CS90, BMST06a], or a type of Source Code Entities (SCEs) [ASGA12, ASGA15].
 - *Scan Time*: Time needed for an initial scan pattern to occur [UNMM06, SFM12].
 - *Detection Time*: The time taken to detect or find defects and bugs, and it is related to studies on scan time [UNMM06, SFM12, NHMG16, MNH+17].
 - *Response Time*: Time taken by the subjects to start answering the task [YKM07, SUGGR15].
 - *Inspection Time*: Time spent looking on a certain AoIs [MGRB15, MDPVRDVI16].
- (2) *Performance*: It can also be referred to as accuracy, which is based on the subjects' performance, and is measured by the correct answers given by subjects during a task of debugging or comprehension. Performance and accuracy was measured in different ways such as the *Percentage of Correct Answer (PCA)* provided by a participant [SSVdP+12, SSGA12, SMS+13, ASGA15], all scores summed from each task [SM10b], or a scale from 1 to 5 (wrong to correct) to measure accuracy [TFSL14]. While in most of the studies

Table 8. The Most Common Dependent Variables in the Study Papers

| Dependent variable | Task | Number of Papers | Papers |
|--------------------|---------------------------|------------------|---|
| Time | Program Comprehension | 17 | [CS90] [CSW02] [NS04] [Dub09] [SM10a] [BSB11] [SSGA12] [BDL+13] [Dcl13] [Loh14] [FBM+14] [BSS+14] [ASB+15] [BBB+15] [JF15] [BdP16] [PSGJ17] |
| | Debugging | 13 | [RLCdB02] [BT04a] [BT04b] [BT07a] [BT07b] [SFM12] [HN12] [Bed12] [TFSL14] [GWMP16] [NHMG16] [PLS+16] [MNH+17] |
| | Non-code Comprehension | 10 | [UNMM06] [YKM07] [JGSH09] [PG10] [SM10b] [SSVdP+12] [SMS+13] [SJAP13] [DSLS+14] [SUGGR15] |
| | Collaborative programming | 2 | [MGRB15] [DB17] |
| | Traceability | 1 | [ASGA12] |
| Performance | Program Comprehension | 11 | [CSW02] [NS04] [NS05] [NS06] [AC06] [Dub09] [SM10a] [SSGA12] [Dcl13] [BDL+13] [Loh14] |
| | Debugging | 10 | [RCdBL02] [BT04b] [BT07b] [Bed12] [SFM12] [CL13] [TFSL14] [JF15] [GWMP16] [NHMG16] [PLS+16] [MNH+17] |
| | Non-code Comprehension | 8 | [YKM07] [SM10b] [SSVdP+12] [SMS+13] [CTKT13] [SJAP13] [DSLS+14] [SUGGR15] |
| | Collaborative programming | 2 | [SB04] [PBT09] |
| | Traceability | 2 | [WSSK14] [ASGA12] |
| Visual Effort | Program Comprehension | 5 | [SM10a] [SSGA12] [BDL+13] [FBM+14] [JF15] |
| | Debugging | 2 | [SFM12] [TFSL14] |
| | Non-code Comprehension | 6 | [YKM07] [JGSH09] [SM10b] [SSVdP+12] [SMS+13] [SJAP13] |
| | Collaborative Programming | 0 | N/A |
| | Traceability | 0 | N/A |
| Viewing Pattern | Program Comprehension | 3 | [BMST06a] [BdP16] [PSGJ17] |
| | Debugging | 11 | [RLCdB02] [BT04a] [BT04b] [BT07a] [BT07b] [HN12] [Bed12] [HLL+13] [GWMP16] [PLS+16] [MNH+17] |
| | Non-code Comprehension | 0 | N/A |
| | Collaborative Programming | 1 | [DB17] |
| | Traceability | 0 | N/A |

accuracy was related to the participants' performance, there was one study where accuracy was related to the eye-tracking devices and their performance during a computer programming experiment [NS04].

- (3) *Viewing Pattern*: It includes patterns such gaze path, viewing path, number of switches between different areas and switching frequency which can measure the dynamics attention allocation [8]. It is a common variable for studies that focused on visual attention of multiple representations and visualization tools, and it is measured by the tracking device.

Viewing pattern try to identify the most associated areas of a code, its representations or AoIs [RLCdB02, BMST06a]. Viewing patterns and switching behaviour were notably

used in the studies that focused on validating the RFV, and determining its effect on the participant's behaviour and attention [BT07a, BT07b].

- (4) *Visual Effort*: Indirectly measured with fixation related variables such as count and duration [36, 69, 72, 73]. It is the amount of visual attention the subjects spend on an AoI [66, 70, 80], or the visual effort spent to find the answer [73].

Independent Variable. The independent variables mostly were related to the aim of the study and they are related to the hypotheses, and may have impact on or/and affect the dependent variables [1, 56]. While full details of the independent variables can be seen in Table 12, the most common independent variables we found in the study papers were related to the following:

- Programming experience in tasks that tried to study the differences between participants with different expertise [CS90, CSW02, RLCdB02, BMST06a, BT07b, HN12, SSVdP+12, Bed12, MGRB15, BBB+15].
- Type of representation in experiments that studied the effect of representation on participants' performance [RLCdB02, PG10, SMS+13, ASB+15].
- The different types of source code entities and code elements [ASGA12, ASGA15].
- Visualization tools and their effect on participants' performance [NS05, NS06, SJAP13, DcI13].
- Identifier style [SM10a, SSGA12, BDL+13].
- Gender of the participants in studies that examined if gender has an effect on computer programming [SSGA12, HLL+13].
- Errors or defects, their type and presence mostly in debugging studies [RLCdB02, UNMM06, SFM12].
- Algorithm or code type and language in studies that used multiple codes or languages [AC06, HLL+13, TFSL14, LWH+16, JF15].
- The RFV restricting condition in studies that validated RFV [BT04b, BT07b].

Mitigating Variables. Mitigating variables are related to both independent and dependent variables, and they can have an impact on the way independent variables affect dependent variables [1, 66, 68]. The details of the mitigating variables are shown in Table 12. In the survey papers we studied, we found mitigating variables to be mostly related to the experience of the participants and their knowledge of a programming language, representation of source code or diagrams, or a tool. Examples include the participants' experience [AC06, SM10a, ASGA12, SFM12, BDL+13, SMS+13, SJAP13, TFSL14], level of study [ASGA12, SSGA12, SMS+13], UML knowledge [JGSH09, SMS+13], design patterns knowledge [JGSH09, PG10], and the participants' style preference of identifiers [SM10a, SSGA12].

4 DISCUSSION

As shown in Section 3, we performed a quantitative analysis on 63 papers that used an eye-tracker in an experiment related to computer programming. In this section, we will discuss some of the quantitative findings reported in the mapping section, and highlight potential issues or problems to be further explored by researchers.

Trend. RQ1 confirmed our claim on the rising use and adaptation of eye-tracking technology in programming research. Considering that the first paper to use an eye-tracker in this field was published in 1990, it took over a decade for similar research in this area to emerge again. In the early 2000s, although a considerable number of publications started using eye-trackers to study the way programmers examine and comprehend programming materials, another portion of the

research published in that period focused on validating attention tracking technologies. RFV technology was validated and tested multiple times in order to determine its accuracy in measuring the viewing habits of programmers. RFV technology was not used in any studies after 2007, as its limitations and effect on the viewing habits of participants were questioned and reported in multiple publications [14, 15]. All the experiments reported after 2007 relied on the use of eye-trackers. In the second decade of the 21st century, the use of eye-trackers in programming research saw a noticeable hike, evident by the reported mapping result in RQ1, where more than 62% (39 papers) of the papers mapped in this survey were published in the last 5 years. This result, although it does not provide a qualitative argument to the validity of eye-tracking in programming research, can be quantitatively interpreted as an argument in favor of using eye-trackers in programming studies, as researches realized its potential and advantages.

Categorization. RQ2 focused on categorizing the collected research papers to identify areas that might be suitable for qualitative systematic literature reviews, or where more primary studies are needed. We identified five classes based on the materials used (code or non-code), the programming task given to participants (comprehension, debugging, or tractability), or the number of programmers working on the task (individual or collaborative programming). More than 80% of the reported studies used a source code or a programming language, while only 11 studies used diagrams or non-code stimuli. Code comprehension and debugging tasks have been studied in more than 70% of the collected papers, while collaborative programming and traceability had the least number of publications, with five and three papers, respectively. The lack of studies on the latter type of problems can be due to the difficulty of setting up the eye-tracking environment for such tasks [27, 54]. In code and non-code comprehension tasks, as well as debugging tasks, participants are shown a code or a diagram that was selected or modified to fit on one screen, and avoid any scrolling that might affect the eye-tracking metrics, such as scan-paths or heat maps. In order to avoid any noise in eye-tracking data, the stimulus needed to be shown in a full view on one display, and in most reported experiments, participants did not use a mouse to navigate through stimuli. This limitation and possible noise to the eye-tracking data posed a challenge in setting up an environment suitable to perform collaborative programming and traceability experiments. In collaborative tasks, the viewing patterns of pair programmers are monitored on either a distributed or shared display. As eye-trackers are designed to follow the eye movement of a single participant sitting in front of a display, collaborative research required a special setup and modified equipment [54]. In one particular case, [27] stated:

“We would have liked to use two of the much cheaper Tobii EyeX trackers, but it is not possible to put more than one EyeX on a computer.” (p. 6248)

Instead of using another Tobii EyeX tracker, [27] had to use Tobii TX300 which costs around 300 times more.

Materials Preference and Targeted Audience. RQ3 and RQ4 reported on the materials and participants used in an eye-tracking experiment in programming research. While students were subjects of eye-tracking experiments in programming, there does not seem to be an agreement on the optimal sample size for programming study with eye-trackers, and in most cases, the sample size and diversity was highlighted as an issue or cited as a possibility for future improvements. The most common sample size used in the studies was 15, followed by 20, and the most common participants were students and faculty members, which reflects that the targeted audience for the majority of these studies is in the academic field.

The most commonly used stimuli in the experiments were written using the Java programming language. This can be directly related to the fact that Java is the language used in a majority of introductory programming courses. Since a vast majority of the reported experiments were conducted in an institute of higher learning and used students as subjects, Java seems to be a rational choice for an experiment in programming research, which will allow the findings to be generalized and adapted by others. While researchers in the reported studies focused more on the students' reading patterns and viewing habits, fewer researchers tried to compare multiple programming languages. Only one paper compared C++ and Python programming languages to find if the students' performance differs between the two languages [80]. The selection of Python programming language was interesting as it follows a different syntax and structure of programming than the traditional Java and C++ languages. However, it was stated by the authors that the sample size for the Python group was smaller due to the fact that it is not a language students learned as an offered course in their institution. Such comparison and possible findings from comparing two different languages can have an impact on the programming language selected to teach introductory programming courses. Another comparison that lacked in terms of quantity of publication was comparing the differences between reading a natural text and a source code, with only two papers published on such a comparison. Such comparison between the reading patterns of natural text and a source code can possibly be linked to the research of comparing programming languages with different structure and syntax, in order to select a language that is easier and more suitable for introductory courses [20, 50]. One recent paper used Scratch programming tools to examine the differences between kids and teens [PSGJ17].

While most of the papers we reviewed tried to examine the usefulness of code representation, no other paper examined an alternative way of programming, such as Scratch. Another noticeable lack or absence in use of materials is the evaluation of the methods used to teach students the fundamentals of programming and problem solving skills, such as pseudocode and flowcharts. Only one paper used pseudocode and flowcharts as stimuli in an eye-tracking experiment [3]. However, no experiment showed the effect of using these methods on the students' ability to perform a programming task. As these methods are some of the first techniques students learn in introductory programming courses, their effectiveness and usefulness remains untested by an eye-tracking experiment.

