



# Systematic literature review of empirical studies on mental representations of programs<sup>☆</sup>

Leah Bidlake<sup>a,\*</sup>, Eric Aubanel<sup>a</sup>, Daniel Voyer<sup>b</sup>

<sup>a</sup> Faculty of Computer Science, University of New Brunswick, 550 Windsor St, Fredericton, NB, E3B 5A3, Canada

<sup>b</sup> Department of Psychology, University of New Brunswick, Canada

## ARTICLE INFO

### Article history:

Received 22 July 2019

Revised 26 February 2020

Accepted 27 February 2020

Available online 29 February 2020

### Keywords:

Mental representations

Program comprehension

Systematic literature review

## ABSTRACT

Programmers are frequently tasked with modifying, enhancing, and extending applications. To perform these tasks, programmers must understand existing code by forming mental representations. Empirical research is required to determine the mental representations constructed during program comprehension to inform the development of programming languages, instructional practices, and tools. To make recommendations for future work a systematic literature review was conducted that summarizes the empirical research on mental representations formed during program comprehension, how the methods and tasks have changed over time, and the research contributions.

The data items included in the systematic review are empirical studies of programmers that investigated the comprehension and internal representation of code written in a formal programming language. The eligibility criteria used in the review were meant to extract studies with a focus on knowledge representation as opposed to knowledge utilization.

The results revealed a lack of incremental research and a dramatic decline in the research meaning that newly developed or popularized languages and paradigms have not been a part of the research reviewed. Accordingly, we argue that there needs to be a resurgence of empirical research on the psychology of programming to inform the design of tools and languages, especially in new and emerging paradigms.

© 2020 Elsevier Inc. All rights reserved.

## 1. Introduction

Programmers are frequently tasked with modifying, enhancing, and extending applications. To perform these tasks, programmers must first understand the existing code. During the comprehension process, programmers form mental representations of the code they are working with (Détienne, 2001). Empirical research is required to determine how these mental representations are created and the form that they take as a starting point to developing programming languages and tools that fit with the underlying representations. Such an approach would assist programmers in the comprehension process and promote the formation of accurate mental representations. In the study conducted by Wiedenbeck et al. (1993) experts were able to build mental representations with more abstract characteristics than non-experts to support comprehension-related programming tasks

such as maintenance and debugging. The findings reported by Wiedenbeck et al. (1993) are compatible with the notion that the ability of programmers to develop accurate mental representations of code is important for supporting programming tasks critical to the development of high-quality software (see also Petre and Blackwell, 1999).

The focus on program comprehension is of particular interest as many of the tasks performed by programmers, such as maintaining and modifying existing applications, require the understanding of code (Corritore and Wiedenbeck, 1999). Empirical research is required to develop an understanding of the comprehension strategies used by programmers. With the growing number of programming languages and tools under development with the goals of increasing productivity and ease of learning and use, there is a great need for this empirical research.

Before conducting more empirical research on mental representations, we need to know the current state of affairs in this area. Knowledge of past empirical research can assist in developing methods that can more accurately target and capture the mental representations under study. To make recommendations for future work on mental representations we must have an understanding of prior research methods and findings. By discovering

<sup>☆</sup> This research did not receive any specific grant from funding agencies in the public, commercial, or not-for-profit sectors.

\* Corresponding author.

E-mail addresses: [leah.bidlake@unb.ca](mailto:leah.bidlake@unb.ca) (L. Bidlake), [aubanel@unb.ca](mailto:aubanel@unb.ca) (E. Aubanel), [voyer@unb.ca](mailto:voyer@unb.ca) (D. Voyer).

what research has already been done, we can also learn what gaps remain.

Accordingly, the purpose of the literature review conducted here is to provide an analysis of the research that has been done to advance the body of knowledge related to the psychology of programming. Our review provides a comprehensive overview of the empirical studies performed to date that have contributed to the understanding of the comprehension strategies used by programmers to understand code, and the mental representations created in this process. The work presented here will assist researchers in the selection of appropriate tasks and methods for future work. This review will also assist researchers who want to perform incremental research by building on the work that has already been done in the domain. Empirical studies related to problem solving, program design and development, debugging skills, programmer performance, and learning programming languages that have no bearing on mental representations are not included in this review. The research questions addressed by our review are as follows:

1. To date, what empirical research has examined mental representations during program comprehension?
2. How have the methods used in the studies changed over time?
3. How have the tasks used in the studies to stimulate the comprehension process changed over time?
4. What are the contributions that have resulted from empirical research on program comprehension?

## 2. Background

Research in program comprehension encompasses both the study of the cognitive processes used by programmers to understand code and how programming languages and tools support these cognitive processes (Storey, 2006). The cognitive component of program comprehension that is of interest here is the abstract mental representations that are formed during program comprehension. These mental representations, often referred to as mental models, are founded in the theories of text comprehension (Pennington, 1987a).

There have been a variety of differing approaches to mental models including how they are defined and inferred (Cañas and Antolí, 1998). Cañas and Antolí suggested that the reason for this is because researchers from different disciplines study mental models using different tasks and are sometimes interested in different aspects of the representations. The unifying definition for a mental model, proposed by Cañas and Antolí (1998), is a dynamic representation formed in working memory as a result of using knowledge from long term memory and the environment. Mental model representations are described by Cañas and Antolí (1998) as both a process and the result of a simulation process elicited by a task.

The mental model approach to program comprehension is based on the propositional or text-based model and the situation model that were first developed to describe text comprehension (Détienne, 2001). The program model, formed by programmers when applying structural knowledge to the code resulting in a surface level representation, corresponds to the propositional or text-based model formed during text comprehension. The situation model, developed to describe the abstract representation of text, corresponds to the domain model formed during program comprehension when domain knowledge is used to form an understanding of the real world situation represented by the program. The mental model approach to program comprehension involves the construction of both the program model and the domain model (Détienne, 2001).

Empirical research on program comprehension has a direct impact on the development of programming languages and tools. Boshernitsan et al. (2007) developed iXj, a tool that uses a vi-

sual language to allow programmers to specify and execute code changes. The design of iXj was guided by the Cognitive Dimensions framework developed by Green (1989) to provide programmers with visual representations that reflect their own mental model of the source code. Tubaishat (2001) developed a theoretical model, Conceptual Model for Software Fault Localization (CMSFL), from empirical research on programming knowledge and plans. The CMSFL model was then used as the basis for developing the BUG-DOCTOR, an Automated Assistant Fault Localization (AASFL) tool that assists programmers with software fault localization. Arab (1992) developed a tool for formatting and documenting Pascal programs to assist programmers to write more readable and easier to understand programs. The development of this tool was influenced by empirical research that identified formatting and documenting as important factors in program comprehension.

Other common approaches to empirical research on programming languages and tools have been usability and comparative studies. Usability studies examine the human use of programming languages, whereas comparative studies compare the features of different programming languages. The usability studies reviewed by Hornbæk (2006) employed a variety of measures and definitions of usability, and often compared previous or competing versions. Comparative studies, including those performed by Prechelt (2000) and Nanz and Furia (2015), analyzed performance of programs written in programming languages from procedural, functional, and scripting paradigms, and also compared languages within these paradigms. Performance in these studies was measured using quantitative analysis of the program features such as lines of code, runtime, memory usage, and reliability. Programmer effort was also analyzed by Prechelt (2000) who measured the time required to write a program. Comparative studies on parallel programming languages have taken a similar approach in comparing solutions to problem sets or algorithms written in different parallel programming languages (Feo, 2015; Chamberlain et al., 2000). The study conducted by Feo (2015) also compared qualitative measures of how the programmers who wrote the solutions felt about how easy or difficult the program was to write, and what they felt was good or bad about the programming language.

The usability and comparative studies mentioned here are comparing programming languages that are already widely used. Studying a language when it is already widely used offers limited advances unless one is willing to modify the language as a result of the empirical observations. These studies did not take into consideration the user and their ability to understand and provide maintenance for programs written in different programming languages. Comparative studies are also limited in that they only give the user choice between existing languages; they do not direct the development of new languages that may require a complete departure from the current approaches to provide the user with languages and tools that align with their cognitive processes. The shift from comparative studies to research on program comprehension, especially in parallel programming where there is a significant lack of theory (Mattson and Wrinn, 2008), is necessary to inform the development of programming languages, instructional practices, and tools.

## 3. Method

To answer the research questions posed in this study, we conducted a systematic literature review following the guidelines provided by Kitchenham and Charters (2007) and in accordance with PRISMA guidelines (Moher et al., 2009). Systematic reviews are performed to collect literature relevant to answer specific research questions by using a well-defined and documented protocol in order to provide a fair evaluation and enable study audit and

reproducibility (Kitchenham and Charters, 2007). In our systematic literature review we adopted as a search strategy automated search and backward snowballing technique. The automated search strategy included a search of databases from each of the relevant fields: psychology and computer science. All available publication types (e.g., theses, journal articles, conference papers, books) were included in the systematic review because the research fields involved in this study use some publication types more frequently than others. For example, computer science publications are more often conference papers whereas psychology publications tend to be journal articles. The inclusion of theses makes it more likely that unpublished work will be considered. To understand how research on program comprehension has developed over time, the review was not limited to a particular time period. The abstracts or full texts of documents extracted by the database searches were screened using specific inclusion and exclusion criteria (i.e., eligibility criteria). The eligibility criteria were used to ensure that the documents were relevant to the research questions presented earlier, which are centred around empirical studies on program comprehension. The automated search was followed by backward snowballing, a technique that uses the reference list of data items to identify additional data items (Wohlin, 2014).

The systematic review also included an analysis using the full text of empirical studies that met the eligibility criteria. The analysis consisted of identifying the year of the study and the tasks used to stimulate the comprehension process, and summarizing the methods used in the study and the findings. The summaries were then used to create categories that reflected the most prevalent methods and findings of the research included in our review. Categorizing the methods and findings assisted us in answering the research questions presented earlier by allowing us to make observations about how the research on mental representations of programs has changed over time and the resulting contributions.

### 3.1. Information sources

Five databases were included in the search process: Computer Source Index, ERIC, IEEE Xplore Digital Library, PsycINFO, and Scopus. Records in the Computer Source Index (formerly Computer Science Index) database are related to the current trends and advances in computer science. The ERIC database records are primarily related to the field of education, with publication dates from 1966 to present. The IEEE Xplore Digital Library database contains scientific and technical content published by the Institute of Electrical and Electronics Engineers (IEEE) and its publishing partners, from the fields of electrical engineering, computer science, and electronics. The publication dates range from 1872 to present. The records in the PsycINFO database refer to literature from the behavioural sciences and mental health fields of study, with publication dates ranging from 1887 to present. The Scopus database contains records from the sciences, health sciences, and social sciences, with publication dates ranging from 1966 to present.

