1    **Topic:** Human-Computer Interaction, Computer Systems

2

3    **Development and validation of a Descriptive Cognitive Model for**

4    **predicting usability issues in a Low Code Development Platform**

5    Carlos Silva[1], Joana Vieira[1], José C. Campos[2,3], Rui Couto[2,3], António N. Ribeiro[2,3]

6    [1]Center for Computer Graphics, Guimarães, Portugal

7    [2]Department of Informatics, University of Minho, Braga, Portugal

8    [3]HASLab/INESC TEC, Braga, Portugal

9

10    **Corresponding author:** Rui Couto, rmscouto@gmail.com

13    <div align="center">**PRÉCIS**</div>

14    This study proposes and evaluates a Descriptive Cognitive Model (DCM) for the identification

15    of initial usability issues in a low-code development platform (LCDP). By applying the

16    proposed DCM we were able to predict the interaction problems felt by first-time users of the

17    LCDP.

18    <div align="center">**ABSTRACT**</div>

19    **Objective**: Development and evaluation of a Descriptive Cognitive Model (DCM) for the identification of three

20    types of usability issues in a low-code development platform (LCDP).

21    **Background:** LCDPs raise the level of abstraction of software development by freeing end-users from

22    implementation details. An effective LCDP requires an understanding of how its users conceptualize

23    programming. It is necessary to identify the gap between the LCDP end-users' conceptualization of

24    programming, and the actions required by the platform. It is also relevant to evaluate how the conceptualization

25    of the programming tasks varies according to the end-users' skills.

26    **Method:** DCMs are widely used in the description and analysis of the interaction between users and systems.

27    We propose a DCM which we called PRECOG that combines task-decomposition methods with knowledge-

based descriptions and criticality analysis. This DCM was validated using empirical techniques to provide the best insight regarding the users' interaction performance. Twenty programmers (10 experts, 10 novices) were observed using a LCDP and their interactions were analyzed according to our DCM.

**Results:** The DCM correctly identified several problems felt by first-time platform users. The patterns of issues observed were qualitatively different between groups. Experts mainly faced interaction related problems, while novices faced problems attributable to a lack of programming skills.

**Conclusion:** Applying the proposed DCM we were able to predict three types of interaction problems felt by first time users of the LCDP.

**Application:** The method is applicable when it is relevant to identify possible interaction problems, resulting from the users' background knowledge being insufficient to guarantee a successful completion of the task at hand.

**Keywords**: End-User Development, Low-Code Development Platforms, Descriptive Cognitive Models, Usability, Human-Computer Interaction

# INTRODUCTION

Low-code development platforms (LCDP) address the need for increased productivity in software development. By raising the abstraction level at which software is developed, they automate low-level and routine development tasks, effectively contributing to solve the problem of global shortage of professional software developers. Forrester's Low-Code Market Forecast predicts low-code platforms will reach over 15 billion US dollars in 2020 (Marvin, 2018). At the same time, they lower the entry barrier to software development. As these low-level tasks become automated, developers are not required to carry them out (or even know how to carry them out). Low-level technical details are effectively hidden by the platform. If the entry level becomes low enough, we can say these platforms become End-User Development (EUD) platforms (Fischer, Giaccardi, Ye, Sutcliffe, & Mehandjiev, 2004). At that point, no special programming skills are needed to use them. Other terms have been used to describe related concepts with varying levels of scope, such as End-User Programming (EUP), End-User Software Engineering (EUSE) and Meta-Design (see Barricelli, Cassano, Fogli, & Piccinno (2019) for a recent systematic review of the literature).

55    Whether considering LCDP or EUD, the users' prior knowledge plays a relevant role in the

56    learning and using of a platform, as it will affect the way users approach the platform (Dijkstra, 1982).

57    In the case of LCDP, there is the double challenge of supporting users with little or no knowledge of

58    programming, while also supporting expert programmers. Indeed, understanding individual

59    differences and expectations, and identifying the sources of variation among different users will help

60    this type of platforms to be more broadly adopted (Blackwell, 2017). Since low-code development

61    platforms aim at reducing the learning burden while providing powerful tools to address a wide range

62    of problems, a trade-off must be established between the scope of application and the learning costs

63    of the platforms and their languages. This necessarily implies building an understanding of how

64    different types of users approach the platforms.

65    Descriptive Cognitive Models (DCM) can be used to study the interaction between one

66    interactive system and its users, in particular to analyze how the interplay between the users' cognitive

67    processes and the user interfaces' design might lead to faulty interactions or use errors (Nielsen,

68    1994). Its applicability to reason about the act of programming has long been explored (cf. Blackwell,

69    Petre & Church, 2019). Nevertheless, in spite of relevant Human-Computer Interaction (HCI)

70    findings and developments since the 1980s and recent developments in both LCDP and EUP, there is

71    still a considerable number of relevant gaps in current knowledge about how people reason during

72    programming and development tasks (Sajaniemi, 2008). According to Myers, Pane, and Ko (2004),

73    conventional programming languages require the user or programmer to make "*tremendous*

74    *transformations*" (pp.48) from what he or she intends to accomplish, to what he or she should code.