Eye-trackers. The early work on using eye-tracking in computer programming was done on the basis that switching behavior during program comprehension can be effectively monitored using an eye-tracker [7], and the cognitive processes of the subjects can be studied with the support of their focused attention [47]. Eye-trackers used in programming research have some limitations. For example, the use of large materials that require scrolling can introduce noise to the collected data, hence inaccurate eye-tracking metrics. Another issue with eye-trackers is related to their limitation on tracking the eye movement of a single participant at a time, which may affect the studies of collaborative programming. In 2008, [54] stated that

“The currently available eye-tracking technologies do not provide solutions for studies on collaborative aspects of visual attention.” ([54], p. 42)

Another issue that eye-trackers had in early years of research is their intrusiveness when mounted on the participants' head, and the effect it might have on creating a comfortable work environment for the subject. An important study evaluating the performance of three types of eye-tracking devices was conducted by [47], using Tobii 1750, ASL 504, and the head-mounted ASL 501 tracker. Given the time needed to set up each device, [47] found that the head-mounted ASL 501 required twice as much time as the stationary devices. Although less than 10% of invalid

data was reported by all the devices, in terms of accuracy, Tobii was the most accurate, followed by the ASL 504 and ASL 501, respectively. The head-mounted version of the eye-tracking device was found to be obtrusive, less accurate, and required a longer setup time. As the RFV was abandoned by researchers in this field after 2007, the use of eye-trackers was focused on its accuracy and effectiveness in measuring eye movements, as well as its usability in an experimental environment. The early version of head-mounted eye-trackers became obsolete in later years, with Tobii devices leading the way in producing non-intrusive, accurate, and easy-to-use eye-trackers. However, one of the issues research might still face when deciding to use an eye-tracker in their study is the cost of such devices. Even though eye-trackers have evolved considerably over the years in terms of design and abilities, they are still not easy to be obtained by all researchers due to high prices. This study focuses on the quantitative reporting of eye-trackers in programming research, and a qualitative assessment of the effectiveness and usability of eye-trackers in programming studies needs to be explored and discussed in a systematic review of primary studies.

Metrics and Variables. Eye-tracking metrics, reported in RQ6, can be a variation of calculations based on duration, count, and sequence of fixations collected by an eye-tracker. A noticeable issue in the reported experiments is the lack of a terminology and inconsistent names of metrics, where a considerable number of papers reported the same measure with different names. While some papers measured fixation duration, others called it fixation time or fixation intensity. Similarly, some researchers used fixation count to evaluate their variables, while others used average fixation count or fixation rate. In this study, we can only report on the quantitative values of the eye-tracking metrics used by researchers in programming experiments, while leaving the terminology and guidelines of using these metrics to be reported in a qualitative review. RQ6 also maps the set of variables used in the selected studies, and can help identify any possible relationships or links between the dependent or independent variables used in the reported studies. It can be seen in the mapping of RQ6 that, while dependent variables were related in most parts to the performance and behavior of the participants, independent variables were more related to participants' background and environmental setup of the experiment. The most common dependent variable used in the studies was time. Time in its variations was used in every class of tasks to identify and measure the performance and viewing habits of participants. Albeit measuring detection, scan, response, fixation, or completion time, this most commonly measured variable can be affected by a variety of independent variables:

- Measuring *Completion time* can be affected by the experience of participants, their programming background, or programming language preference.
- Since *Fixation time* as a dependent variable measures the time participants spend focusing on a specific AoI, it is safe to assume that the selected AoI can be an independent variable that could affect fixation time. This can be verified by examining studies where the defined AoIs were SCEs, and while fixation time was a dependent variable, the different types of SCEs and code elements were the independent variables in those studies.
- *Scan time* and *Detection time* were used mostly in debugging studies where the type or the presence of errors and defects was an independent variable.
- *Completion time* did not show an associated independent variable that stands out from the reported set. It was measured along with various independent variables such as code or algorithm type variable, experience variable, or representation type variable.

One noticeable dependent variable is switching behavior, which was mostly used in debugging studies. Switching variable was commonly used in studies that used multiple representations or examined a visualization tool in order to measure its effectiveness in helping programmers

improve their performance. These tools usually represent the code in multiple forms and include a visualization field that aims at improving the understanding and simplifying the logic flow of the source code. Switching as a dependent variable was also used in studies that validated the effects of RFV on participants behavior, and where the RFV display blurring conditions were an independent variable. Visual effort as a dependent variable measures participants' strain and focused attention to understand a specific AoI. This variable was most commonly used in both code and non-code comprehension tasks. Comprehension tasks tried to understand the viewing habits of participants, and measure the parts of a code or a diagram that required more focused attention and effort to understand. An independent variable that was associated with visual effort is the participants' preference, noticeably in studies where identifier style or programming language preference was set as independent variables. In other words, visual effort as a dependent variable can be affected by the participants' familiarity with the materials and stimulus used in the experiment.

Out of the 63 mapped studies, only 13 study (25%) stated mitigating variables that had a lesser impact or milder effect on the relation between dependent and independent variables. The reported mitigating variables in the mapped studies were related to the participants' background and the materials used. In seven of the studies where visual effort and performance were the dependent variables, participants' experience was reported as a mitigating variable [SM10a, ASGA12, SFM12, BDL+13, SMS+13, SJAP13, TFSL14]. This can be confusing, as one of the most commonly used independent variables in the studies was participants' level of experience, and it has been associated with visual effort and performance variables in other studies. The way to determine if experience is an independent or a mitigating variable can be related to the experimental setup, objective of the study, and the analysis of the dependent variables. This can be further clarified by looking into two examples: When [JGSH09, SMS+13] tried to measure visual effort (dependent variable), then the representation of the stimuli was an independent variable that can affect the participants' effort, while the participants' experience and familiarity with this representation were mitigating variables. In other words, while the type of representation may affect visual effort, a participant who is familiar with the way the stimuli were represented may not be impacted by the representation, and hence require lesser visual effort than those who are not familiar with it, and vice versa. Similarly, when the identifier style was an independent variable that may affect the visual effort and performance (dependent variables) of the participants in [SM10a, SSGA12], then the participants' preference for identifier style and their experience were mitigating variables.

While the link between eye-tracking data and the cognitive load is still an ongoing topic of research, some researchers try to provide possible interpretations of the metrics calculated from eye-tracking data. In the reported studies, variables such as performance and accuracy were mostly calculated by evaluating the answers participants provided, or the ability to find bug(s) in source code. On the other hand, other variables were evaluated based on the eye-tracking metrics and the measures related to eye-tracking data, such as visual effort and viewing patterns. In a recent review, [25] provided a possible interpretation of eye-tracking metrics in relation to information visualization. Similarly, [83] discusses the relation between eye-tracking data and the cognitive load in relation to visual computing. Based on the work of [25] and [83], and from the literature reviewed in this study, we will try to briefly touch on the relation between some variables and the eye-tracking metrics and their relation to programming, while leaving an extended study on this subject for possible future qualitative review. For instance, viewing patterns variable can be studied using the eye-tracking metrics of attention switching and scan-path. While attention switching provides information on the areas that attracted the most attention from participants, it also can be used to interpret a possible link or strong correlation between two AoIs. As mentioned in RQ6, attention switching was mostly used in studies where researchers try to examine the effectiveness of multiple representations, or the effect different presentations of a source-code or a diagram can

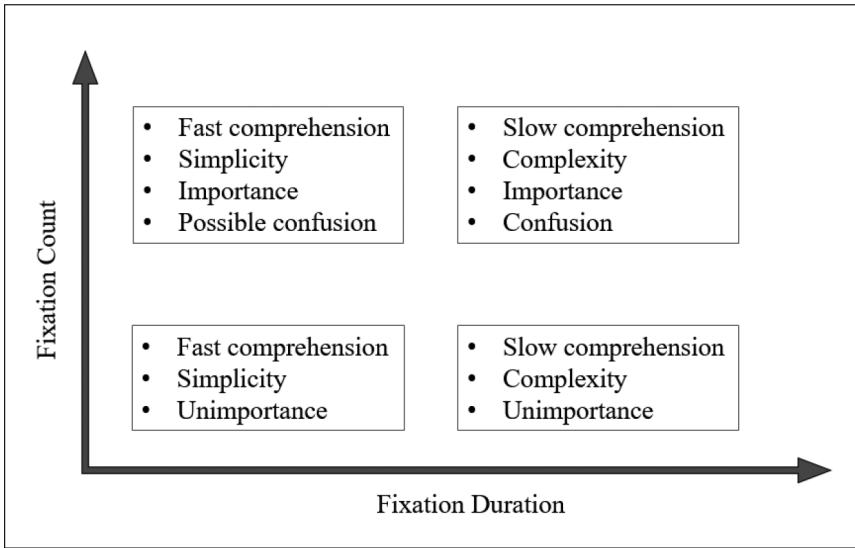


Fig. 14. Possible interpretations of visual effort based on fixation metrics.

have on the participants' performance. Scan-path provides details about the order in which the participants viewed the stimuli, hence the way a participant follows the logic of a source-code or how they jump from one AoI to another can be highlighted and studied. Hence, attention switching and scan-path can be used to evaluate the viewing patterns variable.

The variable of visual effort in programming studies can be interpreted using fixation duration and fixation count metrics. Figure 14 shows the possible interpretations of visual effort on an AoI based on the fixation metrics. While a high fixation duration on an AoI can mean difficulty to understand, complex section, impotence, or notability, a low fixation duration can mean simplicity of the AoI, unimportance, or the participants' fast comprehension. Similarly, visual effort can relate to fixation count, that is, the more frequently a participant keeps visiting an AoI, the more significant and meaningful that area could be. Thus, the relationship between these measures indicates that interpretation of results based on these matrices needs to be done carefully to avoid misleading discussions.

The interpretation of visual effort variable can be simplified as follows:

- Low fixation count and low time indicate less effort [SSGA12, SMS+13].
- High fixation count and more time indicate more effort [TFSL14, SM10a].
- Long fixations on an AoI mean more effort is needed, while AoIs with shorter fixations are more efficient as they require less effort [ASGA15].

5 THREATS TO VALIDITY

Possible threats to the validity of this mapping study are publication selection, data extraction inexactness, and misclassification.

Publication Selection. In conducting this study, we focused our search for papers on some of the most common electronic databases that are considered leaders in computer science and computer engineering, using sequences of keywords as a search query. We did not examine printed materials, gray literature, or publications that were not written in English in our survey. In principle, the

inclusion of such materials can provide more information and data for the survey, and allow more general results to be drawn. Our research questions were derived from our aim to conduct an eye-tracking study, and therefore, the papers selected for the survey followed a pre-defined set of questions to avoid any bias. Selection of publications involved the first and second authors of this article, with inclusion and exclusion reasons stated as discussed in Section 2.2. To avoid possible threats to the statistics and frequencies reported in this survey, duplicated experiments that were published in multiple articles or conferences were not included in the analysis. The second author initially compared the repeated experiments in order to choose one paper to report in this survey, and the first author confirmed the selection.

Extraction of Information and Data. Data and information included in this survey were extracted by two researchers, one article at a time. Both researchers had to reach a level of agreement on the collected data. Any inaccuracy in the data can be a result of the absence of any guideline or methodology for conducting such experiments in this field of study. The most commonly cited guideline in the surveyed papers was [39]. However, not every experiment followed the same style of reporting, which may cause inaccuracy in data collection.

Categorization of Papers. The papers included in this survey were categorized into one of five classes: code comprehension, debugging, non-code comprehension, collaborative programming, and traceability. The process of classifying a paper was based on either the task participants were asked to perform (comprehend, debug, or trace), the material used (source code or non-code), or the setup of the environment (single or pair). Each paper's classification was confirmed and checked by the first author, while disagreements were discussed until a decision was made. In some cases the classification was easy, especially for articles that used non-code stimuli, did traceability study, or used two programmers at the same time in the experiment. However, some of the difficulties in classifying an article were due to the similarities between comprehension and debugging tasks. Generally, in order to debug a code, programmers need to understand the code first, and then try to find any bug(s) in the program. In some experiments, researchers may ask participants to assess the correctness of a code, however, the main objective of the study may not be the debugging performance, but rather the way participants view and try to understand the code. In any confusion similar to this, we relied on the title, the objective of the experiment stated by the authors, or the analysis of the collected data in rare cases. We looked into the data analysis to see if the author considered the debugging performance in their analysis, or just studied the viewing habit and the comprehension of the code.