The keyword, title, and abstract information were used to perform the database searches. The search string was developed using the main components of the search: mental representations, program comprehension, and programmers. Synonyms of these components were then identified from general knowledge acquired through of an initial investigation into this topic. All combinations of the synonyms were searched according to the requirements of the specific database. The combined words could be at most separated by two words and the wildcard symbol (\*) was used to include alternative spellings of words. The following search string was used for database keyword searches of Computer Source, ERIC, and PsycINFO:

**Program Comprehension:** (((code OR program\* OR software) N2 (understanding OR comprehen\* OR stud\* OR analy\* OR maint\*

OR modif\* OR recall\* OR sort\* OR categor\* OR debug\* OR classif\* OR \*copy\*)) AND

**Mental Representations:** (((cognit\* OR mental\* OR knowledge OR program\* OR situation\*) N2 (model\* OR represent\* OR plan\* OR structur\* OR map\* OR chunk\* OR slic\*)) OR schema\*)) AND

**Programmers:** (programmer\* OR coder\*)

The search string for the IEEE Xplore Digital Library had to be modified slightly since it was the only database that imposed a limit on the number of wildcard symbols and keywords. Because of the multidisciplinary nature of the Scopus database, its search was limited by subject area to psychology, computer science, and engineering. The date of the last search of all databases was November 8, 2018. There were no restrictions put on the publication dates or publication types (e.g., theses, journal articles, conference papers, books) when performing the database searches.

Landman et al. (2017) found that when using the IEEE Xplore database the number of results were reduced when adding an OR to their search queries. Given their concern regarding inconsistencies, multiple search queries were tested using the IEEE Xplore database. When additional OR statements were added to the search query used in the current review, the number of results increased as expected.

To mitigate systematic bias in reviews that can be caused by publication bias, Kitchenham and Charters (2007) suggest contacting experts and researchers working in the area that may know of unpublished work. A request for unpublished work of relevance was sent to the Psychology of Programming Interest Group (PPIG) discussion group. The PPIG workshop is the primary venue for research in the psychological aspects of programming and its discussion group is subscribed to by leading researchers in the domain. As a result of this request, we received four unpublished papers.

Internet search engines were not used as an information source for our review due to their unsystematic nature and lack of quality control. Google and Google Scholar both use personalization filters that affect the results that are returned by searches (Pariser, 2011). The use of filters by search engines creates unpredictable and inconsistent searches. The lack of standardization could create bias since searches are filtered based on what the search engine determines the user wants resulting in an incomplete retrieval of data. Internet search engines cannot exclude results that are from predatory publishers.

### 3.2. Search strategy

The search results for each of the databases were as follows: Computer Source returned 65 records, ERIC returned 50 records, IEEE Xplore Digital Library returned 280 records, PsycINFO returned 92 records, and Scopus returned 879 records. In total, 1,366 data items were retrieved through database searches. There were 280 duplicate data items removed, leaving 1,158 unique data items as indicated in Fig. 1.

### 3.3. Eligibility criteria

A preliminary screening was performed on the 1,158 unique data items that were extracted from the five database searches. The data items were screened by reading their titles and abstracts. In cases where it was not evident if the criteria were met from the abstract, then the full text was accessed and reviewed. The flow diagram outlining the results of each screening stage is found in Fig. 1.

The screening process and selection of the data items to include in our review were performed by the first author. The guidelines provided by Kitchenham and Charters (2007) suggest that a single researcher discuss included and excluded papers with other experts. Following these guidelines, the first author met regularly

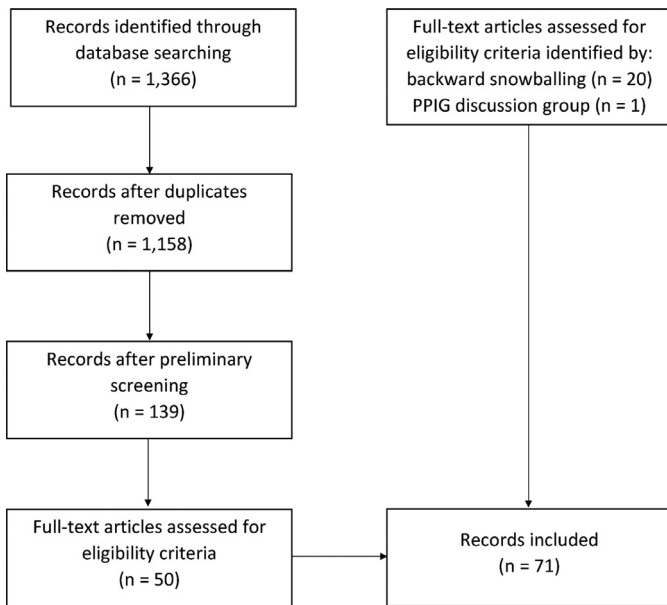


Fig. 1. PRISMA flow diagram.

to discuss the selection of data items with the second and third author who are experts in computer science and psychology respectively. In addition to this, the inter-rater reliability was measured post-hoc by selecting a sample of included and excluded data items that were re-evaluated by the second author. The initial post-hoc inter-rater reliability of the two authors was 49 agreements out of 52 randomly selected papers (92.5%) for the exclusion/inclusion decisions. Refinement of the inclusion criteria allowed an agreement on two of the papers whereas further discussion of the specific papers produced agreement on the remaining papers.

For data items to be included in the review the following inclusion criteria were used:

1. Had to be an empirical study.
2. Study had to have participants that were programmers, meaning that they had prior knowledge of the programming language used in the study.
3. Study had to use coded programs or code fragments that were written in a formal programming language.
4. Study had to contribute to the understanding of the mental representations created during the comprehension process.
5. Study had to contribute to the understanding of the process of building mental representations.

The following exclusion criteria were used to exclude data items from the review:

1. Programming aptitude studies.
2. Usability studies.
3. Studies that involved the teaching and learning of a programming language (e.g., teaching methods, educational software).
4. Studies on the problem solving process of programming.
5. Studies on analytic and predictive techniques for programmer performance.
6. Studies on program design and implementation.
7. Studies on debugging strategies and skills.

After a preliminary screening using the inclusion and exclusion criteria the data items were reduced from 1,158 to 139 (see Fig. 1). In cases where an author republished their study using the same data and analysis, only the most recent publication was kept. Although studies on debugging strategies and skills

were excluded, studies that used debugging as a task to stimulate the comprehension process to study programmers' mental representations were included. For example, the study conducted by Stone et al. (1990) was eliminated because they investigated the debugging skills of programmers and how these skills can be improved. The study by Petre and Blackwell (1999) was excluded because they investigated the visual imagery produced by programmers during the software design phase. Their study did not include the writing or understanding of code. The study by Kamma and Jalote (2013) was eliminated as it examined programmer productivity. Another data item eliminated from the review was a study conducted by Yeh (2014) that did not involve programming or code and investigated the cognitive processes used by participants when solving a software design problem. The criteria for inclusion at this stage required that studies had contributed to the understanding of program comprehension and mental representations.

The documents that were not eliminated during the preliminary screening were then analyzed by reading the full text to determine if they met the eligibility criteria. As a result, 50 documents were included. In the guidelines provided by Kitchenham and Charters (2007) a search of databases alone is not sufficient for a complete systematic review. One of the manual search strategies suggested by Kitchenham and Charters (2007) is to search the reference lists from relevant studies, this is also known as backward snowballing (Wohlin, 2014). The 50 included studies from the automated search were used to perform one iteration of backward snowballing. The data items identified through backward snowballing were also subject to the same eligibility criteria as the data items from the automated search and as a result 20 additional data items were included as indicated in Fig. 1. The four unpublished data items received from the PPIG discussion group were also subject to the same eligibility criteria and resulted in the inclusion of one data item (see Fig. 1).

As a result of the screening process, the data items included in our review are empirical studies of programmers that investigated their comprehension and internal representation of code written in a formal programming language. The eligibility criteria used here are meant to extract studies with a focus on knowledge representation as opposed to knowledge utilization (Atwood and Ramsey, 1978).

### 3.4. Data extraction

The data extraction process was performed by the first author in consultation with the second and third author. The tasks used to stimulate the comprehension process were extracted from the data items by analyzing the description of the study's methods and procedures. Tasks that were assigned to participants during the study to engage them with the code were included in the Task column of Table A.2. The methods and procedures of each study were analyzed to determine how the mental representations held by participants were measured and a summary was written by the first author describing each study's method. Categories describing the methods were formed using the method summaries. After categorizing the methods of all the data items, the most common categories were included in Table A.2. A similar process was used to determine the contributions of each data item. The results of each study were analyzed to determine the main findings of the study and a summary of the findings was written for each study by the first author. Categories describing the types of mental representations formed by the participants of the studies were extracted from the summaries of findings and the most common categories were included in Table A.2.



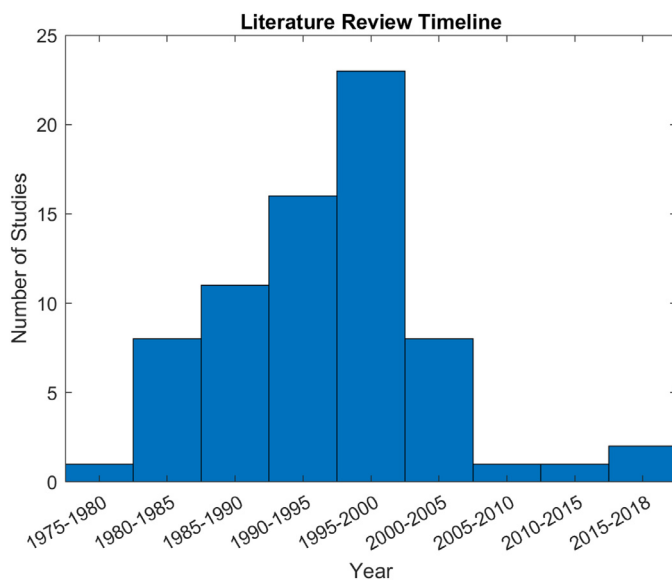


Fig. 2. Frequency of empirical studies on program comprehension.

## 4. Results

The data extracted from the 71 documents identified by the systematic review are provided in Table A.2. The items in the table are listed in chronological order to allow the reader to follow the development of the literature.