75    Visual modelling languages, typically adopted by low-code development platforms, aim to mitigate

76    this problem, but their actual effectiveness is still subject to debate.

77    The distance between the mental and the physical spaces in software development was the

78    motivation behind the current work. More specifically, the long-term goal of this work is to support

79    lowering the learning curve of a specific LCDP to the point that non-programmers (i.e., end-users)

80    might use it to develop software (in practice, turning it into an EUD platform). The challenge then, is

81    how to reduce the learning effort of users without reducing the scope of the possible application

82    domains. As a contribution to this long-term goal, the work described in this paper aimed at

83  understanding the difficulties faced by potential programmers with different expectations and

84  academic backgrounds when using a specific LCDP. To achieve this, we developed a new descriptive

85  cognitive model with the purpose of predicting usability issues in a LCDP.

86  **THE LCDP – Low-Code Development Platform**

87  A low-code development platform supports the development of software applications resorting

88  to minimal code writing. Its objective is to empower different kinds of users, by allowing them to

89  easily and quickly create applications: experienced users (e.g., programmers) are able to create

90  software by writing considerably less code, while users without prior experience will require less

91  formal training to start creating applications.

92  Due to non-disclosure agreement conditions, we are not authorized to name the LCDP under

93  study, and for that reason it will henceforth be referred to simply as the LCDP. The LCDP under

94  study allows developers to create both full stack web applications, and mobile applications. It

95  provides a set of predefined templates to bootstrap the development process, which creates the base

96  application. Developers can then expand the application on top of that. The development process

97  itself is performed by resorting to high level development languages, mainly visual languages, similar

98  to Unified Modeling Language (UML) diagrams (Fowler, M., & Kobryn, 2004). The platform also

99  allows developers to graphically edit the interfaces and automatically generate pages and components

100  (e.g., through drag and drop interactions). With this LCDP, it is possible to develop enterprise-grade

101  level applications thanks to the integration mechanisms provided, for instance, with web services,

102  databases or external systems (e.g., SAP).

103  Different languages with different abstraction levels are provided to define different

104  components of the system. The definition of some aspects of the system, such as navigation between

105  screens, the behavior of the screens and buttons, is done through a statechart-like language, as they

106  are adequate for control-flow modeling. These diagrams have a simple syntax, which has the objective

107  of being easily understood by a large audience. Some more complex aspects, such as data retrieval

108  from a database, resort to a Domain Specific Language (DSL), which is more powerful, but

109  simultaneously more complex. The platform also takes advantage of widely known formats, such as

spreadsheets, in order to speed up the development process. This empowers those users who are non-experienced developers, but have had previous contact with these technologies, to more easily understand the platform. Once finished, the applications are converted into standard technologies (familiar web, back-end and mobile languages), and deployed into a cloud environment. The applications become immediately available once published. Examples of this type of platforms include Appian, Google App Maker, Microsoft PowerApps, MIT App Inventor, Nintex Workflow Cloud, OutSystems, Sysdev Kalipso or Zoho Creator.

**Descriptive Cognitive Models**

At the dawn of HCI as an independent discipline, Richard Young wrote that "*for an interactive device to be satisfactory, its intended users must be able to form a 'conceptual model' of the device which can guide their actions and help them interpret its behavior*" (Young, 1981). Since then, it is commonly agreed that knowledge about how users perceive and interact with a computerized environment is of the foremost importance in the design of computer systems that emphasize usefulness and usability (Silva, 2013). The development process of an interactive system greatly benefits from putting the human, the user, in a central position during discussion and design (Dix, Finlay, Abowd, & Beale, 2004; ISO, 2010). In order to better understand how the user conceptualizes and interacts with a system, the discipline of HCI often resorts to models.

Descriptive Cognitive Models (DCMs) are widely used in the study and development of interfaces. Their analytical processes have since long been applied by experts, analysts and developers in order to obtain insight on how the interaction flow, the design features, or the information content of an interface might lead to performance deficits, faulty interactions or use errors (Nielsen, 1994). Although a comprehensive set of DCMs have been developed since the 1980s, it is usually a combination of different models tailored to a specific application case that provides the best result (i.e., insights on how users are thinking about the system and the interaction process). For the purpose of predicting usability issues in LCDP we have selected task decomposition models for a recursive decomposition of our main task into sub-tasks; knowledge-based analysis to comprehend the user's

136 knowledge about the objects and actions involved in a given task; and risk assessment to analyze and

137 evaluate the risk associated with the identified issues.

## Task Decomposition

139     The process of describing the interaction process is often referred to as *Task Analysis* and

140 consists of detailed descriptions and analysis of how people perform their jobs or tasks. It details what

141 they do, what they act on and what they need to know. Identifying the elements and the goals of the

142 task is an essential step to examine the skills necessary to perform a given job. Task decomposition

143 can be performed either in the design phase of a new system or to suggest changes in an existing

144 system.