6 CONCLUSION

This systematic mapping study reported on papers that used eye-tracking technology in experiments related to programming. Using online electronic databases that are most commonly used by computer scientists and software engineers, we were able to identify 63 papers published between 1990 and June 2017, which used eye-tracking in programming research. Although the first two studies in our reporting, [CS90] and [CW02], had more than a 10 year gap between them, the use of this technology has seen a rise in recent years, with more than 60% of the 63 reported papers published after 2011. We categorized the collected studies into five classes based on the programming area and task: code comprehension, debugging, non-code comprehension, collaborative programming, and requirements traceability. We found that the majority of the reported experiments were conducted on code comprehension and debugging, with fewer papers reporting on collaborative programming and requirements traceability. The fewer number of studies on collaborative programming and traceability studies can be due to some limitations in the eye-tracking technology. In the case of collaborative programming, additional setup and modification to the

environment is required as current eye-tracking devices are not designed to track the visual attention of more than one participant at a time. As for traceability studies, which require large source codes, current eye-tracking devices require a stimuli to be presented on a full screen with no scrolling, which might cause some noise in the eye-tracking data. Researchers took into consideration the possibility of noise in the eye-tracking data due to large source codes or diagrams. Therefore, this limitation of eye-tracking devices may have affected the selection of stimuli. An important part of conducting an eye-tracking experiment related to programming was the selection of participants and materials used to evaluate the participant visual effort and performance. We found the most commonly used materials by researchers in an experiment related to programming was the Java programming language, followed by the UML class diagrams. We relate this common usage of Java in eye-tracking experiments to the participants' familiarity with the language as it is the most common language students learn in their introductory programming course. In this survey, we found that only one paper did a comparison of two different programming languages (C++ and Python) in order to evaluate their effect on students' performance [TFSL14]. Similarly, we found that only one paper studied students' behavior when using flowcharts and pseudocode [ASB+15], which are fundamental tools that novice programmers learn in their introductory programming course. In terms of eye-tracking metrics, an important finding of this study is the lack of a methodology in conducting such experiments, which resulted in terminological inconsistency in the names and formulas of metrics used, which needs to be addressed in a future qualitative review.

APPENDIXES

A FIGURES AND TABLES

Table 9. Details of Programming Languages and Materials Used for Each Paper (Sorted by Year)

| Paper | Task | Materials | Details of the materials used. |
|-----------|---------------------------|-------------|---|
| [CS90] | Program Comprehension | Pascal | Short but complex algorithm. |
| [CSW02] | Program Comprehension | Algorithm | N/A |
| [RCdBL02] | Debugging | Java | 3 debugging sessions for each participant to find as many errors as they could. |
| [RLCdB02] | Debugging | Java | 5 debugging sessions; 1 was a warm-up and four main sessions followed. |
| [BT04a] | Debugging | Java | 3 debugging sessions; 1 was a warm-up performed |
| [BT04b] | Debugging | Java | 3 programs; 1 was a warm-up. |
| [NS04] | Program Comprehension | N/A | 6 short programs in total; two for each eye-tracker. |
| [SB04] | Collaborative Programming | Java | 3 programs that can fit on one screen, and suitable for novice programmers. |
| [NS05] | Program Comprehension | Pascal | 7 short programs, between 11 and 29 Lines of Code (LoC). |
| [NS06] | Program Comprehension | Java | 2 Java programs with 63 and 58 LoC. |
| [BMST06a] | Program Comprehension | Java | 3 short Java programs, with 15, 34, and 38 (LoC). |
| [AC06] | Program Comprehension | Java | 12, recursive and non-recursive versions of six algorithms. |
| [Gu06] | Comprehension | UML diagram | 2 different class diagrams from two different programs. |
| [UNMM06] | Program Comprehension | C | 6 small-scale programs. Each program had 1 logic defect. |

(Continued)

Table 9. Continued

| Paper | Task | Materials | Details of the materials used. |
|------------|---------------------------|------------------------------|--|
| [BT07a] | Debugging | Java | 1 Java program that contained four non-syntactic errors. |
| [BT07b] | Debugging | Java | Similar to the experiment from [RCdBL02]. |
| [YKM07] | Comprehension | UML diagram | 6 modules of a program presented using 3 different types of layouts. |
| [Dub09] | Program Comprehension | Scala | 3 algorithms are implementations of relational algebra operators. |
| [PBT09] | Collaborative Programming | Software Development Project | N/A |
| [JGSH09] | Comprehension | UML diagram | Diagrams from 3 open-source programs. |
| [PG10] | Comprehension | UML diagram | 2 diagrams of 15 and 40 classes. |
| [SM10b] | Comprehension | UML diagram | 8 diagrams investigating four design patterns on 3 systems: JUnit, JHotdraw, and Qt. |
| [SM10a] | Program Comprehension | Phrases | 8 phrases. |
| [BSB11] | Program Comprehension | NLT and Java | 11 small programs with a varying complexity with multiple choice questions about the programs. |
| [ASGA12] | Traceability | Java | 6 short source-codes; each can fit on one screen. |
| [SFM12] | Debugging | C | 4 C language source code snippets, each loaded with a single logical defect. |
| [JN12] | Collaborative Programming | Java | 1 Java program of 300 lines. |
| [HN12] | Debugging | Java | 1 Bubble sort program consisting of two classes and loaded with three bugs. |
| [SSGA12] | Program Comprehension | Java | 3 small Java programs with 30, 36, and 40 LoC. |
| [SSVdP+12] | Comprehension | UML diagram | 3 UML class diagrams. |
| [Bed12] | Debugging | Java | 3 programs, 1 warm-up (data not considered for analysis), and 2 main programs with about 100 LoC each. |
| [BDL+13] | Program Comprehension | C++ | 1 code with two functions that had 6 and 15 identifiers. |
| [CTKT13] | Comprehension | ER Diagram | 1 Entity Relationship Diagram seeded with 17 defects. |
| [DcI13] | program Comprehension | .NET | 1 e-commerce application (around 1000 LoC) in Microsoft ASP.NET that is designed to manage the product inventory and orders received by a small-scale business enterprise. |
| [SMS+13] | Comprehension | TROPOS | 3 requirements comprehension tasks. |
| [SJAP13] | Debugging | Java | 1 open source system with around 60 KLoC. |
| [CL13] | Debugging | C# and Java | 27 trials, 9 for each error type: lexical, logical, and syntactic. |
| [HLL+13] | Debugging | C | 2 programs with 100 LoC, one iterative and the other recursive; each had three semantic or syntactic bugs. |

(Continued)

Table 9. Continued

| Paper | Task | Materials | Details of the materials used. |
|--------------|---------------------------|--------------------------|--|
| [TFSL14] | Debugging | C++ and Python | 10 codes, 5 in C++ and 5 in Python; each subject was tested on either C++ or Python programs but not both. |
| [WSSK14] | Traceability | N/A | N/A. |
| [Loh14] | Program Comprehension | Java | 1 source code of the Apache River project. |
| [FBM+14] | Program Comprehension | C# | 8 codes out of 10 comprehension tasks (2 warm-up not considered). |
| [DSLS+14] | Comprehension | UML diagram | Class diagrams of 2 different programs: JTable and JFreeChart. |
| [BSS+14] | Program Comprehension | Java | 1 program that calculated the area of a rectangle. |
| [ASGA15] | Traceability | Java | 6 pieces of code with varying lengths. |
| [ASB+15] | Program Comprehension | Pseudocode and Flowchart | 2 algorithms: one presented in a pseudocode and the other a flowchart. |
| [SUGGR15] | Comprehension | Graphics | 1 program: JHotDraw. |
| [MD15] | Program Comprehension | C++ | 4 short codes. |
| [MGRB15] | Collaborative programming | Java | 1 program. |
| [RLMM15] | Program Comprehension | Java | 67 methods from 6 different applications. |
| [BBB+15] | Program Comprehension | NLT and Java | 17 natural language trials and 101 source-code trials from novices, and 21 Java trials from experts. |
| [JF15] | Program Comprehension | N/A | 2 programs from the image processing domain. |
| [BdP16] | Program Comprehension | .NET | Snippets of code were presented either in black-and-white, or in color. |
| [GWMP16] | Debugging | Requirements Documents | 2 documents: one was 11 pages seeded with 30 realistic faults and the other 14 pages with 34 faults. |
| [LWH+16] | Debugging | C | 2 programs: one iterative and the other recursive. Each program has three bugs. |
| [MDPVRDVI16] | Program Comprehension | N/A | |
| [NHMG16] | Debugging | C | 8 short codes; each has 15 lines. |
| [PLS+16] | Debugging | Java | 2 Java programs seeded with bugs. |
| [BSL+17] | Debugging | Java | N/A |
| [DB17] | Collaborative Programming | C# | N/A |
| [MNH+17] | Debugging | Java | N/A |
| [PIS17] | Program Comprehension | NLT and C++ | 7 source codes and 3 natural language texts. |
| [PSGJ17] | Program Comprehension | Scratch | N/A |

Table 10. Details of Participants Used

| Paper | Participants | Size | Split | Split Details | Gender |
|------------|--------------|------|---------------------------------------|---|-----------|
| [CS90] | S&FM | 19 | Novice vs. Experienced | 10 low-experience, 9 high-experience | N/A |
| [CSW02] | S&FM | 19 | Novice vs. Experienced | 9 novice, 10 experienced | N/A |
| [RCdBL02] | S&P | 5 | Experienced programmers | 4 students, 1 professional | N/A |
| [RLCdB02] | S | 49 | Less experienced vs. More experienced | N/A | N/A |
| [BT04a] | S& FM | 10 | Programming experience varied | N/A | 1 F, 9 M |
| [BT04b] | S&FM | 18 | Programming experience varied | N/A | 3 F, 15 M |
| [NS04] | S | 12 | Programming experience varied | N/A | 4 F, 8 M |
| [SB04] | P | 10 | Professional programmers | 4 for the first phase, 6 for the second | 1 F, 9 M |
| [NS05] | S | 12 | Students | N/A | 5 F, 7 M |
| [NS06] | S | 16 | Animation group vs. Static group | N/A | 5 F, 11 M |
| [BMST06a] | S | 18 | Novice VSvs. Intermediates | 8 novice, 8 intermediate | 3 F, 13 M |
| [AC06] | S | 15 | N/A | N/A | N/A |
| [Gu06] | S | 12 | N/A | N/A | N/A |
| [UNMM06] | S | 5 | N/A | N/A | N/A |
| [BT07a] | S&FM | 18 | Novice vs. Experienced | 6 novice, 8 expert | N/A |
| [BT07b] | S&FM | 19 | Novice vs. Experienced | 10 novice, 8 expert | 3 F, 15 M |
| [YKM07] | S&FM | 12 | Performance | 3 UADA, 1 UEDI, 3 UEDK, and 5 UEDE | N/A |
| [Dub09] | S | 12 | N/A | N/A | N/A |
| [PBT09] | N/A | 2 | Pair of programmers | N/A | N/A |
| [JGSH09] | S | 24 | Students | 7 post-graduate, 17 graduate students | N/A |
| [PG10] | S | 24 | N/A | N/A | N/A |
| [SM10b] | S&FM | 15 | Novice vs. Experienced | 7 novice, 8 expert | 2 F, 13 M |
| [SM10a] | S&FM | 15 | Identifier Styles | 6 camel-case, 7 underscore, and 2 had no preference | N/A |
| [BSB11] | S | 15 | Programming experience varied | N/A | N/A |
| [ASGA12] | S | 26 | N/A | N/A | 7 F, 19 M |
| [SFM12] | S | 15 | Novice vs. Experienced | 8 expert, 1 excluded, and 7 novice | 2 F, 12 M |
| [JN12] | S | 82 | Pair programmers | 40 pairs | N/A |
| [HN12] | S | 19 | Programming experience varied | 10 High, 10 Low | 2 F, 17 M |
| [SSGA12] | S | 24 | Gender | 7 camel-case, 17 underscore | 9 f, 15 M |
| [SSVdP+12] | S&P | 21 | Novice vs. Experienced | 12 students, 9 practitioners | 1 F, 20 M |
| [Bed12] | S&FM | 18 | Novice vs. Experienced | 6 novice, 8 expert | N/A |