### 4.1. Research timeline

Empirical studies used to develop and validate theories on program comprehension first started to emerge in the 1970's (Fig. 2), with the earliest study found by the systematic review conducted in 1976. The timeline depicted in Fig. 2 shows the growth and subsequent decline in the number of empirical studies that have been conducted on program comprehension. During the 1970's and 1980's, 20 studies were conducted. The number of studies grew to 39 in the 1990's when it peaked. In more recent years the work in this area has dropped off almost completely with only 12 studies conducted since 2000 (Fig. 2). Throughout this timeline, various tasks have been used to stimulate the comprehension process that results in the formation of mental representations.

### 4.2. Tasks

When examining the Task column in Table A.2, it appears that, in 23 of the studies, more than one task type was assigned to the participants. All studies that assigned more than one task type to participants included studying code as one of the task types. There were 52 instances where studying code was one of the task types assigned to programmers and in 29 of these instances studying code was the only task type assigned. The task of studying code was often given to participants without revealing the purpose and was most commonly followed by comprehension or recall questions. Studying code was used considerably more than any other task type (see Table 1).

The modification tasks described in the studies included in our review tended to be specific changes in functionality that did not have a far reaching impact on other parts of the code. The modification task used by Navarro-Prieto and Cañas (2001) required programmers to modify a calculation used in the code. Boehm-Davis et al. (1987) required programmers to perform modifications to one or more specific locations in the code.

Table 1

Frequency of tasks used for stimulating the comprehension process.

Task	Number of studies
Study	52
Modification	12
Maintenance	10
Debug	7
Reuse	4
Classify	3
Documentation	3
Write or Reconstruct Code	3
Hand Execution	2
Enhancement	1
Recopy	1

The maintenance tasks used in the studies included in our review required programmers to have an overall understanding of the code and often used industrial code (von Mayrhauser and Vans, 1994; 1995; 1996; 1998). Studies that involved maintenance tasks often observed programmers in the workplace that were familiarizing themselves with code for the purpose of future maintenance or to perform a specific maintenance task (von Mayrhauser and Vans, 1994; 1995; 1996). There were also variations in the type of maintenance tasks and if the participants all performed the same maintenance task. In the study conducted by von Mayrhauser and Vans (1995) all participants performed a maintenance task but not necessarily the same type and Parkin (2004) compared the results of two different types of maintenance tasks. Other studies had programmers all performing the same maintenance task, including the study by Vans et al. (1999) where professional programmers were all performing corrective maintenance tasks and the study by von Mayrhauser and Vans (1998) where all participants performed adaptation tasks. Studies using maintenance tasks were the only studies that involved programmers with varying levels of prior familiarity with the code used in the study (von Mayrhauser and Vans, 1995; 1998; Vans et al., 1999).

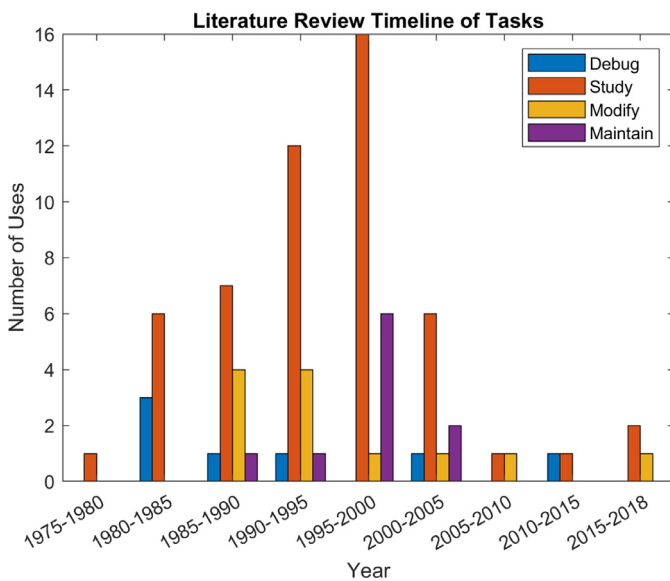
In some cases, the tasks assigned to participants were specific to the programming paradigm of the language used in the study. For example, a reuse task was used to study the mental representations formed by object oriented programmers. Programmers were given code that could be reused in their assigned task either by copying the code to use as a template or by using object oriented design principles of reuse such as inheritance (Burkhardt and Détienne, 1995; Burkhardt et al., 1997; 1998; 2002). The reuse task specifically targets the use of object oriented design principles.

We found that the studies included in our review assigned tasks that varied in their level of realism. Of the most common tasks used to stimulate the comprehension process (see Fig. 3), studying was the only unrealistic task. Despite the observation made by Wiedenbeck et al. (1993) that experts found studying to be an unnatural task, it continued to be the most widely used task (see Fig. 3). Tasks that were considered to be more natural by experts were tasks that had a concrete objective such as debugging or determining the effects of modifications (Wiedenbeck et al., 1993).

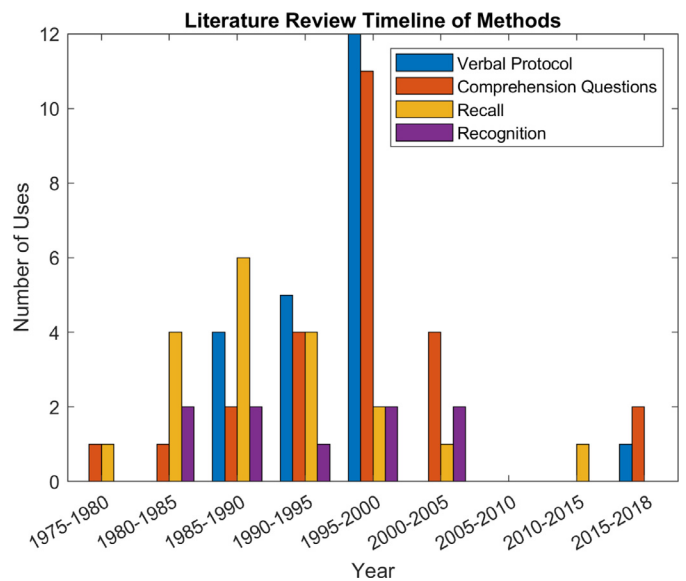
There have been more studies conducted using maintenance tasks than debugging or modification (see Table A.2). However, debugging and modification have been used throughout the timeline and have continued to be used in more recent studies as seen in Fig. 3. Maintenance was more commonly used than debugging and modification for a period of time between 1995-2005 but has not been used since (see Fig. 3).

### 4.3. Research methods

To determine how the methods used in the studies included in our review have changed over time, the method of each study



**Fig. 3.** Frequency of most common tasks for stimulating the comprehension process.



**Fig. 4.** Frequency of most common methods used for determining the mental representations formed during the comprehension process.

has been categorized in Table A.2 for analysis. Through inspection of the Method column in Table A.2, it is apparent that a variety of techniques and measures have been used, both independently and in various combinations, to determine the mental representations formed by participants during the comprehension process. The most common technique was the use of comprehension questions (25 studies). Comprehension questions were also frequently used in addition to other measures to determine the mental representations formed by participants. Each study developed its own set of comprehension questions to measure the strength of different mental representations that may have been formed by the participants. Although some studies tested for the same types of representations (program and domain models studied by Pennington (1987b), Corritore and Wiedenbeck (1991, 1999), Burkhardt et al. (1997), Burkhardt et al. (2002), Wiedenbeck and Ramalingam (1999), Wiedenbeck et al. (1999), Mosemann and Wiedenbeck (2001)) or knowledge structures (i.e., data flow, control flow, state, and function studied by Pennington (1987b), Teasley (1994), Shaft and Vessey (1995), Snyder (1995), Ramalingam and Wiedenbeck (1997), Khazaei and Jackson (2002)) the questions used varied between these studies. The use of verbal protocol analysis while participants performed a task and recall of code were also common techniques for measuring program comprehension (22 and 19 studies respectively). Recognition of code was used in nine studies, summarizing or describing the function of the code was used in five studies, and fill in the blank and sorting were each used three times, whereas all other techniques were each used only once. The recall of code required participants to reproduce the program used in the study either verbatim or a functionally equivalent version. The recognition of code required participants to determine if a given code fragment was from the program used in the study.

Comprehension questions form the method that has been used most consistently throughout the timeline of our review, whereas other techniques have been introduced and in some cases, discontinued at different points along the timeline (Fig. 2). The use of methods depicted in Fig. 2 shows that recall was used more frequently in earlier studies but dropped off around the time that verbal protocol analysis was introduced. The use of verbal protocol analysis first appeared in a study by Letovsky (1987) and was popular until it was last used by Vans et al. (1999) and did not reappear until a study by Nosál and Porubán (2015) (Fig. 4).

Some of the more recent studies have introduced novel techniques for analyzing mental representations formed by programmers such as software that performed screen capture of their actions and monitoring of documents they accessed (Corritore and Wiedenbeck, 2001) and analysis of eye movement data (Fan, 2010).

The data collected in our review and found in the Method column of the extended online appendix<sup>1</sup> indicate that the types of programmers compared in the studies have changed over time. From the first study in 1976 until 1990, the only types of programmers that were compared in the studies were programmers of varying levels of expertise (Shneiderman, 1976; Adelson, 1981; McKeithen et al., 1981; Ehrlich and Soloway, 1984; Soloway and Ehrlich, 1984; Adelson, 1984; Barfield, 1986; Schmidt, 1986; Bateson et al., 1987; Boehm-Davis et al., 1987; Vessey, 1987; Vihmallo and Vihmallo, 1988; Davies, 1990b; Guerin and Matthews, 1990). Programmers were categorized as expert, intermediate, or novice to indicate their expertise in the programming language used in the study or their expertise in the domain relevant to the program. These studies used different definitions and measures to categorize participants' level of expertise. The first occurrence in the data of comparing mental representations formed by programmers trained in different programming paradigms or languages was in 1990. Since that time, comparing programmers with different backgrounds has become a more common type of comparison having been used in studies conducted by Robertson and Yu (1990), Green and Navarro (1995), Corritore and Wiedenbeck (1999), Wiedenbeck and Ramalingam (1999), Wiedenbeck et al. (1999), Corritore and Wiedenbeck (2001), Navarro-Prieto and Cañas (2001), and Khazaei and Jackson (2002). The most common comparison of programming paradigms was between procedural programming and object oriented programming.

#### 4.4. Research contributions

To analyze the contributions that have resulted from the research included in our review, the findings of each study have been categorized in Table A.2. Examination of the Findings column in Table A.2 reveals considerable variation in the description of the representations formed by programmers and the strategies

<sup>1</sup> <http://www.cs.unb.ca/~lbidlak1/ExtendedOnlineAppendix.pdf>.

used during the program comprehension process. The studies conducted by Shneiderman (1976), McKeithen et al. (1981), Ehrlich and Soloway (1984), Mynatt (1984), Barfield (1986, 1997), Guerin and Matthews (1990), Furman (1998), and Fan (2010) found that programmers used chunking to develop a mental representation of the code. Chunking involves the grouping of lines of code together during the comprehension process. The strategy used by programmers when chunking differed between studies. Guerin and Matthews found that programmers chunked by identifying functions in the code; Barfield concluded that programmers grouped sequential lines of code that fit logically together; Furman found that programmers used the visual structure of the code to form chunks; and Fan determined that programmers used beacons to recognize code chunks.