### *Hierarchical Task Analysis (HTA)*

146     The purpose of HTA is to decompose a task into all its sub-tasks in a way that displays the

147 hierarchical relation between them. It is one of the most predominant examples of a task

148 decomposition methodology. The outputs of HTA are a hierarchy of tasks and sub-tasks, together

149 with plans describing in what order and under what conditions sub-tasks are performed. For examples

150 and further details please see Dix *et al.* (2004) and, for a review on different ways of presenting an

151 HTA and a proposal on an updated notation, see Huddlestone and Stanton (2016).

## Knowledge-based analysis

153     The aim of a Knowledge-based approach to task analysis "*is to understand the knowledge*

154 *needed to perform a task*" (Dix et al., 2004). The main goal of this type of analysis is to build general

155 knowledge taxonomies for each task, after listing all objects and actions. Programming is a

156 knowledge-based activity, and for the purpose of this study we will focus on analyses designed to

157 predict difficulties from interface specifications, namely the External-Internal Task Mapping

158 Analysis.

### *External-Internal Task Mapping Analysis (ETIT)*

160    The ETIT model attempts to deal with the mismatch between the way the user thinks, and the

161    way a system is designed, stating that this mismatch continues until the user learns how to translate

162    what he or she wants to do into the system's terms. ETIT is a contribution of one of the most seminal

163    authors in HCI, Thomas P. Moran (Moran, 1983). In his paper, Moran points out the need for the

164    users to map between the task they are performing and their conceptual model of the machine. Thus,

165    ETIT was conceptualized as a way of assessing the (1) complexity of learning of a naïve user or (2)

166    transfer of knowledge between different systems. In the first case, which is the focus of the current

167    work, ETIT assumes two different spaces: 1) the external task space (i.e., the naïve user's mental

168    model of the task) and 2) the internal task space (the system's commands that allow it to perform the

169    task). The relation between both spaces will be an indicator of the difficulty found in learning how to

170    use the system. According to Moran (1983), when people start using a system, they know they must

171    convert the tasks they have to perform into the system's language and concepts, i.e., they must learn

172    to translate what they want to do into the system's terms. In this model, this translation is represented

173    by *mapping rules*.

174    The ETIT analysis has three parts:

175    1. *An external task space (concepts and tasks described in those concepts)* – whenever a person needs

176    to perform a task, this task is formulated in the external domain/real world, not in the system's terms

177    (Moran, 1983). This means that people formulate their own mental model of the task, using their own

178    known concepts, words and logic;

179    2. *An internal task space (concepts and tasks described in those concepts)* – the system's commands

180    and interaction flow that allows the user to perform the task;

181    3. *A mapping* from the external task space to the internal task space.

182    While the external space is rich and diverse, systems are not. Systems usually abstract a small

183    set of primitive concepts, converting the external task space into smaller internal task spaces.

184    Particularly relevant here is that, while ETIT was conceived as a tool for system design, it can also

185 be used as a competence model of the user, because it makes it explicit what is the knowledge

186 necessary to execute a task.

## Risk Assessment

188     The term "risk assessment" is usually connected to occupational health and safety but in this

189 case, the methodology will be applied to the identification of usability and interaction issues that have

190 the potential to compromise the application under development. The aim of such a process is firstly

191 to identify the issue, and then to mitigate its consequences by adding control measures (Amir-Heidari

192 & Ebrahemzadih, 2015). The steps to be applied in our descriptive cognitive model consist of:

193 ●     *hazard identification* - finding, listing and characterizing issues

194 ●     *risk analysis* - determining the likelihood of the issue

195 ●     *risk evaluation* - comparing an estimated issue against predetermined risk criteria to

196 determine the significance or criticality of the issue

197     After performing a risk assessment, a remedy analysis can be performed where error reduction

198 strategies are defined.

199

## Proposed model - PRECOG: low-code development Platform descRiptivE COGnitive model

202     The descriptive cognitive model we propose was named PRECOG – low-code development

203 Platform descRiptivE COGnitive model (Figure 1). From an analysis of the summarized methods, it

204 became clear that these might interplay to account for an informative Descriptive Cognitive Model.

205 In order to apply the proposed model, the analyst should start by performing an HTA of the particular

206 development use case under analysis. The output of the HTA will provide a list of sub-tasks of the

207 use case that need to be further analyzed. This will be done according to an adapted version of the

208 ETIT. In this adapted version, sub-tasks will be described from the perspective of a naïve user's mental

model, using its own terms and concepts (External task space, henceforth the Knowledge-Based Description - KBD), and from the perspective of an internal task space (henceforth System-Based Description - SBD) which details the steps needed in the LCDP, in order to accomplish the sub-task.