(Continued)

Table 10. Continued

| Paper | Participants | Size | Split | Split Details | Gender |
|--------------|--------------|------|-------------------------------------|--|------------|
| [BDL+13] | S&FM | 15 | Identifier Styles | 40% camel-case, 47% underscore, and 13% no preference | 2 Female |
| [CTKT13] | P | 4 | N/A | N/A | N/A |
| [Dcl13] | P | 13 | N/A | 5 control group, 8 experiment group | N/A |
| [SMS+13] | S | 28 | N/A | N/A | 12 F, 16 M |
| [SJAP13] | S | 20 | Novice vs. Experienced | 10 novice, 10 expert | N/A |
| [CL13] | S | 9 | N/A | N/A | 2 F, 7 M |
| [HLL+13] | S | 25 | Gender | 12 female, 13 male | 12 F, 13 M |
| [TFSL14] | S | 38 | Experience and programming language | C++ group (17 novice, 8 expert); Python group (10 novice, 3 expert) | N/A |
| [WSSK14] | S&FM | 8 | N/A | 4 expert developers | N/A |
| [Loh14] | S&P&P | 30 | N/A | N/A | N/A |
| [FBM+14] | P | 15 | N/A | N/A | 1 F, 14 M |
| [DLS+14] | S&P | 23 | N/A | N/A | N/A |
| [BSS+14] | P | 2 | N/A | N/A | N/A |
| [ASGA15] | S | 14 | N/A | N/A | N/A |
| [ASB+15] | S | 26 | Effectiveness | 13 effective, 13 non-effective | 30 F, 73 M |
| [SUGGR15] | S&P | 20 | Experience | 4 novice, 10 intermediate, and 6 expert | 6 F, 14 M |
| [MD15] | N/A | 3 | N/A | N/A | N/A |
| [MGRB15] | S | 10 | Pair programmers | 5 pairs | N/A |
| [RLMM15] | P | 10 | N/A | N/A | N/A |
| [BBB+15] | S&P | 20 | Novice vs. Expert | 14 novice students, 6 professional | 8 F, 12 M |
| [JF15] | S&FM | 20 | N/A | N/A | 3 F, 17 M |
| [BdP16] | S&P | 38 | Representation and Experience | Black-and-white 14, and color 17 and 4 experts | N/A |
| [GWMP16] | S | 13 | N/A | N/A | N/A |
| [LWH+16] | S | 38 | Performance and program type | Iterative (10 high performance, 28 low performance); recursive (12 high performance, 26 low performance) | 16 F, 22 M |
| [MDPVRDVI16] | S | 19 | Representation style | 13 Experiment1 static representation, 6 Experiment, 2 dynamic representation | N/A |
| [NHMG16] | S | 24 | Experience | 15 novice and 9 experts | N/A |
| [PLS+16] | S | 20 | N/A | N/A | N/A |
| [BSL+17] | S | 56 | N/A | N/A | 10 F, 46 M |
| [DB17] | P | 2 | N/A | N/A | 1 F, 1 M |
| [MNH+17] | S | 20 | N/A | N/A | N/A |
| [PIS17] | S | 33 | Experience | Novice and non-novice | 15 F, 18 M |
| [PSGJ17] | S | 44 | Age | Kids and teens | 12 F, 32 M |

Table 11. Details of the Eye-tracking Devices Used in Each Experiment

| Paper | Tracker | Specifications | Duration |
|-----------|----------------------------------|---|--|
| [CS90] | ASL | N/A | N/A |
| [CSW02] | ASL | N/A | N/A |
| [RCdBL02] | RFV | A modified version RFV. | N/A |
| [RLCdB02] | RFV | N/A | 10 minutes debugging session for each program. |
| [BT04a] | RFV and Tobii 1750 | Sampling-rate set to 30Hz for Tobii. | 10 minutes debugging session for each program. |
| [BT04b] | RFV and Tobii 1750 | Sampling-rate set to 30Hz for Tobii. | 10 minutes debugging session for each program. |
| [NS04] | Tobii 1750, ASL 504, and ASL 501 | ASL 504 and ASL 501 Head-Mounted Optics, used with PlanAni animator. | N/A |
| [SB04] | ISCAN RK-726PCI | N/A | 10 minutes debugging session for each program. |
| [NS05] | Tobii | N/A | N/A |
| [NS06] | Tobii | 1, 280 × 1, 024 resolution display with PlanAni animator. | Around 2 hours session for the experiment. |
| [BMST06a] | Tobii 1750 | Sampling rate 50Hz, 17" TFT display resolution of 1, 024 × 768 and using Jeliot 3. | N/A |
| [AC06] | ASL | N/A | 1 hour or less. |
| [Gu06] | EyeLink II | N/A | N/A |
| [UNMM06] | EMR-NC | EMR-NC with sampling rate of 30Hz, a 21" LCD with resolutions set at 1, 024 × 768. | A 5 minute review session, or the subject finds all the defects. |
| [BT07b] | RFV 2.1 and Tobii 1750 | Restricted Focus Viewer (RFV) 2.1 and Tobii 1750 with sampling rate of 30Hz on 17" display resolution of 1, 280 × 1, 024. | 10 minutes to debug the program. |
| [BT07a] | RFV 2.1 and Tobii 1750 | Same setup from [BT07b]. | 10 minutes. |
| [YKM07] | Tobii 1750 | Sampling rate of 50MHZ and less than 0.5 degrees error rate. | Between 10 and 20 minutes. |
| [Dub09] | Tobii 1750 | N/A | Between 45 minutes and 1 hour. |
| [PBT09] | ASL 504 and Tobii 1750 | Complex setup of a system to capture eye movement data of pair programmers. | N/A |
| [JGSH09] | EyeLink II | N/A | An experiment took about 1 hour for each subject. |
| [PG10] | EyeLink II | A high-resolution EyeLink II and fast data rate of 500 samples per second. | Between 12 and 15 minutes. |
| [SM10b] | Tobii 1750 | Temporal resolution was set at 50Hz, on a 17" TFT-LCD screen resolution 1, 024 × 768. | The study was designed to be completed in 20 minutes or less. |
| [SM10a] | Tobii 1750 | Temporal resolution was set at 50Hz, on a 17" TFT-LCD screen resolution 1, 024 × 768. | The experiment took 13 minutes on average. |
| [BSB11] | Tobii T120 | Sampling rate of 120Hz, with Tobii Studio 2.0.6. | N/A |
| [ASGA12] | FaceLAB | Two 27" LCD monitors, one with screen resolution 1, 920 × 1, 080 | Average 20 minutes, including setting up eye-tracking system. |

(Continued)

Table 11. Continued

| Paper | Tracker | Specifications | Duration |
|------------|-------------------------------|--|--|
| [SFM12] | Tobii 1750 | 17" LCD screen, resolution was set to 1,024 × 768. | 5 minutes or less. |
| [JN12] | Tobii 1750 | Two synchronized 50Hz Tobii eye-trackers. | Approximately 45 minutes. |
| [HN12] | Tobii T60 XL | With 60Hz sampling rate and Tobii Studio 2.1 software. | A 15 minute time limit. |
| [SSGA12] | FaceLAB | Screen resolution is 1,920 × 1,080, and Taupe system to analyze the collected data. | 25 minutes in average, 11 hours total time of eye-tracking. |
| [SSVdP+12] | EyeLink II | With <i>Taupe</i> system to analyze the data. | N/A |
| [Bed12] | Tobii 1750 | Sampling rate of 50Hz. | N/A |
| [BDL+13] | Tobii 1750 | Temporal resolution was set at 50Hz, on a 17" TFT-LCD screen resolution 1,024 × 768, and ClearView software. | N/A |
| [CTKT13] | Tobii | 50Hz sampling rate, on a 17" monitor with a resolution of 1,024 × 728. | N/A |
| [DcI13] | Tobii 1750 | N/A | N/A |
| [SMS+13] | FaceLAB | Two 27" LCD monitors. | 10 minutes for each session. Total eye-tracking time was 2.85 hours. |
| [SJAP13] | Tobii X60 | 60Hz sampling rate, a 24" monitor, and 1,920 × 1,080 screen resolution. | All tasks were selected to be solved in approximately half an hour. |
| [CL13] | Tobii T60 | N/A | N/A |
| [HLL+13] | Eyelink 1000 | Screen resolution was 1,024 × 768. | N/A |
| [TFSL14] | Mirametrix S2 | An average accuracy of 0.5–1.0 degrees with 60Hz rate. | N/A |
| [WSSK14] | Tobii X60 | 24" LCD monitor. | N/A |
| [Loh14] | SMI iView X Hi-Speed 1250 | Using a scan rate of 500Hz. | About 90 minutes per subject. |
| [FBM+14] | Tobii TX300 | Using a 300Hz tracking frequency, and 96 dpi 1920 × 1080 23" monitor. | 1.5 hour experiment. |
| [DSLS+14] | FaceLAB | N/A | N/A |
| [BSS+14] | SMI RED-m | 120Hz eye-tracker using the Ogama tracking software. | N/A |
| [ASGA15] | FaceLAB | Same setup from [ASGA12]. | N/A |
| [ASB+15] | SMI iView X Hi-Speed 500/1250 | With 500Hz time resolution, and the <i>SMI BeGaze 2.4</i> software to analyze the data. | N/A |
| [SUGGR15] | Tobii T60 XL | Sampling rate of 60Hz, on a 24" flat-panel screen. | 7–23 minutes to complete the experimental tasks. |
| [MD15] | Tobii X2-30 | At 30 samples/sec. | N/A |
| [MGRB15] | Tobii X60 | With Tobii Studio 3.0.2 software. | Time limited to 15 minutes for the session. |
| [RLMM15] | Tobii T300 | Sampling rate 120Hz, and a resolution of 1,920 × 1,080. | A 1 hour session. |
| [BBB+15] | SMI RED-m | Sample rate set at 120Hz, and <i>Ogama</i> tool to record the data. | N/A |

(Continued)

Table 11. Continued

| Paper | Tracker | Specifications | Duration |
|--------------|------------------------------|---|--|
| [JF15] | Eye Tribe | Sampling rate of 60Hz, 9 points calibration process, and <i>Ogama</i> tool to record the data on a 1, 280 × 1, 024 screen resolution. | N/A |
| [BdP16] | Tobii TX300 | Sampling rate of 300Hz, with Tobii Studio. | N/A |
| [GWMP16] | Eyelink 1000 | Using sampling rate of 250Hz, | N/A |
| [LWH+16] | Eyelink 1000 | Sampling rate set at 1,000Hz, on a 22" LCD monitor with 1, 024 × 768 resolution. | Time limited to 10 minutes for the session. |
| [MDPVRDVI16] | N/A | N/A | N/A |
| [NHMG16] | SMI | Sampling rate 250Hz. | Between 20 and 45 minutes. |
| [PLS+16] | Tobii X60 | | 20 minutes to debug the program |
| [BSL+17] | GazePoint GP3 | On a 24" monitor with 1,920×1,080 pixels, and using GP3 software. | N/A |
| [DB17] | Tobii TX300 and Tobii EyeX | TX300 on a 23" and 1,920×1,080 monitor, and Tobii EyeX on a 24" and 1,920×1,200 monitor with sampling rate 30Hz. | 15 minutes for each task. |
| [MNH+17] | Tobii EyeX | On a EyeInfo Framework, with accuracy error < 0.4 degrees. | Participants spent 15 minutes to debug 2 programs. |
| [PIS17] | Tobii X60 | With Tobii studio V2.3. | N/A |
| [PSGJ17] | SMI RED 250 and Tobii mobile | 4 SMI trackers, and one Tobii with sampling rate 60Hz. | N/A |