Another group of studies included in our review found that programmers formed mental representations at varying levels of abstraction. Studies performed by Pennington (1987b), Bergantz and Hassell (1991), Burkhardt et al. (1997), Burkhardt et al. (2002), Ramalingam and Wiedenbeck (1997), Corritore and Wiedenbeck (1999), Wiedenbeck and Ramalingam (1999), Wiedenbeck et al. (1999), Mosemann and Wiedenbeck (2001), and Parkin (2004) support the two model theory that programmers form low level program models and high level situation or domain models during the comprehension process. Studies conducted by von Mayrhauser and Vans (1994, 1995, 1996, 1998), Vans (1996), von Mayrhauser et al. (1997), and Vans et al. (1999) found that programmers formed mental models at three levels of abstraction and switched between them during the comprehension process, with the program model as the lowest level, the situation model as the intermediate level, and the domain model as the highest level of abstraction.

The knowledge structures and representations that compose the models at different levels of abstraction were also investigated by a number of studies. Pennington (1987b) found that mental models were developed by programmers using a bottom-up approach (concrete to abstract) since control flow representations were used initially in the comprehension process to form program models whereas data flow and functional representations were used later to form situation models. Shaft and Vessey (1995) determined that the direction in which the representations were developed depended on the expertise of the programmer in the application domain; in an unfamiliar domain programmers developed representations in a bottom-up direction (data flow, control flow, state, and function), but in familiar domains programmers developed representations in the opposite direction (top-down). Contrary to Pennington's findings, Bergantz and Hassell (1991) concluded that control flow representations did not influence the comprehension process and that the representations used by programmers differed depending on their level of expertise: less experienced programmers developed more data flow relationships whereas more experienced programmers developed more function relationships in their mental representations. Teasley (1994) concluded that different types of knowledge structures are all acquired at a similar rate and there was no strong evidence to indicate that programmers use a bottom-up approach. Navarro-Prieto and Cañas (2001) compared programmers with backgrounds in different programming paradigms and found that procedural programmers had better developed control flow representations than data flow representations whereas visual programmers developed both representations equally well. Khazaei and Jackson (2002) compared the representations formed by programmers who had experience in both event driven and object oriented paradigms when understanding programs written in the different paradigms and found that programmers formed stronger control flow, data flow, and functional models when understanding event driven programs compared to object oriented programs. Snyder (1995) determined

that the representations formed by programmers were dependent on the task they were assigned, and that modification tasks required the programmer to develop relationships between four representations: data flow, control flow, state, and function.

Contributions made by research on mental representations formed by programmers during program comprehension have been important in the development of tools and languages that are more intuitive and align with programmers' internal representations. The lack of research on mental representations of parallel programmers reinforces the need to return to this research approach to develop tools and languages for parallel programmers instead of relying on usability and comparative studies. Research that analyses the usability of tools and languages and comparative studies is not informed by theories of program comprehension or mental representations of programmers and is unable to provide insight into how programmers internalize and represent code.

## 5. Discussion

The purpose of our study was to determine the extent of empirical research that examined mental representations formed during program comprehension to date, how the tasks and methods used in the research have changed over time, and the contributions that have resulted from the research.

### 5.1. Decline in research

To determine the extent of empirical research that has examined mental representations formed during program comprehension to date, a systematic review was performed and summarized in Table A.2. The timeline (Fig. 2) that resulted from the systematic review indicates that recently there has been a dramatic decline in research on program comprehension. The declining number of studies focusing on program comprehension and mental representations of programmers in recent years is possibly due to a change in focus. For example, there were a number of studies conducted in recent years that focused on strategies and tools for teaching programming (Oliveira Aureliano, 2013; Dillon, 2013; Lane, 2005), learning second or subsequent programming languages (Scholtz and Wiedenbeck, 1990), and designing programs (Yeh, 2014; Basili and Reiter, 1981) that did not meet the criteria for our review. For example, Whalley and Kasto (2014) conducted a study that examined the progression of learning and the development of cognitive structures during the learning process. Another study that was eliminated measured the amount of programmer effort required to write programs using different parallel programming models (Hochstein et al., 2008). Cañas and Antolí (1998) stated that research of mental representations was blocked due to lack of agreement in definition and methodology. The results of the current study found in Table A.2 support this conjecture by demonstrating the lack of agreement on the tasks used to elicit the mental representations and the methodology utilized to measure them. Other reasons for the observed decline may include a movement towards usability and comparative studies where there may be more funding particularly from companies wanting to demonstrate the usability of their programming languages and tools. Usability and comparative studies tend to have more concrete measures such as lines of code, runtime, and memory usage that allow a quantitative analysis of the program features providing more conclusive results.

### 5.2. Discussion of research questions

By examining the Task and Method columns in Table A.2, we were able to determine how the studies included in our review have changed over time. For a period of time, empirical

studies on program comprehension and mental representations were performed on varying levels of programmers using a variety of programming languages. The tasks and methods used to assess and analyze program comprehension and mental representations of the participants were not consistent between studies. While [Vessey \(1987\)](#) criticized the use of debugging as a task to stimulate program comprehension, it has been used in studies throughout the timeline of our review by [Weiser \(1981, 1982\)](#), [Adelson \(1984\)](#), [Gilmore and Green \(1988\)](#), [Davies \(1990a\)](#), [Romero and Du Boulay \(2004\)](#), and [Fan \(2010\)](#). The research performed by [Wiedenbeck et al. \(1993\)](#) compared expert and novice programmers who were given the task of studying code for understanding and found that studying code was an unnatural task for the experts who commented that normally they would have a concrete objective, such as debugging or predicting the effects of modifications, in mind when reading code. Wiedenbeck speculated that depending on the task, different information may be extracted during program comprehension. Wiedenbeck's observation underlines the importance of ensuring the task is appropriate for the participants in the study. In some cases researchers also expressed concerns over the validity of methods used to measure comprehension including comprehension questions ([Shaft and Vessey, 1995](#)) and recall ([Guerin and Matthews, 1990](#)) that are used in a number of other studies. The focus of the research has also shifted over time. Comparing expert and novice programmers was the focus early on in the research, although no common definition or measure of expertise was used in these studies. The focus has now shifted to comparing programmers with backgrounds in different programming paradigms and languages.

The contributions that resulted from the empirical research on program comprehension have been summarized in the Findings column of [Table A.2](#). Although there are some researchers who build on their own work in an incremental fashion ([Vessey, 1987](#)) or use aspects of other studies as a model ([Guerin and Matthews, 1990](#)), this building process is limited and often deviates from the previous work in such a way that it is hard to connect the findings. Even studies that attempt to support the findings of other studies end up with inconclusive ([Corritore and Wiedenbeck, 1999](#)) or even contradictory results ([Wiedenbeck and Ramalingam, 1999](#)).

### 5.3. General observations

Another finding that emerged from the systematic review was that despite the widespread use of expert and novice as categories to describe programmers, no common definition or measure of expertise has been developed or adopted. Even studies that did not compare expert and novice programmers often used these categories to identify the group of programmers that were used as participants. There was also no consistency between the use of expert and experienced when describing programmers and these terms were often used interchangeably in the same study ([Barfield, 1986](#); [Vessey, 1987](#); [Fix et al., 1993](#)).

## 6. Threats to validity

One threat to validity is the possibility that data items may be missing from our review. Every effort was made to find all relevant data items by selecting databases from both psychology and computer science domains and the multidisciplinary database Scopus. The database selection and development of the search strings was done in consultation with an Information Services Librarian who has experience developing search strings for systematic literature reviews to mitigate the threat of incomplete search terms. There was no restriction on the publication dates which also reduced the threat of missing data items. Backward snowballing was

used to find additional data items that may have been missed by the automated search. To reduce the threat of publication bias as suggested by [Kitchenham and Charters \(2007\)](#), a request was sent to the PPIG discussion group for unpublished work. Although unpublished work might be considered as potentially posing a threat to validity as a result of its perceived poorer quality, the PPIG discussion group is subscribed to by experts in the psychology of programming field, many of whom are authors of published work included in our review.

The first author developed the research questions in consultation with the second and third author in order to mitigate the threat of inappropriate research questions. Preliminary research was also conducted by the first author to gain an understanding of what research questions may be able to provide insight for future empirical studies on mental representations. Another possible threat to validity is the appropriateness of the inclusion and exclusion criteria and the reliability of inclusion decisions. The first author developed the eligibility criteria in consultation with the second and third author, who are experts in high performance computing and cognitive psychology, respectively. To ensure consistency and reduce bias in the inclusion decisions, the first author met regularly with the second and third author to discuss the selection of data items. The results of the post-hoc inter-rater reliability between the first and second author indicates they were in agreement on 92.5% of the data items. This high inter-rater reliability demonstrates consistency and unbiased inclusion decisions during the screening process.

The data extraction process was performed by the first author, introducing a threat to validity as a result of subjective interpretation of the data. To mitigate this threat, the first author discussed the data extraction process with the second and third author and developed the categories for the methods and findings in consultation with them. The first author also held regular meetings with the two co-authors to discuss ambiguous data items as a way to ensure a valid coding of study properties.

## 7. Future Work

The decline in research on mental representations is not an indication that this topic lacks relevance or importance, but an indication of a shift in focus. Research topics that have gained more interest include usability and comparative studies which allow tools and languages to be compared and ranked as more or less superior without the risk of finding that the tool or language does not coincide with the mental representation of the user. [Ericsson et al. \(2006\)](#) found that there has also been more focus on organizational settings where programmers work in teams ([Gren et al., 2017](#); [Teh et al., 2012](#); [Dingsøyr and Dybå, 2012](#)). The study of programming teams would be attractive to software companies as they may anticipate a more immediate return on their investment into this area of research. However, the use of mental model theory to design programming tools and languages has demonstrated measurable benefits. For example, [Sulír and Nosál' \(2015\)](#) studied mental model overlapping to develop source code annotations that allow programmers to share their mental models. From their study they found that the annotations improved program comprehension and reduced maintenance time. The program slicing tool developed by [Korel and Rilling \(1998\)](#) that assists with program comprehension was developed based on research that demonstrated the use of program slicing to improve the process of program understanding. In our view, the current trend in which the usability of languages is assessed after they have already been in use is counterproductive as it ignores cognitive processes involved in programming. Accordingly, we argue that there needs to be a resurgence of empirical research on the psychology of



programming to inform the design of tools and languages, especially in new and emerging programming paradigms.