The KBD relies on user's data prior to any interaction with the system. With only a brief description of each task, participants should describe how they would reach each task's final goal using their current knowledge and familiar development tools. This information provides the analysts with the participant's mental model of the tasks in hand, before interaction with the system under evaluation.
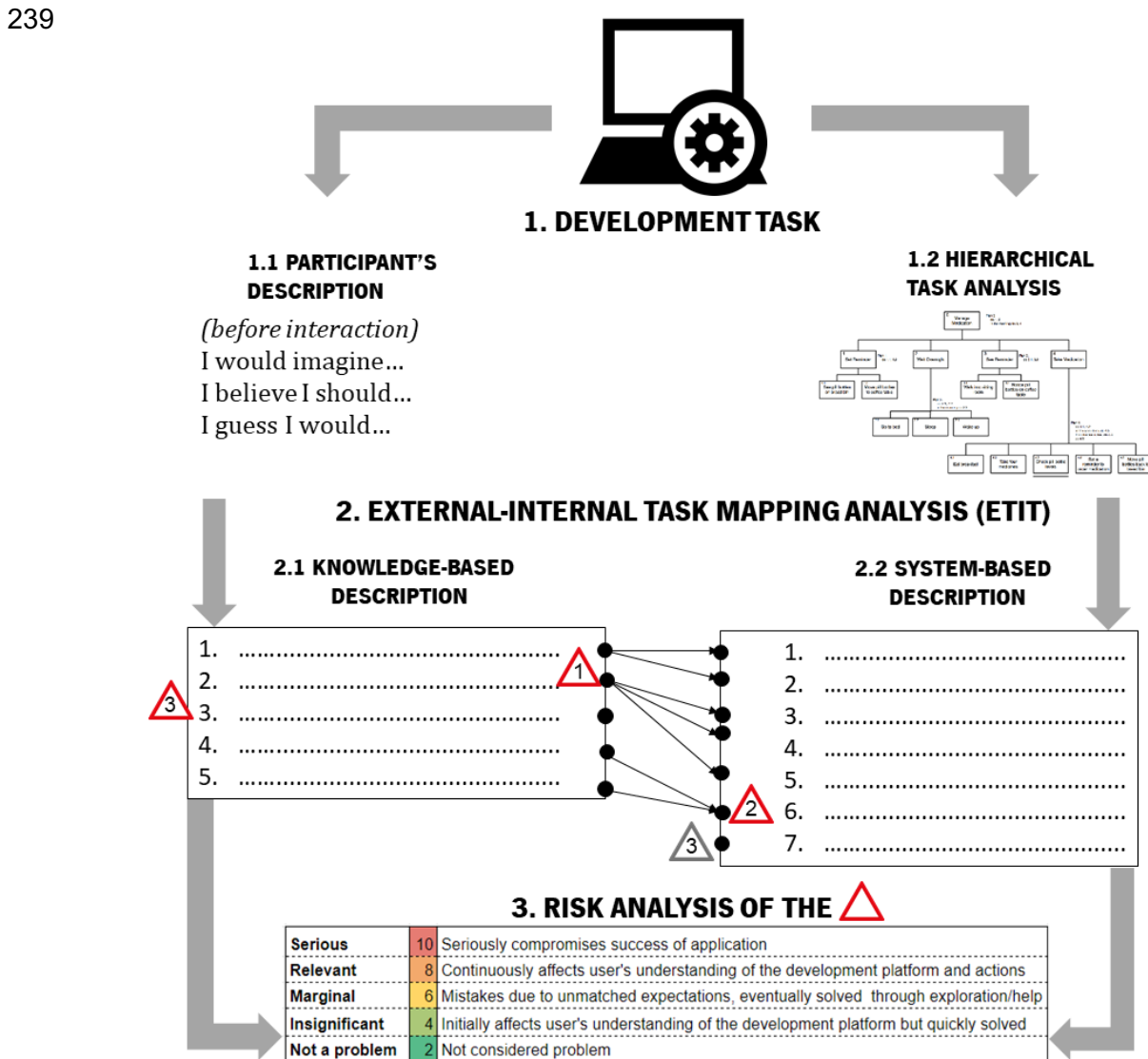
The SBD, on the other hand, is a step-by-step description of the actions performed by an expert user of the system. For the purpose of comparability, the evaluator should define the same decomposition stop-condition for both descriptions. One fitting criterion for a LCDP HTA stop-condition is a discernable user interaction capable of being recorded by the platform (e.g., a drag-and-drop action; the selection of an item from a drop-down menu; the establishment of a new connection in a state-chart; the writing of an expression to define a condition). User interactions at a more atomic level, such as mouse movement, hovering, or typing of a specific character, do not need to be specified in HTAs for LCDP.

Mapping rules should be established between the naive user's mental model (KBD) and the system-based description (SBD), in order to identify possible conflicts. According to our predictions, three types of conflicts might be uncovered by looking at the mapping between the two spaces:

1. *Under decomposition conflict* - occurs when a procedure that is considered by the user as a single step in the knowledge-based description requires multiple steps in the system-based description. This type of conflicts might lead to an underestimation of the sub-task's complexity.