Table 12. All Variables Listed by the Study Papers

| Paper | Dependent Variables | Independent Variables | Mitigating Variables |
|-----------|---|--|----------------------|
| [CS90] | Number of fixations and fixation time | Subjects experience | N/A |
| [CSW02] | Performance and time | Expertise and line number | N/A |
| [RCdBL02] | Debugging performance | N/A | N/A |
| [RLCdB02] | Accuracy, fixation time, and switching frequency | Visualization modality (textual or graphical), visualization perspective (data structure or control-flow), type of error (data structure or control-flow), and programming experience (less experienced or more experienced) | N/A |
| [BT04a] | Accuracy, fixation time, and switching frequency | N/A | N/A |
| [BT04b] | Debugging performance, fixation time, and switching frequency | RFV restricting condition and measuring tool | N/A |
| [NS04] | Accuracy and time | The eye-tracking device used | N/A |
| [SB04] | Debugging performance | N/A | N/A |
| [NS05] | Performance and gaze location | Visualization tool | N/A |
| [NS06] | Performance and gaze location | Visualization tool | N/A |

(Continued)

Table 12. Continued

| Paper | Dependent Variables | Independent Variables | Mitigating Variables |
|------------|---|---|--|
| [BMST06a] | Number of fixations, fixation time, and number of switches | Experience | |
| [AC06] | Accuracy and number of lines | Algorithm type: recursive and non-recursive | Expertise |
| [Gu06] | N/A | N/A | N/A |
| [UNMM06] | Scan time and defect detection time | Presence of defects | N/A |
| [BT07b] | Debugging performance, fixation time, and switching frequency | RFV restricting condition and level of experience | N/A |
| [BT07a] | Fixation time and switching frequency | N/A | N/A |
| [YKM07] | Accuracy, response time, and effort | Layout (Orthogonal, Three-cluster, and Multiple-cluster) | N/A |
| [Dub09] | Performance and time | Code density and presence or absence of grounding hints | N/A |
| [PBT09] | Accuracy | Type of eye-tracker | N/A |
| [JGSH09] | Effort, number of fixations, and fixation time | Design (no pattern, canonical pattern and modified pattern) and task (program comprehension task and modification task) | UML knowledge and design patterns knowledge |
| [PG10] | Fixation time | Representations of variables and tasks | JHotDraw knowledge and design pattern knowledge |
| [SM10b] | Accuracy, time, relevance, and visual effort | Reading behaviour | N/A |
| [SM10a] | Accuracy, find time, and visual effort | Identifier style (Style) with two treatments: camel-case or underscore | Reading time, experience, phrase length, phrase origin, and style preference |
| [ASGA12] | Accuracy and fixation time | Source code entities (SCEs), e.g., class names, method names | Level of study and programming experience |
| [SFM12] | Scan time, defect detection time, accuracy, and visual effort | Defects | Experience |
| [JN12] | Speech, selection, and gaze cross-recurrence | Selection sharing | N/A |
| [HN12] | Gaze time and switching patterns | Programming experience, familiarity with jGRASP, and debugging performance | N/A |
| [SSGA12] | Accuracy, time (Speed), and visual effort | Gender and identifiers style | Study level and style preference |
| [SSVdP+12] | Accuracy, time, and effort | Status (practitioner, student) and Expertise (expert, novice) | Question precision (precise or not precise) |
| [Bed12] | Debugging performance, fixation time, switching frequency, and type of switches | Expertise | N/A |

(Continued)

Table 12. Continued

| Paper | Dependent Variables | Independent Variables | Mitigating Variables |
|-----------|--|--|--|
| [BDL+13] | Performance, time, and visual effort | Identifier style | Experience |
| [CTKT13] | Performance and defect detection difficulty | Search pattern | N/A |
| [Dc13] | Accuracy and task completion time | Visualization tool | N/A |
| [SMS+13] | Accuracy, time, and effort | Representation type (Graphic vs. text) | Level of study, level of experience in object-oriented modeling, level of UML knowledge, English language proficiency, and linguistic distance |
| [SJAP13] | Effectiveness (accuracy), efficiency (time), and visual effort | 3D visualization tool (SeeIT 3D, No- SeeIT 3D) | Expertise (novices and experts), University (X, Y, and Z) |
| [CL13] | Debugging performance (accuracy and reaction time), eye gaze, and mouse cursor behavior (frequency and duration) | N/A | N/A |
| [HLL+13] | Gaze sequences | Gender (male, female) and problem type (Iterative, Recursive) | N/A |
| [TFSL14] | Effectiveness (accuracy), efficiency (time), and visual effort | Programming language (C++, Python) | Expertise |
| [WSSK14] | Performance | N/A | N/A |
| [Loh14] | Performance and time | Comprehension skill, expression type, activation of the target relation, and type of comprehension questions | N/A |
| [FBM+14] | Time and effort | Task difficulty | N/A |
| [DSLS+14] | Time and accuracy | Different variants of design pattern | N/A |
| [BSS+14] | Fixation time | N/A | N/A |
| [ASB+15] | Number of fixations and fixation time | Performance (effective and the non-effective) and type of representation (flowchart or pseudocode) | The content of the task (the instruction) and the answer (selection) |
| [SUGGR15] | Performance, accuracy, and response time | Expertise, and type of representation (Scatter plot, Treemap, Hierarchical Dependency Graph (HDG)) | N/A |
| [MD15] | | | |
| [MGRB15] | Inspection time and number of fixations | Level of knowledge and Llevel of practical experience | N/A |
| [RLMM15] | Gaze time, fixation, and regression counts | N/A | N/A |
| [BBB+15] | Fixation order, fixation time, and fixation location | Experience level (novice and expert) and materials (natural language text or source code) | N/A |

(Continued)

Table 12. Continued

| Paper | Dependent Variables | Independent Variables | Mitigating Variables |
|--------------|--|---|----------------------|
| [JF15] | Accuracy, completion time, and visual effort | Regularity of code (Regular vs. Non-Regular) | N/A |
| [BdP16] | Time, Viewing Pattern | Syntax highlighting | N/A |
| [GWMP16] | Time, Performance, Visual Effort | N/A | N/A |
| [LWH+16] | Visual attention, gaze paths | Program type (iterative and recursive) and Performance (low and high) | N/A |
| [MDPVRDVI16] | Time, Performance, Visual Effort | Experience, knowledge | N/A |
| [NHMG16] | Time, Performance, Visual Effort | Experience | N/A |
| [PLS+16] | Time, Viewing Pattern | Performance, Bug type | N/A |
| [BSL+17] | Performance | Knowledge of the language, Experience, and Familiarity with the type of error | N/A |
| [DB17] | Time, Viewing Pattern | N/A | N/A |
| [MNH+17] | Time, Performance, Viewing Pattern | Program variability (with or without) | N/A |
| [PIS17] | Time, Viewing Pattern | Experience, Representations (NLT or SC) | N/A |
| [PSGJ17] | Time, Viewing Pattern | Age (Kids or Teenager) | N/A |

B MAPPING STUDY PAPERS

[AC06] Christoph Aschwanden and Martha Crosby. Code scanning patterns in program comprehension. In *Proceedings of the 39th Hawaii International Conference on System Sciences*, 2006.

[ASB + 15] Magdalena Andrzejewska, Anna Stolinska, Wladyslaw Blasiak, Pawel Pkeczkowski, Roman Rosiek, Bozena Rozek, Mirosława Sajka, and Dariusz Wcislo. Eye-tracking verification of the strategy used to analyse algorithms expressed in a flowchart and pseudocode. *Interactive Learning Environments*, page 1–15, 2015.

[ASGA12] Nawazish Ali, Zohreh Sharafi, Y. Gueheneuc, and Giuliano Antoniol. An empirical study on requirements traceability using eye-tracking. In *Proceedings of the 28th IEEE International Conference on Software Maintenance (ICSM)*, pages 191–200. IEEE, 2012.

[ASGA15] Nasir Ali, Zohreh Sharafi, Yann-Gaël Guéhéneuc and Giuliano Antoniol. An empirical study on the importance of source code entities for requirements traceability. *Empirical Software Engineering* 20(2):442–478, 2015.

[BBB + 15] Teresa Busjahn, Roman Bednarik, Andrew Begel, Martha Crosby, James H. Paterson, Carsten Schulte, Bonita Sharif, and Sascha Tamm. Eye movements in code reading: Relaxing the linear order. In *Proceedings of the 2015 IEEE 23rd International Conference on Program Comprehension*, pages 255–265. IEEE Press, 2015.

[BDL + 13] Dave Binkley, Marcia Davis, Dawn Lawrie, Jonathan I. Maletic, Christopher Morrell, and Bonita Sharif. The impact of identifier style on effort and comprehension. *Empirical Software Engineering* 18(2):219–276, 2013.

[BdP16] Tanya Beelders and Jean-Pierre du Plessis. The influence of syntax highlighting on scanning and reading behavior for source code. In *Proceedings of the Annual Conference of the South African Institute of Computer Scientists and Information Technologists*, page 5. ACM, 2016.

[Bed12] Roman Bednarik. Expertise-dependent visual attention strategies develop over time during debugging with multiple code representations. *International Journal of Human-Computer Studies*, 70(2):143–155, 2012.

[BMST06a] Roman Bednarik, NIKO Myller, ERKKI Sutinen, and MARKKU Tukiainen. Analyzing individual differences in program comprehension with rich data. *Technology Instruction Cognition and Learning* 3(3/4):205–232, 2006.

[BSB11] Teresa Busjahn, Carsten Schulte, and Andreas Busjahn. Analysis of code reading to gain more insight in program comprehension. In *Proceedings of the 11th Koli Calling International Conference on Computing Education Research*, pages 1–9. ACM, 2011.

[BSL + 17] Titus Barik, Justin Smith, Kevin Lubick, Elisabeth Holmes, Jing Feng, Emerson Murphy-Hill, and Chris Parnin. Do developers read compiler error messages? In *Proceedings of the 39th International Conference on Software Engineering*, pages 575–585. IEEE Press, 2017.

[BSS + 14] Teresa Busjahn, Carsten Schulte, Bonita Sharif, Andrew Begel, Michael Hansen, Roman Bednarik, Paul Orlov, Petri Ihantola, Galina Shchekotova, Maria Antropova, et al. Eye tracking in computing education. In *Proceedings of the 10th Annual Conference on International Computing Education Research*, pages 3–10. ACM, 2014.

[BT04a] Roman Bednarik and Markku Tukiainen. Visual attention and representation switching in Java program debugging: A study using eye movement tracking. In *Proceedings of the 16th Annual Workshop of the Psychology of Programming Interest Group*, pages 159–169, 2004.

[BT04b] Roman Bednarik and Markku Tukiainen. Visual attention tracking during program debugging. In *Proceedings of the 3rd Nordic Conference on Human Computer Interaction*, pages 331–334. ACM, 2004.

[BT07a] Roman Bednarik and Markku Tukiainen. Analysing and interpreting quantitative eye-tracking data in studies of programming: Phases of debugging with multiple representations. In *Proceedings of the 19th Annual Workshop of the Psychology of Programming Interest Group (PPIG'07)*, pages 158–172, 2007.

[BT07b] Roman Bednarik and Markku Tukiainen. Validating the restricted focus viewer: A study using eye-movement tracking. *Behavior Research Methods* 39(2):274–282, 2007.

[CL13] Monchu Chen and Veraneka Lim. Eye gaze and mouse cursor relationship in a debugging task. In *HCI International 2013—Posters Extended Abstracts*, pages 468–472. Springer, 2013.

[CS90] Martha E. Crosby and Jan Stelovsky. How do we read algorithms? A case study. *Computer* 23(1):25–35, 1990.