The study of expert mental representations is important in forming the development of programming languages, instructional practices, and tools. To perform research on expert programmers it is necessary to be able to determine if participants are in fact experts. There has been a lack of agreement among researchers on how expertise should be measured. Siegmund et al. (2013) found that programmer experience can be determined by measuring their self estimation of their experience level compared to their peers and their experience level with logical programming. The study conducted by Baltes and Diehl (2018) found that self-assessment of expertise by programmers was not consistent between programmers with different programming language backgrounds and that years of experience was not related to expertise. The position taken by Parnin et al. (2017) is that expertise cannot be measured using superficial measures such as years of experience but instead using multiple measures such as observing the brain activity of programmers during program comprehension and assessing programming knowledge using concept inventories. To date there remains no standard for measuring programmer expertise so there is a need for more research to develop a standard measure for categorizing programmers based on their expertise. The distinction between expert and experienced also needs to be established. Experience is a measure of time spent working in a particular field or performing a task, however, it does not necessarily translate into expertise, which is a measure of performance (Ericsson et al., 2006). One recommendation for future work is to develop a tool for assessing programmer expertise that is not solely reliant on experience as a gauge.

Empirical research on mental representations formed by programmers during program comprehension has been predominately conducted using sequential code. Studies involving parallel programmers are most often concerned with productivity (Hochstein et al., 2005; Ebcioğlu et al., 2006). However, the mental representations formed by expert parallel programmers during the comprehension of parallel programs is an important area of study to determine how their representations differ from the representations developed during sequential source code comprehension. The comprehension of parallel code requires programmers to mentally execute multiple timelines that are occurring in parallel at the machine level. Therefore, parallel program comprehension may require additional dimensions to construct a mental representation. To explore this research question, program comprehension studies need to be conducted using parallel programmers as participants and assign tasks that stimulate the comprehension process at a level that requires programmers to understand how the code executes in parallel. Possible tasks include identifying the presence of race conditions, rating efficiency or increasing efficiency of parallel programs.

Our literature review suggests that there is no consensus on the method that provides the most accurate account of the mental representations formed by programmers during the comprehension process. In addition, the different methods that have been used only provide an indirect analysis of these mental representations. Our review contains only one study, conducted by Fan (2010), that used eye tracking. Fan used eye tracking data to determine how the program comprehension process is affected by beacons, comments, and task motivation. This author concluded that eye tracking data can be used for tracing and analyzing the program comprehension process.

Work in psychology of programming has been moving towards the use of electroencephalography (EEG) to investigate models of cognition in recent years. For example, Crk et al. (2016) used EEG to determine programmer expertise. However, models of cognition are not the same as the mental representations that are of interest

in our review. In future work it is recommended that eye tracking be used in conjunction with direct questioning to formulate a model of the mental representations formed by programmers during program comprehension.

## 8. Conclusion

Our review contains empirical studies that have been used to build a timeline to provide insight as to how the research on program comprehension and mental representations has evolved. The collection of knowledge contained in our review would be of interest to researchers who want to build on the work done by others in this field and those who want to expand this work to include new programming languages and paradigms. Our review demonstrates that the field of program comprehension is lacking incremental research that builds on previous work, and as a result, the research in this area has been scattered.

The results of our literature review indicate a lack of consistency and agreement in the tasks and methods used for studying mental representations formed during program comprehension. The classifications of programmers based on expertise also varies greatly between studies. There are no common definitions for these classifications and the terms expert and experienced are used interchangeably without considering the difference between these classifications.

Our review also points to gaps in the research on program comprehension. In recent years, the work in this field has declined dramatically and as a result, newly developed or popularized languages and paradigms have not been a part of the research reviewed here. In particular, parallel programming has been neglected in program comprehension research and consequently has developed using mostly informal approaches (Mattson and Wrinn, 2008). Because of the considerable differences between parallel programming and the programming examined in the studies in our review, it is impossible to determine whether the findings summarized in Table A.2 would resemble the comprehension process and mental representations formed by parallel programmers. Therefore, future work should focus on empirical research designed to analyze the mental representations formed by expert parallel programmers during program comprehension to inform the development of tools and languages that support parallel programmers.

## Declaration of Competing Interest

The authors declare that they do not have any financial or non-financial conflict of interests.

## CRediT authorship contribution statement

**Leah Bidlake:** Conceptualization, Data curation, Formal analysis, Investigation, Methodology, Visualization, Writing - original draft, Writing - review & editing. **Eric Aubanel:** Conceptualization, Methodology, Writing - review & editing. **Daniel Voyer:** Conceptualization, Methodology, Writing - review & editing.

## Acknowledgments

I would like to thank Richelle Witherspoon, Information Services Librarian at the University of New Brunswick for her assistance in formulating the search terms and guidance in executing the automated search.

## Appendix A. Summary of Literature

**Table A.2**

Summary of literature included in systematic review. Filled circles indicate that the feature applied and empty circles indicate that it did not apply. An extended table can be found at: <http://www.cs.unb.ca/~lbidlak1/ExtendedOnlineAppendix.pdf>.

Author	Task											Method						Findings				
	Classify	Debug	Documentation	Enhancement	Hand Execution	Maintenance	Modification	Recopy	Reuse	Study	Write or Reconstruct Code	Comprehension Questions	Recall	Recognition	Verbal Protocol	Other	Chunking and Slicing	Mental Model Theory (multilevel)	Top Down vs. Bottom Up	Other		
Shneiderman (1976)	○	○	○	○	●	○	○	○	○	○	○	●	●	○	○	○	●	○	○	○		
Adelson (1981)	○	○	○	○	○	○	○	○	○	○	○	○	●	○	○	○	○	○	○	●		
McKeithen et al. (1981)	○	○	○	○	○	○	○	○	○	○	○	○	●	○	○	○	●	○	○	○		
Weiser (1981)	○	●	○	○	○	○	○	○	○	○	○	○	○	●	○	○	●	○	○	○		
Weiser (1982)	○	●	○	○	○	○	○	○	○	○	○	○	○	●	○	○	●	○	○	○		
Adelson (1984)	○	●	○	○	○	○	○	○	○	○	○	●	○	○	○	○	○	○	○	●		
Mynatt (1984)	○	○	○	○	●	○	○	○	○	○	○	○	●	○	○	○	●	○	○	○		
Ehrlich and Soloway (1984)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	●	○	○	○		
Soloway and Ehrlich (1984)	○	○	○	○	○	○	○	○	○	○	○	○	●	○	○	○	○	○	○	●		
Barfield (1986)	○	○	○	○	○	○	○	○	○	○	○	○	●	○	○	○	○	○	○	○		
Schmidt (1986)	○	○	○	○	○	○	○	○	○	○	○	○	●	○	○	○	○	○	○	●		
Bateson et al. (1987)	○	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	○	○	○	●		
Boehm-Davis et al. (1987)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Letovsky (1987)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Littman et al. (1987)	○	○	○	○	○	●	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Pennington (1987b)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Vessey (1987)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Détienne (1988)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Gilmore and Green (1988)	○	●	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Vihmalo and Vihmalo (1988)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Davies (1990b)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Davies (1990a)	○	●	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Détienne and Soloway (1990)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Guerin and Matthews (1990)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Robertson and Yu (1990)	●	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Bergantz and Hassell (1991)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Corritore and Wiedenbeck (1991)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Koenemann and Robertson (1991)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Koubek and Salvendy (1991)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Wiedenbeck (1991)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Boehm-Davis et al. (1992)	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		

(continued on next page)

Table A.2 (continued)

Author	Task											Method						Findings				
	Classify	Debug	Documentation	Enhancement	Hand Execution	Maintenance	Modification	Recopy	Reuse	Study	Write or Reconstruct Code	Comprehension Questions	Recall	Recognition	Verbal Protocol	Other	Chunking and Slicing	Mental Model Theory (multilevel)	Top Down vs. Bottom Up	Other		
Fix et al. (1993)	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	○	○	○	○	●		
Wiedenbeck et al. (1993)	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	○	○	○	○	●		
Davies (1994)	○	○	○	○	○	○	○	○	○	●	○	○	○	●	○	○	○	○	○	●		
Teasley (1994)	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	○	○	○	○	○		
von Mayrhauser and Vans (1994)	○	○	○	○	○	●	○	○	○	●	○	○	○	○	●	○	○	●	○	○		
Burkhardt and D�tienne (1995)	○	○	○	○	○	○	○	○	●	○	○	○	○	○	●	○	○	○	●	○		
Davies et al. (1995)	●	○	○	○	○	○	○	○	○	●	○	○	○	○	○	●	○	○	○	●		
Green and Navarro (1995)	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	●	○	○	○	●		
Sch�mann (1995)	○	○	○	○	○	○	○	○	○	●	○	●	●	●	○	●	○	○	●	○		
Shaft and Vessey (1995)	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	○	○	○	●	○		
Snyder (1995)	○	○	○	○	○	○	●	○	○	○	○	●	○	○	○	○	○	○	○	●		
von Mayrhauser and Vans (1995)	○	○	○	○	○	●	○	○	○	●	○	○	○	○	○	○	○	●	○	○		
Vans (1996)	○	○	○	○	○	●	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
von Mayrhauser and Vans (1996)	○	○	○	○	○	○	○	○	○	●	○	○	○	○	○	○	○	○	○	○		
Ye and Salvendy (1996)	○	○	○	○	○	○	○	○	○	●	○	○	○	○	○	●	○	○	●	○		
Barfield (1997)	○	○	○	○	○	○	○	●	○	○	○	○	●	○	○	○	●	○	○	○		
Burkhardt et al. (1997)	○	○	●	○	○	○	○	○	●	●	○	●	○	○	○	○	○	○	○	○		
Ramalingam and Wiedenbeck (1997)	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	○	○	○	○	○		
von Mayrhauser et al. (1997)	○	○	○	●	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Burkhardt et al. (1998)	○	○	●	○	○	○	○	○	●	●	○	○	○	○	○	○	○	○	○	○		
Furman (1998)	○	○	○	○	○	○	○	○	○	●	○	○	○	○	○	○	○	○	○	○		
Shaft and Vessey (1998)	○	○	○	○	○	○	○	○	○	●	○	○	○	○	○	○	○	○	○	○		
von Mayrhauser and Vans (1998)	○	○	○	○	○	●	○	○	○	○	○	○	○	○	○	○	○	○	○	○		
Wong et al. (1998)	○	○	○	○	○	○	○	○	○	●	○	○	○	○	○	○	○	○	○	○		
Corritore and Wiedenbeck (1999)	○	○	○	○	○	○	○	○	○	●	○	○	○	○	○	○	○	○	○	○		