2. *Over decomposition conflict* - occurs when a procedure that is considered by the user as having multiple steps in the knowledge-based description requires only one step in the system-based description. This type of conflicts might lead to an overestimation of the sub-task's complexity and failure to identify and take shortcuts during the development task.

3. *No Correspondence conflict* - occurs when there is no link between a step in the knowledge-based description and one (or several) steps in the system-based description. This might occur because the user is not aware of the appropriate steps required by the LCDP or because the system does not include a feature representative of a mental step the user thinks is needed.



*Figure 1* – *The proposed methodology applied to a given subtask of the HTA. Red triangles signal identified issues that occurred during interaction and the numbers inside correspond to the type of error (1 - Under decomposition; 2 - Over decomposition; 3 - No correspondence). Triangles are colored grey if, after empirical tests with a user, the analyst finds that the predicted error did not occur (false positive).*

245       Once one of these types of conflicts is found, a risk analysis of the conflict should be made.

246 This risk analysis is usually performed in safety critical scenarios, which is not the context of the

247 current study. Nevertheless, for the purpose of providing a richer PRECOG analysis, the risk analysis

248 step can help define priorities in the continuous improvement of the system.

249       Figure 1 summarizes the steps of the overall analysis process. First, a development task is

250 selected (in this case, *create a user interface to list and search for books*) and two branches originate

251 from this task. On the top-right hand side of the figure (Step 1.2), an example HTA is presented for

252 the development task. The HTA will generate the System-Based Description (SBD) or Internal Task

253 Space for the selected task. On the left (Step 1.1), the Knowledge-Based Description (KBD) or

254 External Task Space consists in a description made by people who have not interacted with the low-

255 code system, listing how they would expect to perform the task at hand, knowing what they know at

256 the moment. The second step (Step 2 in the figure) consists in comparing the KBD with the SBD per

257 participant and identifying conflicts that might occur in the interaction using the conflict taxonomy

258 described above ((1) *under decomposition*, (2) *over decomposition*, and (3) *no correspondence*). In

259 this example, four issues arise and are represented by the four triangles. The triangles are red if that

260 issue eventually occurred during interaction, and are marked as grey if they did not occur, being

261 classified as a false positive. *Under decomposition* and *over decomposition* markers are placed inside

262 the box, and *no correspondence* markers are placed on the leftmost border of the boxes.

263       The third and last step in Figure 1 consists of a risk analysis of the identified issues. The risk

264 analysis of a conflict includes:

265       • **Frequency** - An ordinal scale from 0 (never) to 5 (frequent)

266       • **Criticality** - All issues are evaluated regarding the level of criticality as described in Figure

267 2, ranging from 2 (not a problem) to 10 (serious) (Figure 2).

268       • **Pure Risk** - The Pure Risk value of the error is a single value representing the weight that

269 should be attributed to an error, and it results from the intersection between a value of Frequency with

270 a value of Criticality from the matrix of Frequency with Criticality adapted from Amir-Heidari and

271 Ebrahemzadih (2015) (Figure 3). The value ranges from 0 to 50, and the higher the value, the more

272 Critical should the issue be considered, with a higher priority for intervention.

273

| Serious | 10 | Seriously compromises the success of the application |
|---|---|---|
| Relevant | 8 | Continuously affects the user's understanding of the development platform and actions |
| Marginal | 6 | Mistakes due to unmatched expectations, eventually solved through exploration/help |
| Insignificant | 4 | Initially affects the user's understanding of the development platform but quickly solved |
| Not a problem | 2 | Not considered a problem |

274 *Figure 2. Descriptors for each level of Criticality, adapted to the analyzed use-cases*

275



276 *Figure 3. Matrix of Frequency with Criticality adapted from Amir-Heidari and Ebrahemzadih*
277 *(2015). The intersection between a value of Frequency with a value of Criticality provides the Pure*
278 *Risk of the issue, a single value representing the weight that should be attributed to the issue. The*
279 *higher the value, the more critical should the issue be considered.*

280 A model of this sort should be developed for each sub-task deemed relevant for analysis. This

281 will provide relevant information about the participants' difficulties in mapping their conceptual

282 model of the task to the system's operational environment.

283

## Applications of the PRECOG model

PRECOG can have two main applications. It can (1) be used retrospectively to understand the root-cause of an issue identified through user testing, or (2) it can be used predictively to understand which potential issues are going to arise.

In the first case, the model is used to analyze and understand the problems identified during user testing. These problems can be tracked by identifying the correspondence conflicts in the previously made ETIT mapping. This helps to identify potential causes for the problems in terms of mismatches between the KBD and the SBD. The number of observations of the different errors provides additional information that is then used in the risk analysis (as the frequency of occurrence of the errors).