[CSW02] Martha E. Crosby, Jean Scholtz, and Susan Wiedenbeck. The roles beacons play in comprehension for novice and expert programmers. In *Proceedings of the 14th Workshop of the Psychology of Programming Interest Group*, pages 58–73, 2002.

[CTKT13] Nergiz Ercil Cagiltay, Gul Tokdemir, Ozkan Kilic, and Damla Topalli. Performing and analyzing non-formal inspections of entity relationship diagram (ERD). *Journal of Systems and Software* 86(8):2184–2195, 2013.

[DB17] Sarah D’Angelo and Andrew Begel. Improving communication between pair programmers using shared gaze awareness. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*, pages 6245–6290. ACM, 2017.

[DcI13] Hacı Ali Duru, Murat Perit Çakır, and Veysi İşler. How does software visualization contribute to software comprehension? A grounded theory approach. *International Journal of Human-Computer Interaction* 29(11):743–763, 2013.

[DSLS + 14] Benoit De Smet, Lorent Lempereur, Zohreh Sharafi, Yann-Gaël Guéhéneuc, Giuliano Antoniol, and Naji Habra. Taupe: Visualizing and analyzing eyetracking data. *Science of Computer Programming* 79:260–278, 2014.

[Dub09] Gilles Dubochet. Computer code as a medium for human communication: Are programming languages improving? In *Proceedings of the 21st Working Conference on the Psychology of Programmers Interest Group*, number LAMP-CONF-2009-001, pages 174–187, 2009.

[FBM + 14] Thomas Fritz, Andrew Begel, Sebastian C. Müller, Serap Yigit-Elliott, and Manuela Züger. Using psycho-physiological measures to assess task difficulty in software development. In *Proceedings of the 36th International Conference on Software Engineering*, pages 402–413. ACM, 2014.

[Gu06] Yann-Gaël Guéhéneuc. Taupe: Towards understanding program comprehension. In *Proceedings of the 2006 Conference of the Center for Advanced Studies on Collaborative Research*, page 1. IBM Corp., 2006.

[GWMP16] Anurag Goswami, Gursimran Walia, Mark McCourt, and Ganesh Padmanabhan. Using eye tracking to investigate reading patterns and learning styles of software requirement inspectors to enhance inspection team outcome. In *Proceedings of the 10th ACM/IEEE International Symposium on Empirical Software Engineering and Measurement*, page 34. ACM, 2016.

[HLL + 13] T.-Y. Hou, Y.-T. Lin, Y.-C. Lin, C.-H. Chang, and M.-H. Yen. Exploring the gender effect on cognitive processes in program debugging based on eye movement analysis. In *Proceedings of the 5th International Conference on Computer Supported Education (CSEDU'13)*, pages 469–473, 2013.

[HN12] Prateek Hejmady and N. Hari Narayanan. Visual attention patterns during program debugging with an IDE. In *Proceedings of the Symposium on Eye Tracking Research and Applications*, pages 197–200. ACM, 2012.

[JF15] Ahmad Jbara and Dror G. Feitelson. How programmers read regular code: A controlled experiment using eye-tracking. In *Proceedings of the 2015 IEEE 23rd International Conference on Program Comprehension*, pages 244–254. IEEE Press, 2015.

[JGSH09] Sebastien Jeanmart, Yann-Gaël Guéhéneuc, Houari Sahraoui, and Naji Habra. Impact of the visitor pattern on program comprehension and maintenance. In *Proceedings of the 2009 3rd International Symposium on Empirical Software Engineering and Measurement*, pages 69–78. IEEE Computer Society, 2009.

[JN12] Patrick Jermann and Marc-Antoine Nussli. Effects of sharing text selections on gaze cross-recurrence and interaction quality in a pair programming task. In *Proceedings of the ACM 2012 Conference on Computer Supported Cooperative Work*, pages 1125–1134. ACM, 2012.

[Loh14] Sebastian Lohmeier. Activation and the comprehension of indirect anaphors in source code. In *Proceedings of the 25th Annual Workshop of the Psychology of Programming Interest Group*, 2007, 2014.

[LWH + 16] Yu-tzu Lin, Cheng-chih Wu, Ting-yun Hou, Yu-chih Lin, Fang-ying Yang, and Chia-hu Chang. Tracking students cognitive processes during program debugging an eye-movement approach. *IEEE Transactions on Education* 59(3):175–186, 2016.

[MD15] Sayani Mondal and Partha Pratim Das. An IDE-agnostic system to capture reading behavior of C++ programs using eyegaze tracker. In *Proceedings of the 2015 5th National Conference on Computer Vision, Pattern Recognition, Image Processing and Graphics (NCVPRIPG)*, pages 1–4. IEEE, 2015.

[MDPVRDVI16] Ana I. Molina-Díaz, Maximiliano Paredes-Velasco, Miguel A. Redondo-Duque, and Jesús Ángel Velázquez-Iturbide. Evaluation experiences of the representation techniques of greedy programs: Application to the GreedEx tool. *IEEE Revista Iberoamericana de Tecnologías del Aprendizaje* 11(3):179–186, 2016.

[MGRB15] Ana Isabel Molina, Jesus Gallardo, Miguel Angel Redondo, and Crescencio Bravo. Assessing the awareness mechanisms of a collaborative programming support system. *Dyna* 82(193):212–222, 2015.

[MNH + 17] Jean Melo, Fabricio Batista Narcizo, Dan Witzner Hansen, Claus Brabrand, and Andrzej Wasowski. Variability through the eyes of the programmer. In *Proceedings of the 25th International Conference on Program Comprehension*, pages 34–44. IEEE Press, 2017.

[NHMG16] Markus Nivala, Florian Hauser, Jürgen Mottok, and Hans Gruber. Developing visual expertise in software engineering: An eye tracking study. In *Proceedings of the 2016 IEEE Global Engineering Education Conference (EDUCON)*, pages 613–620. IEEE, 2016.

[NS04] Seppo Nevalainen and Jorma Sajaniemi. Comparison of three eye-tracking devices in psychology of programming research. In *Proceedings of the 6th Annual Psychology of Programming Interest Group*, pages 170–184, 2004.

[NS05] Seppo Nevalainen and Jorma Sajaniemi. Short-term effects of graphical versus textual visualisation of variables on program perception. In *Proceedings of the 17th Annual Psychology of Programming Interest Group Workshop*, pages 77–91, 2005.

[NS06] Seppo Nevalainen and Jorma Sajaniemi. An experiment on short-term effects of animated versus static visualization of operations on program perception. In *Proceedings of the 2nd International Workshop on Computing Education Research (ICER'06)*, pages 7–16. ACM, 2006.

[PBT09] Sami Pietinen, Roman Bednarik, and Markku Tukiainen. An exploration of shared visual attention in collaborative programming. In *Proceedings of the 21st Annual Psychology of Programming Interest Group Conference (PPIG)*, 2009.

[PG10] Gerardo Cepeda Porras and Yann-Gaël Guéhéneuc. An empirical study on the efficiency of different design pattern representations in UML class diagrams. *Empirical Software Engineering* 15(5):493–522, 2010.

[PLS + 16] Fei Peng, Chunyu Li, Xiaohan Song, Wei Hu, and Guihuan Feng. An eye tracking research on debugging strategies towards different types of bugs. In *Proceedings of the 2016 IEEE 40th Annual Computer Software and Applications Conference (COMPSAC)*, Vol. 2, pages 130–134. IEEE, 2016.

[PIS17] Patrick Peachock, Nicholas Iovino, and Bonita Sharif. 2017. Investigating Eye Movements in Natural Language and C++ Source Code—A Replication Experiment. In *Proceedings of the 11th International Conference on Augmented Cognition. Neurocognition and Machine Learning (AC'17)*, pages 206–218. Springer International Publishing, 2017.

[PSGJ17] Sofia Papavaslopoulou, Kshitij Sharma, Michail Giannakos, and Letizia Jaccheri. Using eye-tracking to unveil differences between kids and teens in coding activities. In *Proceedings of the 2017 Conference on Interaction Design and Children*, pages 171–181. ACM, 2017.

[RBCL03] Pablo Romero, Benedict Boulay, Richard Cox, and Rudi Lutz. Java debugging strategies in multi-representational environments. *Psychology of Programming Languages Interest Group*, pages 1–14, 2003.

[RCdBL02] Pablo Romero, Richard Cox, Benedict du Boulay, and Rudi Lutz. Visual attention and representation switching during Java program debugging: A study using the restricted focus viewer. In *Diagrammatic Representation and Inference*, pages 221–235. Springer, 2002.

[RLCdB02] P. Romero, R. Lutz, R. Cox, and B. du Boulay. Co-ordination of multiple external representations during Java program debugging. In *Proceedings of the IEEE 2002 Symposia on Human Centric Computing Languages and Environments*, pages 207–214, 2002.

[RLMM15] P. Rodeghero, Cheng Liu, P. W. McBurney, and C. McMillan. An eye-tracking study of Java programmers and application to source code summarization. *IEEE Transactions on Software Engineering*, 41(11):1038–1054, Nov. 2015.

[SB04] Randy Stein and Susan E. Brennan. Another persons eye gaze as a cue in solving programming problems. In *Proceedings of the 6th International Conference on Multimodal Interfaces (ICMI'04)*, pages 915, ACM, 2004.

[SFM12] Bonita Sharif, Michael Falcone, and Jonathan I. Maletic. An eye-tracking study on the role of scan time in finding source code defects. In *Proceedings of the Symposium on Eye Tracking Research and Applications (ETRA'12)*, pages 381–384, ACM, 2012.

[SJAP13] B. Sharif, G. Jetty, J. Aponte, and E. Parra. An empirical study assessing the effect of seeit 3D on comprehension. In *Proceedings of the 1st IEEE Working Conference on Software Visualization (VISOFT'13)*, pages 110, Sept. 2013.

[SM10a] B. Sharif and J. I. Maletic. An eye-tracking study on camel-case and under score identifier styles. In *Proceedings of the 2010 IEEE 18th International Conference on Program Comprehension (ICPC)*, pages 196–205, June 2010.

[SM10b] B. Sharif and J. I. Maletic. An eye-tracking study on the effects of layout in understanding the role of design patterns. In *Proceedings of the 2010 IEEE International Conference on Software Maintenance (ICSM)*, pages 110, Sept. 2010.

[SMS + 13] Z. Sharafi, A. Marchetto, A. Susi, G. Antoniol, and Y.-G. Guéhéneuc. An empirical study on the efficiency of graphical vs. textual representations in requirements comprehension. In *Proceedings of the 2013 IEEE 21st International Conference on Program Comprehension (ICPC)*, pages 33–42, May 2013.

[SSGA12] Z. Sharafi, Z. Soh, Y. Guéhéneuc, and G. Antoniol. Women and men—Different but equal: On the impact of identifier style on source code reading. In *Proceedings of the 2012 IEEE 20th International Conference on Program Comprehension (ICPC)*, pages 27–36, June 2012.

[SSVdP + 12] Z. Soh, Z. Sharafi, B. Van den Plas, G. C. Porras, Y. Gueheneuc, and G. Antoniol. Professional status and expertise for UML class diagram comprehension: An empirical study. In *Proceedings of the , 2012 IEEE 20th International Conference on Program Comprehension (ICPC)*, pages 163–172, June 2012.

[SUGGR15] M. Sami Uddin, V. Gaur, C. Gutwin, and C. K. Roy. On the comprehension of code clone visualizations: A controlled study using eye-tracking. In *Proceedings of the 2015 IEEE 15th International Working Conference on Source Code Analysis and Manipulation (SCAM)*, pages 161–170, Sept 2015.

[TFSL14] Rachel Turner, Michael Falcone, Bonita Sharif, and Alina Lazar. An eyetracking study assessing the comprehension of C++ and Python source code. In *Proceedings of the Symposium on Eye Tracking Research and Applications (ETRA'14)*, pages 231–234, ACM, 2014.