(continued on next page)

Table A.2 (continued)

Author	Task											Method						Findings				
	Classify	Debug	Documentation	Enhancement	Hand Execution	Maintenance	Modification	Recopy	Reuse	Study	Write or Reconstruct Code	Comprehension Questions	Recall	Recognition	Verbal Protocol	Other	Chunking and Slicing	Mental Model Theory (multilevel)	Top Down vs. Bottom Up	Other		
Vans et al. (1999)	○	○	○	○	○	●	○	○	○	○	○	○	○	○	●	○	○	○	○	○		
Wiedenbeck and Ramalingam (1999)	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	○	○	○	○	○		
Wiedenbeck et al. (1999)	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	○	○	○	○	○		
Corritore and Wiedenbeck (2001)	○	○	○	○	○	●	○	○	○	●	○	○	○	○	○	●	○	○	●	○		
Mosemann and Wiedenbeck (2001)	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	○	○	○	●	○		
Navarro-Prieto and Cañas (2001)	○	○	○	○	○	○	●	○	○	●	○	○	○	●	○	●	○	○	○	●		
Romero (2001)	○	○	○	○	○	○	○	○	○	●	○	●	●	●	○	○	○	○	○	●		
Burkhardt et al. (2002)	○	○	●	○	○	○	○	○	●	●	○	●	○	○	○	○	○	●	○	○		
Khazaei and Jackson (2002)	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	○	○	○	○	●		
Parkin (2004)	○	○	○	○	○	●	○	○	○	○	○	○	○	○	○	●	○	●	○	○		
Romero and Du Boulay (2004)	○	●	○	○	○	○	○	○	○	○	○	○	○	○	○	●	○	○	○	●		
Sajaniemi and Prieto (2005)	●	○	○	○	○	○	●	○	○	●	○	○	○	○	○	●	○	○	○	●		
Fan (2010)	○	●	○	○	○	○	○	○	○	●	○	○	●	○	○	●	●	○	○	○		
Alardawi and Agil (2015)	○	○	○	○	○	○	○	○	○	●	○	●	○	○	○	○	○	○	○	●		
Nosál and Porubán (2015)	○	○	○	○	○	○	●	○	○	●	○	●	○	○	●	○	○	○	○	●		



## Supplementary material

Supplementary material associated with this article can be found, in the online version, at doi:[10.1016/j.jss.2020.110565](https://doi.org/10.1016/j.jss.2020.110565).

## References

- Adelson, B., 1981. Problem solving and the development of abstract categories in programming languages. *Memory Cognit.* 9 (4), 422–433. doi:[10.3758/BF03197568](https://doi.org/10.3758/BF03197568).
- Adelson, B., 1984. When novices surpass experts: the difficulty of a task may increase with expertise. *J. Exp. Psychol.* 10 (3), 483–495. doi:[10.1037/0278-7393.10.3.483](https://doi.org/10.1037/0278-7393.10.3.483).
- Alardawi, A., Agil, A., 2015. Novice comprehension of object-oriented oo programs: an empirical study. In: 2015 World Congress on Information Technology and Computer Applications (WCITCA), pp. 1–4. doi:[10.1109/WCITCA.2015.7367057](https://doi.org/10.1109/WCITCA.2015.7367057).
- Arab, M., 1992. Enhancing program comprehension: formatting and documenting. *ACM SIGPLAN Notices* 27 (2), 37–46. doi:[10.1145/130973.130975](https://doi.org/10.1145/130973.130975).
- Atwood, M.E., Ramsey, H.R., 1978. Cognitive structures in the comprehension and memory of computer programs: an investigation of computer program debugging. Technical Report. SCIENCE APPLICATIONS INC ENGLEWOOD CO.
- Baltes, S., Diehl, S., 2018. Towards a theory of software development expertise. In: Proceedings of the 2018 26th ACM Joint Meeting on European Software Engineering Conference and Symposium on the Foundations of Software Engineering. ACM, New York, NY, USA, pp. 187–200. doi:[10.1145/3236024.3236061](https://doi.org/10.1145/3236024.3236061).
- Barfield, W., 1986. Expert-novice differences for software: implications for problem-solving and knowledge acquisition. *Behav. Inf. Technol.* 5 (1), 15–29. doi:[10.1080/01449298608914495](https://doi.org/10.1080/01449298608914495).
- Barfield, W., 1997. Skilled performance on software as a function of domain expertise and program organization. *Percept. Motor Skills* 85 (3, Pt 2), 1471–1480. doi:[10.2466/pms.1997.85.3f.1471](https://doi.org/10.2466/pms.1997.85.3f.1471).
- Basili, V.R., Reiter, R.W., 1981. A controlled experiment quantitatively comparing software development approaches. *IEEE Trans. Softw. Eng.* 7 (3), 299–320. doi:[10.1109/TSE.1981.230841](https://doi.org/10.1109/TSE.1981.230841).
- Bateson, A.G., Alexander, R.A., Murphy, M.D., 1987. Cognitive processing differences between novice and expert computer programmers. *Int. J. Man-Mach. Stud.* 26 (6), 649–660. doi:[10.1016/S0020-7373\(87\)80058-5](https://doi.org/10.1016/S0020-7373(87)80058-5).
- Bergantz, D., Hassell, J., 1991. Information relationships in prolog programs: how do programmers comprehend functionality? *Int. J. Man-Mach. Stud.* 35 (3), 313–328. doi:[10.1016/S0020-7373\(05\)80131-2](https://doi.org/10.1016/S0020-7373(05)80131-2).
- Boehm-Davis, D.A., Holt, R.W., Schultz, A.C., 1987. Development and use of cognitive representations of software in a modification task. In: Proceedings of the 1987 IEEE International Conference on Systems, Man and Cybernetics, pp. 990–994.
- Boehm-Davis, D.A., Holt, R.W., Schultz, A.C., 1992. The role of program structure in software maintenance. *Int. J. Man-Mach. Stud.* 36 (1), 21–63. doi:[10.1016/0020-7373\(92\)90051-L](https://doi.org/10.1016/0020-7373(92)90051-L).
- Boshernitsan, M., Graham, S.L., Hearst, M.A., 2007. Aligning development tools with the way programmers think about code changes. In: Proceedings of the SIGCHI conference on Human factors in computing systems. ACM, pp. 567–576. doi:[10.1145/1240624.1240715](https://doi.org/10.1145/1240624.1240715).
- Burkhardt, J., Détienné, F., Wiedenbeck, S., 1998. The effect of object-oriented programming expertise in several dimensions of comprehension strategies. In: Proceedings. 6th International Workshop on Program Comprehension. IWPC'98 (Cat. No.98TB100242), pp. 82–89. doi:[10.1109/WPC.1998.693294](https://doi.org/10.1109/WPC.1998.693294).
- Burkhardt, J.-M., Détienné, F., 1995. An empirical study of software reuse by experts in object-oriented design. In: Human-Computer Interaction: INTERACT'95. In: IFIP Advances in Information and Communication Technology. Springer US, pp. 133–138. doi:[10.1007/978-1-5041-2896-4\\_22](https://doi.org/10.1007/978-1-5041-2896-4_22).
- Burkhardt, J.-M., Détienné, F., Wiedenbeck, S., 1997. Mental representations constructed by experts and novices in object-oriented program comprehension. In: Human-Computer Interaction INTERACT '97. Springer US, pp. 339–346. doi:[10.1007/978-0-387-35175-9\\_55](https://doi.org/10.1007/978-0-387-35175-9_55).
- Burkhardt, J.-M., Détienné, F., Wiedenbeck, S., 2002. Object-oriented program comprehension: effect of expertise, task and phase. *Empir. Softw. Eng.* 7 (2), 115–156. doi:[10.1023/A:1015297914742](https://doi.org/10.1023/A:1015297914742).
- Cañas, J.J., Antolí, A., 1998. The role of working memory in measuring mental models. In: Proceedings of the Ninth European Conference on Cognitive Ergonomics-Cognition and Cooperation. European Association of Cognitive Ergonomics (EACE): Rocquencourt, France.
- Chamberlain, B.L., Deitz, S.J., Snyder, L., 2000. A comparative study of the NAS MG benchmark across parallel languages and architectures. In: Supercomputing, ACM/IEEE 2000 Conference doi:[10.1109/SC.2000.10006](https://doi.org/10.1109/SC.2000.10006), pp. 46–46.
- Corritore, C., Wiedenbeck, S., 1991. What do novices learn during program comprehension? *Int. J. Hum.-Comput. Interact.* 3 (2), 199–222. doi:[10.1080/10447319109526004](https://doi.org/10.1080/10447319109526004).
- Corritore, C.L., Wiedenbeck, S., 1999. Mental representations of expert procedural and object-oriented programmers in a software maintenance task. *Int. J. Hum.-Comput. Stud.* 50 (1), 61–83. doi:[10.1006/ijhc.1998.0236](https://doi.org/10.1006/ijhc.1998.0236).