The second application for PRECOG is to use it as a predictive model, which can provide valuable information without the time-consuming process of data gathering with real users. This application takes advantage of the fact that the mapping of the participant's (knowledge-based) description to the platform's requirements provides a first prediction of the effect of the differences between the user's knowledge and system's requirements to achieve a given task. In this case, instead of evaluating the criticality of the issues that effectively occurred, an expert analyst walks through the LCDP and decides on the likelihood of the identified potential problems, evaluating their probability of occurrence given the platform's design. This evaluation should consider, for instance, visual aids and widgets that are available on the platform, and which might be helpful in solving the issue under analysis. This stage of the process refines the predictive power of the model, as it eliminates false positives identified in the mapping stage. Calculating the Pure Risk of all identified errors is then done resorting to an adapted version of the matrix in Figure 4, using Probability (Never, Low, Medium, High) instead of Frequency. The end result will be a list of potential errors organized by Pure Risk evaluation. Remedy analysis of relevant issues (for instance, all Marginal, Relevant and Serious issues would be further detailed) might then be performed, resorting to either expert evaluation or empirical analysis.

**Application and Validation of PRECOG in Empirical user studies**

In this section, we present an application of the PRECOG model where we will detail the analytical process of constructing the HTA, the adapted ETIT analysis to map user knowledge and system-based descriptions, and finally the risk analysis of the identified potential interaction problems. The presented results correspond to the validation of the PRECOG model through empirical user studies. They allow us to understand the suitability and viability of using the selected techniques for the analysis of interaction conflicts in a low-code development platform, including the impact of the identified problems and the analysis of their root-causes.

# METHOD

The first phase of the application of the PRECOG model consisted in defining, with a professional user and LCDP platform developer, a set of representative tasks which could be performed by different types of users. After all tasks were defined, this expert user performed them, and the performance was later analyzed in detail in order to obtain a HTA of each task. The HTA also provided the basis to develop the System-Based Description to be used in the model for each task.

The second phase of the application was performed after empirical user studies with 20 participants. Besides providing usability metrics of performance, these user studies allowed the authors to gather the Knowledge-Based Description of each participant, collected prior to any contact with the LCDP. The user studies complied with the American Psychological Association Code of Ethics, and an informed consent was obtained from each participant.

With both the Knowledge and System-based descriptions, it was possible to do the mapping between both, applying the PRECOG to each participant, and listing all the potential issues and mistakes that could happen during task execution.

The final phase of the validation effort was carried out after thorough video analysis of each participant's performance and comparison between the model's prediction and the real outcome in the user studies. Having identified all observed issues, risk analysis was applied to understand the issues' frequency and criticality.

## Participants

A total of 20 participants were recruited (Table 1). The recruited population respected the following requirements:

- 10 participants had a software-engineering background (formal education in the past or present) - denoted as **Experts**;
- 10 participants had education in social sciences, economics or finance areas - denoted as **Novices;**
- All participants were over 18 years old, proficient in English, unfamiliar with the LCDP (had never worked with it), and willing to accept the sessions to be recorded (screen and audio).

The recruitment was made via internal mailing lists at the author's institutions as well as personal contacts. Of the 20 participants, 7 were female and 13 were male.

***Table 1 -*** *Characterization of the participants*

|  |  | Novices | Experts |
|---|---|---|---|
| Gender | Female | 6 | 1 |
|  | Male | 4 | 9 |
| Degree |  | Psychology, Economics, Biochemistry, Management, Acoustics | Software Engineering |

The recruitment phase included a questionnaire to understand the participants' experience with programming languages. This questionnaire was custom-made, inspired by the knowledge acquired by students during an informatics degree, and was divided into three sections, each with increasing complexity in terms of computer science skills. The first section tested if the participant was familiarized with EUD tools (specifically, Spreadsheets editing software), and basic computational concepts, such as the concept of formula. The second section tested the capability of the user to understand simple software development concepts, such as interpreting and writing software, elementary data structures (e.g. arrays and binary trees), and query languages. The third level tested if the user had advanced software development skills, such as communication protocols, object

359 oriented concepts and software modelling. Participants should also indicate which programming

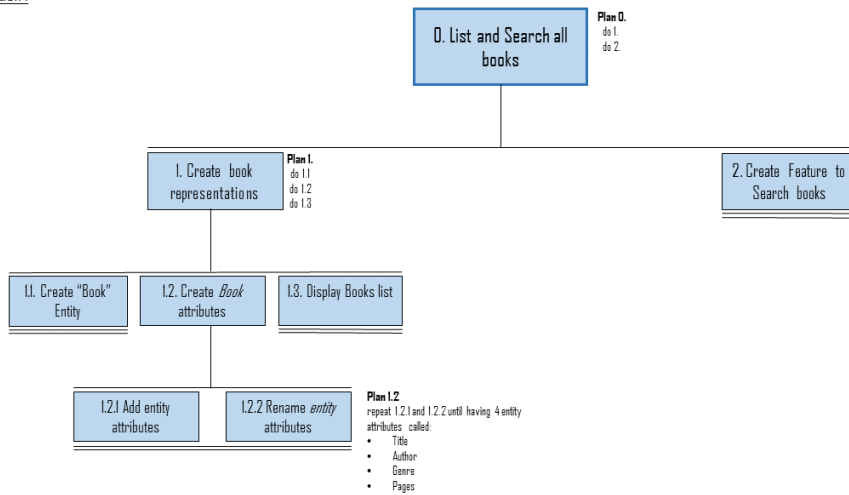360 languages and integrated development environment (IDE), if any, they felt comfortable using.