[UNMM06] Hidetake Uwano, Masahide Nakamura, Akito Monden, and Ken-ichi Matsumoto. Analyzing individual performance of source code review using reviewers eye movement. In *Proceedings of the 2006 Symposium on Eye Tracking Research and Applications (ETRA'06)*, pages 133–140, ACM, 2006.

[WSSK14] Braden Walters, Timothy Shaffer, Bonita Sharif, and Huzefa Kagdi. Capturing software traceability links from developers eye gazes. In *Proceedings of the 22nd International Conference on Program Comprehension (ICPC'14)*, pages 201–204, ACM, 2014.

[YKM07] S. Yusuf, H. Kagdi, and J. I. Maletic. Assessing the comprehension of UML class diagrams via eye-tracking. In *Proceedings of the 15th IEEE International Conference on Program Comprehension (ICPC'07)*, pages 113–122, June 2007.

ACKNOWLEDGMENTS

The authors would like to thank the reviewers and associate editor for their comments which improved this manuscript.

REFERENCES

- [1] Nasir Ali, Zohreh Sharafi, Yann-Gael Gueheneuc, and Giuliano Antoniol. 2015. An empirical study on the importance of source code entities for requirements traceability. *Empirical Software Engineering* 20, 2 (2015), 442–478.
- [2] Lorin W. Anderson, David R. Krathwohl, and Benjamin Samuel Bloom. 2001. *A Taxonomy for Learning, Teaching, and Assessing: A Revision of Bloom's Taxonomy of Educational Objectives*. Allyn & Bacon.
- [3] Magdalena Andrzejewska, Anna Stolińska, Władysław Błasiak, Paweł Pkęczkowski, Roman Rosiek, Bożena Rożek, Mirosława Sajka, and Dariusz Wcisło. 2016. Eye-tracking verification of the strategy used to analyse algorithms expressed in a flowchart and pseudocode. *Interactive Learning Environments* 24, 8 (2016), 1981–1995.
- [4] Stuart Charters and Barbara Kitchenham. 2007. *Guidelines for Performing Systematic Literature Reviews in Software Engineering*. Technical Report. Keele University and Durham University Joint Report.
- [5] Roman Bednarik, Teresa Busjahn, and Carsten Schulte (Eds.). 2014. *Eye Movements in Programming Education: Analyzing the Experts Gaze*. Technical Report. University of Eastern Finland, Joensuu, Finland.
- [6] Roman Bednarik, Niko Myller, Erkki Sutinen, and Markku Tukiainen. 2005. Applying eye-movement tracking to program visualization. In *Proceedings of the 2005 IEEE Symposium on Visual Languages and Human-Centric Computing*. IEEE, 302–304.
- [7] Roman Bednarik, Niko Myller, Erkki Sutinen, and Markku Tukiainen. 2005. Effects of experience on gaze behaviour during program animation. In *Proceedings of the 17th Annual Psychology of Programming Interest Group Workshop*. 49–61.
- [8] Roman Bednarik, Niko Myller, Erkki Sutinen, and Markku Tukiainen. 2006. Analyzing individual differences in program comprehension. *Technology Instruction Cognition and Learning* 3, 3/4 (2006), 205.
- [9] Roman Bednarik, Niko Myller, Erkki Sutinen, and Markku Tukiainen. 2006. Program visualization: Comparing eye-tracking patterns with comprehension summaries and performance. In *Proceedings of the 18th Annual Psychology of Programming Workshop*. 66–82.
- [10] Roman Bednarik and Markku Tukiainen. 2004. Visual attention and representation switching in Java program debugging: A study using eye movement tracking. In *Proceedings of the 16th Annual Workshop of the Psychology of Programming Interest Group*. 159–169.
- [11] Roman Bednarik and Markku Tukiainen. 2004. Visual attention tracking during program debugging. In *Proceedings of the 3rd Nordic Conference on Human-Computer Interaction*. ACM, 331–334.
- [12] Roman Bednarik and Markku Tukiainen. 2005. Effects of display blurring on the behavior of novices and experts during program debugging. In *CHI'05 Extended Abstracts on Human Factors in Computing Systems*. ACM, 1204–1207.
- [13] Roman Bednarik and Markku Tukiainen. 2006. An eye-tracking methodology for characterizing program comprehension processes. In *Proceedings of the 2006 Symposium on Eye Tracking Research and Applications*. ACM, 125–132.

- [14] Roman Bednarik and Markku Tukiainen. 2007. Analysing and interpreting quantitative eye-tracking data in studies of programming: Phases of debugging with multiple representations. In *Proceedings of the 19th Annual Workshop of the Psychology of Programming Interest Group (PPIG'07)*. 158–172.
- [15] Roman Bednarik and Markku Tukiainen. 2007. Validating the restricted focus viewer: A study using eye-movement tracking. *Behavior Research Methods* 39, 2 (2007), 274–282.
- [16] Roman Bednarik and Markku Tukiainen. 2008. Temporal eye-tracking data: Evolution of debugging strategies with multiple representations. In *Proceedings of the 2008 Symposium on Eye Tracking Research and Applications*. ACM, 99–102.
- [17] Mordechai Ben-Ari, Roman Bednarik, Ronit Ben-Bassat Levy, Gil Ebel, Andrés Moreno, Niko Myller, and Erkki Sutinen. 2011. A decade of research and development on program animation: The Jeliot experience. *Journal of Visual Languages and Computing* 22, 5 (2011), 375–384.
- [18] Alan F. Blackwell, Anthony R. Jansen, and Kim Marriott. 2000. Restricted focus viewer: A tool for tracking visual attention. *Theory and Application of Diagrams*. Springer, 162–177.
- [19] Crescencio Bravo, Rafael Duque, and Jess Gallardo. 2013. A groupware system to support collaborative programming: Design and experiences. *Journal of Systems and Software* 86, 7 (2013), 1759–1771.
- [20] Teresa Busjahn, Roman Bednarik, Andrew Begel, Martha Crosby, James H. Paterson, Carsten Schulte, Bonita Sharif, and Sascha Tamm. 2015. Eye movements in code reading: Relaxing the linear order. In *Proceedings of the 2015 IEEE 23rd International Conference on Program Comprehension*. IEEE Press, 255–265.
- [21] Teresa Busjahn, Roman Bednarik, and Carsten Schulte. 2014. What influences dwell time during source code reading?: Analysis of element type and frequency as factors. In *Proceedings of the Symposium on Eye Tracking Research and Applications*. ACM, 335–338.
- [22] Teresa Busjahn, Carsten Schulte, and Andreas Busjahn. 2011. Analysis of code reading to gain more insight in program comprehension. In *Proceedings of the 11th Koli Calling International Conference on Computing Education Research*. ACM, 1–9.
- [23] Teresa Busjahn, Carsten Schulte, Bonita Sharif, Andrew Begel, Michael Hansen, Roman Bednarik, Paul Orlov, Petri Ihantola, Galina Shchekotova, Maria Antropova, and others. 2014. Eye tracking in computing education. In *Proceedings of the 10th Annual Conference on International Computing Education Research*. ACM, 3–10.
- [24] Teresa Busjahn, Carsten Schulte, Sascha Tamm, and Roman Bednarik (Eds.). 2015. *Eye Movements in Programming Education II: Analyzing the Novice's Gaze*. Technical Report. Freie Universität Berlin, Department of Mathematics and Computer Science, Berlin, Germany. 1–41 pages.
- [25] Zoya Bylinskii, Michelle A. Borkin, Nam Wook Kim, Hanspeter Pfister, and Aude Oliva. 2015. Eye fixation metrics for large scale evaluation and comparison of information visualizations. In *Workshop on Eye Tracking and Visualization*. Springer, 235–255.
- [26] Martha E. Crosby and Jan Stelovsky. 1990. How do we read algorithms? A case study. *Computer* 23, 1 (1990), 25–35.
- [27] Sarah D'Angelo and Andrew Begel. 2017. Improving communication between pair programmers using shared gaze awareness. In *Proceedings of the 2017 CHI Conference on Human Factors in Computing Systems*. ACM, 6245–6290.
- [28] Benoit De Smet, Lorent Lempereur, Zohreh Sharafi, Yann-Gaël Guéhéneuc, Giuliano Antoniol, and Naji Habra. 2014. Taupe: Visualizing and analyzing eye-tracking data. *Science of Computer Programming* 79 (2014), 260–278.
- [29] Fadi P. Deek and James A. McHugh. 1998. A survey and critical analysis of tools for learning programming. *Computer Science Education* 8, 2 (1998), 130–178.
- [30] Tore Dyba, Torgeir Dingsoyr, and Geir K. Hanssen. 2007. Applying systematic reviews to diverse study types: An experience report. In *Proceedings of the 1st International Symposium on Empirical Software Engineering and Measurement (ESEM'07)*. IEEE, 225–234.
- [31] Joseph H. Goldberg and Xerxes P. Kotval. 1999. Computer interface evaluation using eye movements: Methods and constructs. *International Journal of Industrial Ergonomics* 24, 6 (1999), 631–645.
- [32] Anabela Gomes and António José Mendes. 2007. An environment to improve programming education. In *Proceedings of the 2007 International Conference on Computer Systems and Technologies*. ACM, 88.
- [33] Anabela Gomes and António José Mendes. 2007. Learning to program—Difficulties and solutions. In *Proceedings of the International Conference on Engineering Education (ICEE)*, Vol. 2007.
- [34] Yann-Gael Guéhéneuc, Huzefa Kagdi, and Jonathan I. Maletic. 2009. Working session: Using eye-tracking to understand program comprehension. In *Proceedings of the IEEE 17th International Conference on Program Comprehension (ICPC'09)*. IEEE, 278–279.
- [35] Victor Henning and Jan Reichelt. 2008. Mendeley—A last.fm for research?. In *Proceedings of the IEEE 4th International Conference on eScience (eScience'0)*. IEEE, 327–328.
- [36] Ahmad Jbara and Dror G. Feitelson. 2015. How programmers read regular code: A controlled experiment using eye tracking. In *Proceedings of the 2015 IEEE 23rd International Conference on Program Comprehension*. IEEE, 244–254.