- Corritore, C.L., Wiedenbeck, S., 2001. An exploratory study of program comprehension strategies of procedural and object-oriented programmers. *Int. J. Hum.-Comput. Stud.* 54 (1), 1–23. doi:[10.1006/ijhc.2000.0423](https://doi.org/10.1006/ijhc.2000.0423).
- Crk, I., Kluthe, T., Stefik, A., 2016. Understanding programming expertise: an empirical study of phasic brain wave changes. *ACM Trans. Comput.-Hum. Interact.* 23 (1), 1–29. doi:[10.1145/2829945](https://doi.org/10.1145/2829945).
- Davies, S., 1990a. The nature and development of programming plans. *Int. J. Man-Mach. Stud.* 32 (4), 461–481. doi:[10.1016/S0020-7373\(05\)80143-9](https://doi.org/10.1016/S0020-7373(05)80143-9).
- Davies, S.P., 1990b. Plans, goals and selection rules in comprehension of computer programs. *Behav. Inf. Technol.* 9 (3), 201–214. doi:[10.1080/01449299008924237](https://doi.org/10.1080/01449299008924237).
- Davies, S.P., 1994. Knowledge restructuring and the acquisition of programming expertise. *Int. J. Hum.-Comput. Stud.* 40 (4), 703–726. doi:[10.1006/ijhc.1994.1032](https://doi.org/10.1006/ijhc.1994.1032).
- Davies, S.P., Gilmore, D.J., Green, T.R.G., 1995. Are objects that important? Effects of expertise and familiarity on classification of object-oriented code. *Hum.-Comput. Interact.* 10 (2–3), 227–248. doi:[10.1207/s15327051hci1002&3\\_3](https://doi.org/10.1207/s15327051hci1002&3_3).
- Détienné, F., 1988. Une application de la théorie des schémas à la compréhension de programmes. = applying schema theory to program understanding. *Le Travail Humain* 51 (4), 335–350.
- Détienné, F., 2001. Software Design-Cognitive Aspect. Springer Science & Business Media.
- Détienné, F., Soloway, E., 1990. An empirically-derived control structure for the process of program understanding. *Int. J. Man-Mach. Stud.* 33 (3), 323–342. doi:[10.1016/S0020-7373\(05\)80122-1](https://doi.org/10.1016/S0020-7373(05)80122-1).
- Dillon, E.C.J., 2013. Measuring the effects of low assistive vs. moderately assistive environments on novice programmers. ProQuest Information & Learning Ph.D. thesis.
- Dingsøyr, T., Dybå, T., 2012. Team effectiveness in software development: human and cooperative aspects in team effectiveness models and priorities for future studies. In: 2012 5th International Workshop on Co-operative and Human Aspects of Software Engineering (CHASE), pp. 27–29. doi:[10.1109/CHASE.2012.6223016](https://doi.org/10.1109/CHASE.2012.6223016).
- Ebcioğlu, K., Sarkar, V., El-Ghazawi, T., Urbanic, J., Center, P.S., 2006. An experiment in measuring the productivity of three parallel programming languages. In: Proceedings of the Third Workshop on Productivity and Performance in High-End Computing, pp. 30–36.
- Ehrlich, K., Soloway, E., 1984. An empirical investigation of the tacit plan knowledge in programming. In: Human factors in computer systems, 16, pp. 113–134.
- Ericsson, K., Charness, N.E., Feltovich, P.J., Hoffman, R.R., 2006. The cambridge handbook of expertise and expert performance. Cambridge Handbooks in Psychology. Cambridge University Press doi:[10.1017/CBO9780511816796](https://doi.org/10.1017/CBO9780511816796).
- Fan, Q., 2010. The effects of beacons, comments, and tasks on program comprehension process in software maintenance Ph.D. thesis.
- Feo, J.T., 2015. A comparative study of parallel programming languages: the salishan problems. Special Topics in Supercomputing, v. 6. Elsevier Science.
- Fix, V., Wiedenbeck, S., Scholtz, J., 1993. Mental representations of programs by novices and experts. In: Conference on Human Factors in Computing Systems - Proceedings, pp. 74–79. doi:[10.1145/169059.169088](https://doi.org/10.1145/169059.169088).
- Furman, S.M., 1998. Improving software comprehension. ProQuest Information & Learning Ph.D. thesis.
- Gilmore, D.J., Green, T.R., 1988. Programming plans and programming expertise. *Q. J. Exp. Psychol. A* 40 (3-A), 423–442. doi:[10.1080/02724988843000005](https://doi.org/10.1080/02724988843000005).
- Green, T.R.G., 1989. Cognitive dimensions of notations. In: Proceedings of the Fifth Conference of the British Computer Society, Human-Computer Interaction Specialist Group on People and Computers V. Cambridge University Press, New York, NY, USA, pp. 443–460.
- Green, T.R.G., Navarro, R., 1995. Programming plans, imagery, and visual programming. In: Human-Computer Interaction: Interact '95. In: IFIP Advances in Information and Communication Technology. Springer US, pp. 139–144. doi:[10.1007/978-1-5041-2896-4\\_23](https://doi.org/10.1007/978-1-5041-2896-4_23).
- Gren, L., Torkar, R., Feldt, R., 2017. Group development and group maturity when building agile teams: a qualitative and quantitative investigation at eight large companies. *J. Syst. Softw.* 124, 104–119. doi:[10.1016/j.jss.2016.11.024](https://doi.org/10.1016/j.jss.2016.11.024).
- Guerin, B., Matthews, A., 1990. The effects of semantic complexity on expert and novice computer program recall and comprehension. *J. Gen. Psychol.* 117 (4), 379–389. doi:[10.1080/00221309.1990.9921144](https://doi.org/10.1080/00221309.1990.9921144).
- Hochstein, L., Basili, V.R., Vishkin, U., Gilbert, J., 2008. A pilot study to compare programming effort for two parallel programming models. *J. Syst. Softw.* 81 (11), 1920–1930. doi:[10.1016/j.jss.2007.12.798](https://doi.org/10.1016/j.jss.2007.12.798).
- Hochstein, L., Carver, J., Shull, F., Asgari, S., Basili, V., Hollingsworth, J.K., Zelkowitz, M.V., 2005. Parallel programmer productivity: a case study of novice parallel programmers. In: Supercomputing, 2005. Proceedings of the ACM/IEEE SC 2005 Conference doi:[10.1109/SC.2005.53](https://doi.org/10.1109/SC.2005.53), pp. 35–35.
- Hornbæk, K., 2006. Current practice in measuring usability: challenges to usability studies and research. *Int. J. Hum.-Comput. Stud.* 64 (2), 79–102. doi:[10.1016/j.ijhcs.2005.06.002](https://doi.org/10.1016/j.ijhcs.2005.06.002).
- Kamma, D., Jalote, P., 2013. Effect of task processes on programmer productivity in model-based testing. In: Proceedings of the 6th India software engineering conference. ACM, pp. 23–28. doi:[10.1145/2442754.2442758](https://doi.org/10.1145/2442754.2442758).
- Khazaei, B., Jackson, M., 2002. Is there any difference in novice comprehension of a small program written in the event-driven and object-oriented styles? In: Proceedings - IEEE 2002 Symposium on Human Centric Computing Languages and Environments, HCC 2002, pp. 19–26. doi:[10.1109/HCC.2002.1046336](https://doi.org/10.1109/HCC.2002.1046336).
- Kitchenham, B., Charters, S., 2007. Guidelines for performing systematic literature reviews in software engineering. Technical Report. Keele University and University of Durham.
- Koenemann, J., Robertson, S.P., 1991. Expert problem solving strategies for program comprehension. In: Proceedings of the SIGCHI Conference on Human Factors in Computing Systems. ACM, pp. 125–130. doi:[10.1145/108844.108863](https://doi.org/10.1145/108844.108863).
- Korel, B., Rilling, J., 1998. Program slicing in understanding of large programs. In: Proceedings. 6th International Workshop on Program Comprehension. IWPC'98 (Cat. No.98TB100242), pp. 145–152. doi:[10.1109/WPC.1998.693339](https://doi.org/10.1109/WPC.1998.693339).