361 **Defining the System-based description**

362 The participants' use-case consisted in creating a web application to manage books. The

363 development of the application ("My Books") was divided into five tasks, which could be performed

364 in any order, each aimed at fulfilling one of the following requirements (in order of increasing

365 difficulty):

366     1. The user of the "My Books" application should be able to list and search all books;

367     2. The user should be able to see and edit the details of a book;

368     3. The user should be able to register new books in the application;

369     4. The application should present the user with a homepage with two buttons:

370         a. One that redirects users to the list of books;

371         b. Another that goes to the screen that allows registering a new book;

372     5. When seeing the details of a book, the user should see a list of other books from the

373        same author.

374 We started by developing an HTA for the overall use case, divided into 5 HTA sub-diagrams

375 corresponding to the five different tasks of the use case. Figure 4 presents the HTA diagram for the

376 first task, following the graphical notation of Marshall *et al.* (2003), where the main task description

377 is at the root of the diagram ("0. List and Search all books") and the different sub-tasks at lower levels

378 (e.g., "1. Create book representations").

379



380        *Figure 4 - Example of one HTA diagram for the use case task "List and Search all books".*

381        *When a (sub)task is further decomposed, a plan describing how its subtasks can be combined is*

382        *detailed*

383        After defining the HTA for each sub-task, we defined the level of the HTA tree that better

384        represents distinguishable tasks in the interface and we list all the units of interaction that are required

385        to perform in the LCDP in order to complete that specific sub-task. For instance, to create a list of

386        books with a search function the user of this particular LCDP as to perform the following actions:

387           1.  Go to Data Menu;

388           2.  Select Database;

389           3.  Add Entity named "Book";

390           4.  Add Attributes;

391           5.  Rename Attributes;

392           6.  Add Book Entity to Home Screen representation;

393           7.  Boilerplate generation of search functionality

394

395

**Defining the Knowledge-based Description**

In order to obtain a Knowledge-based Description, each participant was instructed to verbally describe how he or she would perform each task, both in terms of the interface and in terms of the back-end development. Participants did this with the conceptual knowledge they had and using the tools they knew (if any), they should describe *how they would* complete the tasks. This was requested before the participant interacted with the LCDP.

**Material**

The tests were performed in quiet testing rooms. Locations were equipped with a table and three chairs, a laptop computer for the participant, with the LCDP running, ActivePresenter 7 screen and audio capture software (Atomi Systems, 2019), and a video camera (the participant's facial expressions were not captured at any time, the camera focused on the screen of the laptop as a backup measure).

**Procedure**
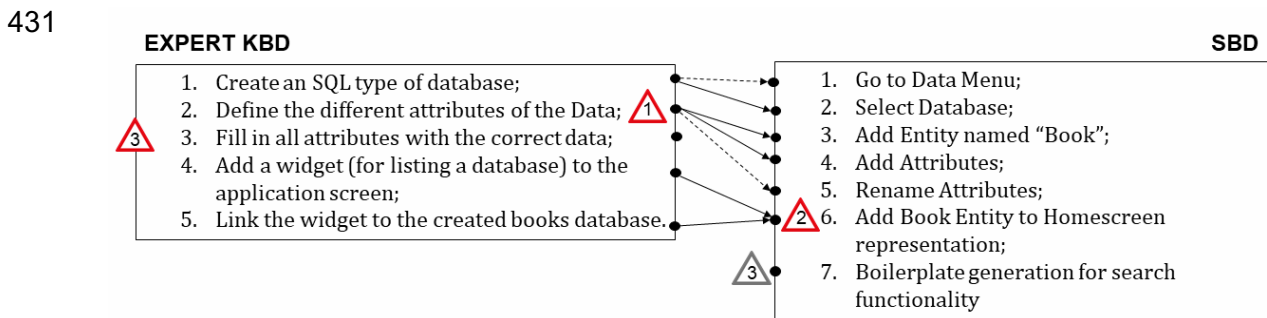
The participant was welcomed by the test moderator and the data logger, who explained that an evaluation was being carried out on how hard or how easy it was to use a particular LCDP. He or she read and signed the informed consent where more detailed information was provided. Each participant was presented with the five tasks. First, the participant was instructed to describe verbally how he or she would perform each task, both in terms of the interface and in terms of back-end development. Regarding this process, it became evident that the test moderator played a role in the success of this phase, which would be used to build the Knowledge-based Description of each participant. The moderator should prompt the participant when he or she becomes quiet, asking specifically about aspects of the application in order to gather as much information as possible.