- [37] Tomoko Kashima, Shimpei Matsumoto, and Shuichi Yamagishi. 2014. Proposal of a method to measure difficulty level of programming code with eye-tracking. In *Human-Computer Interaction. Advanced Interaction Modalities and Techniques*. Springer, 264–272.
- [38] Barbara Kitchenham and Pearl Brereton. 2013. A systematic review of systematic review process research in software engineering. *Information and Software Technology* 55, 12 (2013), 2049–2075.
- [39] Barbara A. Kitchenham, Shari Lawrence Pfleeger, Lesley M. Pickard, Peter W. Jones, David C. Hoaglin, Khaled El Emam, and Jarrett Rosenberg. 2002. Preliminary guidelines for empirical research in software engineering. *IEEE Transactions on Software Engineering* 28, 8 (2002), 721–734.
- [40] Tomasz Kocejko, Jacek Ruminski, Adam Bujnowski, and Jerzy Wtorek. 2016. The evaluation of eGlasses eye tracking module as an extension for Scratch. In *Proceedings of the 2016 9th International Conference on Human System Interactions (HSI)*. IEEE, 465–471.
- [41] Martin Konopka. 2015. Combining eye tracking with navigation paths for identification of cross-language code dependencies. In *Proceedings of the 2015 10th Joint Meeting on Foundations of Software Engineering*. ACM, 1057–1059.
- [42] Meng-Lung Lai, Meng-Jung Tsai, Fang-Ying Yang, Chung-Yuan Hsu, Tzu-Chien Liu, Silvia Wen-Yu Lee, Min-Hsien Lee, Guo-Li Chiou, Jyh-Chong Liang, and Chin-Chung Tsai. 2013. A review of using eye-tracking technology in exploring learning from 2000 to 2012. *Educational Research Review* 10 (2013), 90–115.
- [43] Iain Milne and Glenn Rowe. 2002. Difficulties in learning and teaching programming views of students and tutors. *Education and Information Technologies* 7, 1 (2002), 55–66.
- [44] David Montano, Jairo Aponte, and Andrian Marcus. 2009. Sv3D meets eclipse. In *Proceedings of the 5th IEEE International Workshop on Visualizing Software for Understanding and Analysis (VISSOFT'09)*. IEEE, 51–54.
- [45] IT Chan Mow. 2008. Issues and difficulties in teaching novice computer programming. In *Innovative Techniques in Instruction Technology, e-learning, e-assessment, and Education*. Springer, 199–204.
- [46] Gail C. Murphy, Mik Kersten, and Leah Findlater. 2006. How are Java software developers using the Elipse IDE? *IEEE Software* 23, 4 (2006), 76–83.
- [47] Seppo Nevalainen and Jorma Sajaniemi. 2004. Comparison of three eye tracking devices in psychology of programming research. In *Proceedings of the 16th Annual Psychology of Programming Interest Group Workshop*. 151–158.
- [48] Christopher Palmer and Bonita Sharif. 2016. Towards automating fixation correction for source code. In *Proceedings of the 9th Biennial ACM Symposium on Eye Tracking Research and Applications*. ACM, 65–68.
- [49] Oren Patashnik. 1988. BibTeXing. documentation for general BibTeX users. *Electronic Document Accompanying BibTeX Distribution*.
- [50] Patrick Peachock, Nicholas Iovino, and Bonita Sharif. 2017. *Investigating Eye Movements in Natural Language and C++ Source Code—A Replication Experiment*. Springer International Publishing, Cham, 206–218. DOI : http://dx.doi.org/10.1007/978-3-319-58628-1_17
- [51] Kai Petersen, Robert Feldt, Shahid Mujtaba, and Michael Mattsson. 2008. Systematic mapping studies in software engineering. In *Proceedings of the 12th International Conference on Evaluation and Assessment in Software Engineering (EASE'08)*, Vol. 8. 68–77.
- [52] Kai Petersen, Sairam Vakkalanka, and Ludwik Kuzniarz. 2015. Guidelines for conducting systematic mapping studies in software engineering: An update. *Information and Software Technology* 64 (2015), 1–18.
- [53] Razvan Petrusel and Jan Mendling. 2013. Eye-tracking the factors of process model comprehension tasks. In *Advanced Information Systems Engineering*. Springer, 224–239.
- [54] Sami Pietinen, Roman Bednarik, Tatiana Glotova, Vesa Tenhunen, and Markku Tukiainen. 2008. A method to study visual attention aspects of collaboration: Eye-tracking pair programmers simultaneously. In *Proceedings of the 2008 Symposium on Eye Tracking Research and Applications (ETRA'08)*. ACM, New York, 39–42. DOI : <http://dx.doi.org/10.1145/1344471.1344480>
- [55] Sami Pietinen, Roman Bednarik, and Markku Tukiainen. 2010. Shared visual attention in collaborative programming: A descriptive analysis. In *Proceedings of the 2010 ICSE Workshop on Cooperative and Human Aspects of Software Engineering (CHASE'10)*. ACM, New York, 21–24. DOI : <http://dx.doi.org/10.1145/1833310.1833314>
- [56] Gerardo Cepeda Porras and Yann-Gaël Guéhéneuc. 2010. An empirical study on the efficiency of different design pattern representations in UML class diagrams. *Empirical Software Engineering* 15, 5 (2010), 493–522.
- [57] Keith Rayner. 1998. Eye movements in reading and information processing: 20 years of research. *Psychological Bulletin* 124, 3 (1998), 372.
- [58] Keith Rayner. 2009. Eye movements and attention in reading, scene perception, and visual search. *The Quarterly Journal of Experimental Psychology* 62, 8 (2009), 1457–1506.
- [59] V. Renumol, S. Jayaprakash, and D. Janakiram. 2009. Classification of cognitive difficulties of students to learn computer programming. *Indian Institute of Technology, India*.

- [60] Mitchel Resnick, John Maloney, Andrés Monroy-Hernández, Natalie Rusk, Evelyn Eastmond, Karen Brennan, Amon Millner, Eric Rosenbaum, Jay Silver, Brian Silverman, and others. 2009. Scratch: Programming for all. *Communications of the ACM* 52, 11 (2009), 60–67.
- [61] Paige Rodeghero and Collin McMillan. 2015. An empirical study on the patterns of eye movement during summarization tasks. In *Proceedings of the 2015 ACM/IEEE International Symposium on Empirical Software Engineering and Measurement (ESEM'15)*. 1–10. DOI: <http://dx.doi.org/10.1109/ESEM.2015.7321188>
- [62] Paige Rodeghero, Collin McMillan, Paul W. McBurney, Nigel Bosch, and Sidney D'Mello. 2014. Improving automated source code summarization via an eye-tracking study of programmers. In *Proceedings of the 36th International Conference on Software Engineering (ICSE'14)*. ACM, New York, 390–401. DOI: <http://dx.doi.org/10.1145/2568225.2568247>
- [63] Pablo Romero, Benedict Boulay, Richard Cox, and Rudi Lutz. 2003. Java debugging strategies in multi-representational environments. In *Proceedings of the 15th Annual Psychology of Programming Interest Group Workshop*. 421–435.
- [64] Md. Sami Uddin, Varun Gaur, Carl Gutwin, and Chanchal K. Roy. 2015. On the comprehension of code clone visualizations: A controlled study using eye tracking. In *Proceedings of the 2015 IEEE 15th International Working Conference on Source Code Analysis and Manipulation (SCAM'15)*. 161–170. DOI: <http://dx.doi.org/10.1109/SCAM.2015.7335412>
- [65] Zohreh Sharafi. 2011. A systematic analysis of software architecture visualization techniques. In *Proceedings of the 2011 IEEE 19th International Conference on Program Comprehension (ICPC'11)*. 254–257. DOI: <http://dx.doi.org/10.1109/ICPC.2011.40>
- [66] Zohreh Sharafi, A. Marchetto, A. Susi, G. Antoniol, and Yann-Gaël Guéhéneuc. 2013. An empirical study on the efficiency of graphical vs. textual representations in requirements comprehension. In *Proceedings of the 2013 IEEE 21st International Conference on Program Comprehension (ICPC'13)*. 33–42. DOI: <http://dx.doi.org/10.1109/ICPC.2013.6613831>
- [67] Zohreh Sharafi, Zéphyrin Soh, and Yann-Gaël Guéhéneuc. 2015. A systematic literature review on the usage of eye-tracking in software engineering. *Information and Software Technology* 67 (2015), 79–107.
- [68] Zohreh Sharafi, Zéphyrin Soh, Yann-Gaël Guéhéneuc, and G. Antoniol. 2012. Women and men—Different but equal: On the impact of identifier style on source code reading. In *Proceedings of the 2012 IEEE 20th International Conference on Program Comprehension (ICPC'12)*. 27–36. DOI: <http://dx.doi.org/10.1109/ICPC.2012.6240505>
- [69] Bonita Sharif, Michael Falcone, and Jonathan I. Maletic. 2012. An eye-tracking study on the role of scan time in finding source code defects. In *Proceedings of the Symposium on Eye Tracking Research and Applications (ETRA'12)*. ACM, New York, 381–384. DOI: <http://dx.doi.org/10.1145/2168556.2168642>
- [70] Bonita Sharif, G. Jetty, J. Aponte, and E. Parra. 2013. An empirical study assessing the effect of seeing 3D on comprehension. In *Proceedings of the 2013 1st IEEE Working Conference on Software Visualization (VISSOFT'13)*. 1–10. DOI: <http://dx.doi.org/10.1109/VISSOFT.2013.6650519>
- [71] Bonita Sharif and Huzefa Kagdi. 2011. On the use of eye tracking in software traceability. In *Proceedings of the 6th International Workshop on Traceability in Emerging Forms of Software Engineering (TEFSE'11)*. ACM, New York, 67–70. DOI: <http://dx.doi.org/10.1145/1987856.1987872>
- [72] Bonita Sharif and Jonathan I. Maletic. 2010. An eye tracking study on camelcase and under_score identifier styles. In *Proceedings of the 2010 IEEE 18th International Conference on Program Comprehension (ICPC'10)*. 196–205. DOI: <http://dx.doi.org/10.1109/ICPC.2010.41>
- [73] Bonita Sharif and Jonathan I. Maletic. 2010. An eye tracking study on the effects of layout in understanding the role of design patterns. In *Proceedings of the 2010 IEEE International Conference on Software Maintenance (ICSM'10)*. 1–10. DOI: <http://dx.doi.org/10.1109/ICSM.2010.5609582>
- [74] Kshitij Sharma, Patrick Jermann, Marc-Antoine Nüssli, and Pierre Dillenbourg. 2013. Understanding collaborative program comprehension: Interlacing gaze and dialogues. In *Proceedings of Computer Supported Collaborative Learning (CSCL'13)*.
- [75] Judy Sheard, S. Simon, Margaret Hamilton, and Jan Lönnberg. 2009. Analysis of research into the teaching and learning of programming. In *Proceedings of the 5th International Workshop on Computing Education Research Workshop*. ACM, 93–104.
- [76] Dag I. K. Sjøberg, Jo Erskine Hannay, Ove Hansen, Vigdis By Kampenes, Amela Karahasanovic, N.-K. Liborg, and Anette C. Rekdal. 2005. A survey of controlled experiments in software engineering. *IEEE Transactions on Software Engineering* 31, 9 (2005), 733–753.
- [77] Zéphyrin Soh. 2011. Context and vision: Studying two factors impacting program comprehension. In *Proceedings of the 2011 IEEE 19th International Conference on Program Comprehension (ICPC'11)*. 258–261. DOI: <http://dx.doi.org/10.1109/ICPC.2011.37>
- [78] Zéphyrin Soh, Zohreh Sharafi, Bertrand Van den Plas, Gerardo Cepeda Porras, Yann-Gaël Guéhéneuc, and Giuliano Antoniol. 2012. Professional status and expertise for UML class diagram comprehension: An empirical study.

- In *Proceedings of the 2012 IEEE 20th International Conference on Program Comprehension (ICPC'12)*. 163–172. DOI : <http://dx.doi.org/10.1109/ICPC.2012.6240484>
- [79] Koji Torii, Ken-ichi Matsumoto, Kumiyo Nakakoji, Yoshihiro Takada, Shingo Takada, and Kazuyuki Shima. 1999. Ginger2: An environment for computer-aided empirical software engineering. *IEEE Transactions on Software Engineering* 25, 4 (July 1999), 474–492. DOI : <http://dx.doi.org/10.1109/32.799942>
 - [80] Rachel Turner, Michael Falcone, Bonita Sharif, and Alina Lazar. 2014. An eye-tracking study assessing the comprehension of C++ and python source code. In *Proceedings of the Symposium on Eye Tracking Research and Applications (ETRA'14)*. ACM, New York, 231–234. DOI : <http://dx.doi.org/10.1145/2578153.2578218>
 - [81] Hidetake Uwano, Akito Monden, and Ken-ichi Matsumoto. 2008. DRESREM 2: An analysis system for multi-document software review using reviewers' eye movements. In *Proceedings of the 3rd International Conference on Software Engineering Advances (ICSEA'08)*. 177–183. DOI : <http://dx.doi.org/10.1109/ICSEA.2008.49>
 - [82] Hidetake Uwano, Masahide Nakamura, Akito Monden, and Ken-ichi Matsumoto. 2007. Exploiting eye movements for evaluating reviewer's performance in software review. *IEICE Transactions on Fundamentals of Electronics, Communications and Computer Sciences* 90, 10 (2007), 2290–2300.
 - [83] Johannes Zagermann, Ulrike Pfeil, and Harald Reiterer. 2016. Measuring cognitive load using eye tracking technology in visual computing. In *Proceedings of the 6th Workshop on Beyond Time and Errors on Novel Evaluation Methods for Visualization (BELIV'16)*. 78–85.

Received August 2016; revised July 2017; accepted September 2017