- Koubek, R.J., Salvendy, G., 1991. Cognitive performance of super-experts on computer program modification tasks. *Ergonomics* 34 (8), 1095–1112. doi:[10.1080/00140139108964849](#).
- Landman, D., Serebrenik, A., Vinju, J.J., 2017. Challenges for static analysis of java reflection-literature review and empirical study. In: 2017 IEEE/ACM 39th International Conference on Software Engineering (ICSE). IEEE, pp. 507–518. doi:[10.1109/ICSE.2017.53](#).
- Lane, H.C., 2005. Natural language tutoring and the novice programmer. ProQuest Information & Learning Ph.D. thesis.
- Letovsky, S., 1987. Cognitive processes in program comprehension. *The J. Syst. Softw.* 7 (4), 325–339. doi:[10.1016/0164-1212\(87\)90032-X](#).
- Littman, D., Pinto, J., Letovsky, S., Soloway, E., 1987. Mental models and software maintenance. *The J. Syst. Softw.* 7 (4), 341–355. doi:[10.1016/0164-1212\(87\)90033-1](#).
- Mattson, T., Wrinn, M., 2008. Parallel programming: can we please get it right this time? In: Proceedings of the 45th Annual Design Automation Conference. ACM, New York, NY, USA, pp. 7–11. doi:[10.1145/1391469.1391474](#).
- von Mayrhauser, A., Vans, A., 1994. Comprehension processes during large scale maintenance. In: Proceedings of 16th International Conference on Software Engineering, pp. 39–48. doi:[10.1109/ICSE.1994.296764](#).
- von Mayrhauser, A., Vans, A., 1995. Industrial experience with an integrated code comprehension model. *Softw. Eng. J.* 10 (5), 171–182. doi:[10.1049/sej.1995.0023](#).
- von Mayrhauser, A., Vans, A., 1996. Identification of dynamic comprehension processes during large scale maintenance. *IEEE Trans. Softw. Eng.* 22 (6), 424–437. doi:[10.1109/32.508315](#).
- von Mayrhauser, A., Vans, A.M., 1998. Program understanding behavior during adaptation of large scale software. In: Proceedings. 6th International Workshop on Program Comprehension. IWPC'98 (Cat. No.98TB100242), pp. 164–172. doi:[10.1109/WPC.1998.693345](#).
- von Mayrhauser, A., Vans, A.M., Howe, A.E., 1997. Program understanding behaviour during enhancement of large-scale software. *J. Softw. Maint.* 9 (5), 299–327. doi:[10.1002/\(SICI\)1096-908X\(199709\)10:5<299::AID-SMR157>3.0.CO;2-S](#).
- McKeithen, K., Reitman, J., Rueter, H., Hirtle, S., 1981. Knowledge organization and skill differences in computer programmers. *Cognit. Psychol.* 13 (3), 307–325. doi:[10.1016/0010-0285\(81\)90012-8](#).
- Moher, D., Liberati, A., Tetzlaff, J., Altman, D.G., Group, T.P., 2009. Preferred reporting items for systematic reviews and meta-analyses: the prisma statement. *PLOS Med.* 6 (7). doi:[10.1371/journal.pmed.1000097](#).
- Mosemann, R., Wiedenbeck, S., 2001. Navigation and comprehension of programs by novice programmers. In: Proceedings 9th International Workshop on Program Comprehension. IWPC 2001. IEEE, pp. 79–88. doi:[10.1109/WPC.2001.921716](#).
- Mynatt, B.T., 1984. The effect of semantic complexity on the comprehension of program modules. *Int. J. Man-Mach. Stud.* 21 (2), 91–103. doi:[10.1016/S0020-7373\(84\)80060-7](#).
- Nanz, S., Furia, C.A., 2015. A comparative study of programming languages in rosetta code. In: 2015 IEEE/ACM 37th IEEE International Conference on Software Engineering, 1, pp. 778–788. doi:[10.1109/ICSE.2015.90](#).
- Navarro-Prieto, R., Cañas, J.J., 2001. Are visual programming languages better? The role of imagery in program comprehension. *Int. J. Hum.-Comput. Stud.* 54 (6), 799–829. doi:[10.1006/jhc.2000.0465](#).
- Nosál, M., Porubán, J., 2015. Program comprehension with four-layered mental model. In: 2015 13th International Conference on Engineering of Modern Electric Systems (EMES), pp. 1–4. doi:[10.1109/EMES.2015.7158420](#).
- Oliveira Aureliano, V.C., 2013. A methodology for teaching programming for beginners. In: Proceedings of the ninth annual international ACM conference on International computing education research. ACM, pp. 169–170. doi:[10.1145/2493394.2493417](#).
- Pariser, E., 2011. The filter bubble: how the new personalized web is changing what we read and how we think. Penguin.
- Parkin, P., 2004. An exploratory study of code and document interactions during task-directed program comprehension. In: 2004 Australian Software Engineering Conference. Proceedings., 2004, pp. 221–230. doi:[10.1109/ASWEC.2004.1290475](#).
- Parnin, C., Siegmund, J., Peitek, N., 2017. On the nature of programmer expertise. In: PPIG, p. 16.
- Pennington, N., 1987a. Comprehension strategies in programming. In: Empirical Studies of Programmers: Second Workshop. Ablex Publishing Corp., pp. 100–113.
- Pennington, N., 1987b. Stimulus structures and mental representations in expert comprehension of computer programs. *Cognit. Psychol.* 19 (3), 295–341. doi:[10.1016/0010-0285\(87\)90007-7](#).
- Petre, M., Blackwell, A.F., 1999. Mental imagery in program design and visual programming. *Int. J. Hum.-Comput. Stud.* 51 (1), 7–30. doi:[10.1006/jhc.1999.0267](#).
- Prechelt, L., 2000. An empirical comparison of seven programming languages. *Computer* 33 (10), 23–29. doi:[10.1109/2.876288](#).
- Ramalingam, V., Wiedenbeck, S., 1997. An empirical study of novice program comprehension in the imperative and object-oriented styles. In: Papers Presented at the 7th Workshop on Empirical Studies of Programmers, ESP 1997, pp. 124–139. doi:[10.1145/266399.266411](#).
- Robertson, S., Yu, C.-C., 1990. Common cognitive representations of program code across tasks and languages. *Int. J. Man-Mach. Stud.* 33 (3), 343–360. doi:[10.1016/S0020-7373\(05\)80123-3](#).
- Romero, P., 2001. Focal structures and information types in prolog. *Int. J. Hum.-Comput. Stud.* 54 (2), 211–236. doi:[10.1006/jhc.2000.0408](#).
- Romero, P., Du Boulay, B., 2004. Structural knowledge and language notational properties in program comprehension. In: 2004 IEEE Symposium on Visual Languages - Human Centric Computing, pp. 223–225. doi:[10.1109/VLHCC.2004.50](#).
- Sajaniemi, J., Prieto, R.N., 2005. An investigation into professional programmers' mental representations of variables. In: Proceedings - IEEE Workshop on Program Comprehension, pp. 55–64. doi:[10.1109/WPC.2005.8](#).
- Schmidt, A.L., 1986. Effects of experience and comprehension on reading time and memory for computer programs. *Int. J. Man-Mach. Stud.* 25 (4), 399–409. doi:[10.1016/S0020-7373\(86\)80068-2](#).
- Scholtz, J., Wiedenbeck, S., 1990. Learning second and subsequent programming languages: a problem of transfer. *Int. J. Hum.-Comput. Interact.* 2 (1), 51–72. doi:[10.1080/10447319009525970](#).
- Schömann, M., 1995. Knowledge organization of novices and advanced programmers: effect of previous knowledge on recall and recognition of lisp-procedures. In: Cognition and computer programming. Ablex Publishing, pp. 193–217.
- Shaft, T.M., Vessey, I., 1995. Research report-the relevance of application domain knowledge: the case of computer program comprehension. *Info. Sys. Res.* 6 (3), 286–299. doi:[10.1287/isre.6.3.286](#).
- Shaft, T.M., Vessey, I., 1998. The relevance of application domain knowledge: characterizing the computer program comprehension process. *J. Manag. Inf. Syst.* 15 (1), 51–78. doi:[10.1080/07421222.1998.11518196](#).
- Shneiderman, B., 1976. Exploratory experiments in programmer behavior. *Int. J. Comput. Inf. Sci.* 5 (2), 123–143. doi:[10.1007/BF00975629](#).
- Siegmund, J., Kästner, C., Liebig, J., Apel, S., Hanenberg, S., 2013. Measuring and modeling programming experience. *Empir. Softw. Eng.* 19 (5), 1299–1334. doi:[10.1007/s10664-013-9286-4](#).
- Snyder, J.R., 1995. The role of local program information in software maintenance productivity. ProQuest Information & Learning Ph.D. thesis.
- Soloway, E., Ehrlich, K., 1984. Empirical studies of programming knowledge. *IEEE Trans. Softw. Eng.* SE-10 (5), 595–609. doi:[10.1109/TSE.1984.5010283](#).
- Stone, D.N., Jordan, E.W., Wright, M.K., 1990. The impact of pascal education on debugging skill. *Int. J. Man-Mach. Stud.* 33 (1), 81–95. doi:[10.1016/S0020-7373\(05\)80116-6](#).
- Storey, M.-A., 2006. Theories, tools and research methods in program comprehension: past, present and future. *Softw. Qual. J.* 14 (3), 187–208. doi:[10.1007/s11219-006-9216-4](#).
- Sulir, M., Nosál, M., 2015. Sharing developers' mental models through source code annotations. In: 2015 Federated Conference on Computer Science and Information Systems (FedCSIS), p. 997–1006. doi:[10.15439/2015F301](#).
- Teasley, B., 1994. The effects of naming style and expertise on program comprehension. *Int. J. Hum. - Comput. Stud.* 40 (5), 757–770. doi:[10.1006/jhc.1994.1036](#).
- Teh, A., Baniassad, E., Rooy, D.v., Boughton, C., 2012. Social psychology and software teams: establishing task-effective group norms. *IEEE Softw.* 29 (4), 53–58. doi:[10.1109/MS.2011.157](#).
- Tubaishat, A., 2001. A knowledge base for program debugging. In: AICCSA '01 Proceedings of the ACS/IEEE International Conference on Computer Systems and Applications, 2001-January, pp. 321–327. doi:[10.1109/AICCSA.2001.934005](#).
- Vans, A.M., 1996. A multi-level code comprehension model for large-scale software. ProQuest Information & Learning Ph.D. thesis.
- Vans, A.M., von Mayrhauser, A., Somlo, G., 1999. Program understanding behavior during corrective maintenance of large-scale software. *Int. J. Hum.-Comput. Stud.* 51 (1), 31–70. doi:[10.1006/jhc.1999.0268](#).
- Vessey, I., 1987. On matching programmers' chunks with program structures: an empirical investigation. *Int. J. Man-Mach. Stud.* 27 (1), 65–89. doi:[10.1016/S0020-7373\(87\)80044-5](#).
- Vihmal, A., Vihmal, M., 1988. Utilization of subject's background knowledge in computer program comprehension. *Zeitschrift für Psychologie mit Zeitschrift für angewandte Psychologie* 196 (4), 401–413.
- Weiser, M., 1981. Program slicing. In: Proceedings - International Conference on Software Engineering, pp. 439–449.
- Weiser, M., 1982. Programmers use slices when debugging. *Commun. ACM* 25 (7), 446–452. doi:[10.1145/358557.358577](#).
- Whalley, J., Kasto, N., 2014. A qualitative think-aloud study of novice programmers' code writing strategies. In: Proceedings of the 2014 conference on Innovation & Technology in Computer Science Education, pp. 279–284. doi:[10.1145/2591708.2591762](#).
- Wiedenbeck, S., 1991. The initial stage of program comprehension. *Int. J. Man-Mach. Stud.* 35 (4), 517–540. doi:[10.1016/S0020-7373\(05\)80090-2](#).
- Wiedenbeck, S., Fix, V., Scholtz, J., 1993. Characteristics of the mental representations of novice and expert programmers: an empirical study. *Int. J. Man-Mach. Stud.* 39 (5), 793–812. doi:[10.1006/imms.1993.1084](#).
- Wiedenbeck, S., Ramalingam, V., 1999. Novice comprehension of small programs written in the procedural and object-oriented styles. *Int. J. Hum.-Comput. Stud.* 51 (1), 71–87. doi:[10.1006/jhc.1999.0269](#).
- Wiedenbeck, S., Ramalingam, V., Sarasamma, S., Corritore, C.L., 1999. A comparison of the comprehension of object-oriented and procedural programs by novice programmers. *Interact. Comput.* 11 (3), 255–282. doi:[10.1016/S0953-5438\(98\)00029-0](#).
- Wohlin, C., 2014. Guidelines for snowballing in systematic literature studies and a replication in software engineering. In: Proceedings of the 18th International Conference on Evaluation and Assessment in Software Engineering, pp. 1–10. doi:[10.1145/2601248.2601268](#).
- Wong, S.-Y., Cheung, H., Chen, H.-C., 1998. The advanced programmer's reliance on program semantics: evidence from some cognitive tasks. *Int. J. Psychol.* 33 (4), 259–268. doi:[10.1080/002075998400303](#).
- Ye, N., Salvendy, G., 1996. An objective approach to exploring skill differences in strategies of computer program comprehension. *Behav. Inf. Technol.* 15 (3), 139–148. doi:[10.1080/014492996120229](#).

Yeh, K.-C., 2014. Toward understanding the cognitive processes of software design in novice programmers. ProQuest Information & Learning Ph.D. thesis.

**Leah Bidlake** has been a faculty member in the Faculty of Computer Science at the University of New Brunswick since 2016. She received her Bachelor and Master of Computer Science from UNB and is currently working towards her PhD in Computer Science. Her research interest is in psychology of programming.

**Dr. Eric Aubanel** has been a faculty member in the Faculty of Computer Science since 2002. He received his Bachelor of Science from Trent University and his PhD from Queen's University. Eric's current research includes the psychology of programming and the development of tools and algorithms to support

scientific computing on heterogeneous distributed resources, including manycore accelerators. In 2016 he published a graduate textbook on parallel computing, *Elements of Parallel Computing*, with CRC Press.

**Daniel Voyer** completed his undergraduate studies and masters degree at the Université de Montreal and his doctorate in cognitive psychology at the University of Waterloo. He worked at St. Francis Xavier University from 1991 to 2000. He has been a professor at the University of New Brunswick since August 2000. He has published extensively in journals such as *Psychological Bulletin*, *Brain and Cognition*, and *Laterality*. His research expertise lies in hemispheric specialization, individual and sex differences in cognition, auditory perception, and meta-analysis.