Then, the participant was requested to carry out the tasks upon performing a tutorial. The tutorial consisted of an interactive session, where the users were put in contact with the LCDP and explained the basic concepts. The test moderator informed the participants that there were several
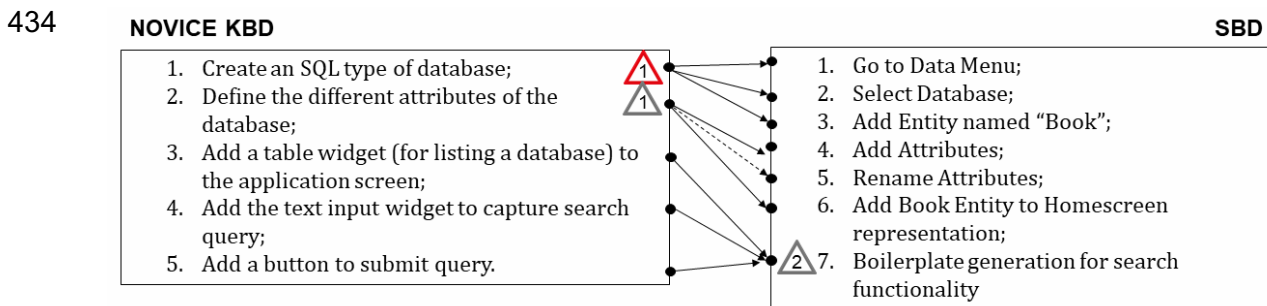
421    ways of concluding the tasks, and that they could search the internet for answers. Each participant

422    was provided with the same instructions. The participants were given a written copy of the

423    instructions and respective memory aides.

424    **Analysis**

425    For each subtask in the Hierarchical Task Analysis deemed relevant, we developed adapted

426    ETIT mappings from the collected Knowledge-based Description to what was the required plan of

427    action in the platform (the System-based Description). Figures 5 and 6 show as an example the ETIT

428    mapping for the sub-task "List and Search all Books", for an expert and a novice user, respectively.

429    The end-result of the mapping exercise allows the analyst to identify what type of use-errors might

430    occur during that sub-task.

431



432    *Figure 5 - Example of the adapted ETIT mapping for the sub-task "List and Search all Books" in*

433    *an Expert user. Dashed arrows correspond to implicit steps.*

434



435    *Figure 6 - Example of the adapted ETIT mapping for the sub-task "List and Search all Books" in a*

436    *Novice user.*

437     As we pointed out earlier, analysts can uncover three types of conflicts by looking at the

438 mapping rules. Figure 5 illustrates all three types of conflicts - *under decomposition* in step 2 of the

439 KBD, *no correspondence* in step 3, and over *decomposition* in step 6 of the SBD.

**440 Identify the root-cause of issues found in empirical tests with real users**

441     The analysis depicted in Figure 1, which shows the complete flow from the definition of the

442 models, through the identification of a KDB-SBD conflict, to risk analysis of that conflict, was the

443 one followed in the present work.

444     Each interaction video of the participants was thoroughly observed, the issues that occurred

445 were identified and their relevance analyzed using the following steps:

446 -   **Root Cause** – After analysis of the interaction that resulted in an error, a root cause

447     was identified.

448 -   **Remedy Analysis** – A recommendation for a way to avoid the error or issue was

449     devised.

450 -   **Evidence** – When available, evidence gathered during data collection with

451     participants (video) was also registered.

452 -   **Pure Risk** – The frequency and criticality of the issues was evaluated by the authors

453     and the LCDP's research and development team.

454     One complete analysis took on average 40 minutes per participant for an experienced analyst.

455     The Frequency of an issue was defined as the frequency with which the issue under analysis

456 was observed in the empirical tests. The scale can be adjusted depending on the size of the sample,

457 without the need to modify our model.

458     To determine Criticality, all observed issues were evaluated by four authors and four

459 professional LCDP developers regarding the level of criticality as described in Figure 2, ranging from

460 2 (not a problem) to 10 (serious). The evaluations were performed individually considering the

461 general criticality of the issue, and not the particular context/participant where it happened. Inter-rater

462 reliability was assessed using a two-way, average measures ICC (intraclass correlation). The resulting

463 ICC was in the "fair" range, ICC=0.50 (Cicchetti, 1994), indicating that evaluators had a fair degree

464 of agreement. The ICC increased to 0.68, in the "good" range, when only one group (authors) was

465 considered. This indicates different evaluation criteria from both groups, something which would be

466 worth exploring in the future. The final Criticality value was the mode of the eight evaluations.

# RESULTS

468     In this section we will present results on the comparison between the PRECOG's predictions

469 and the interaction issues observed during the empirical usability tests. Moreover, we will also address

470 the nature of the use-errors that were predicted and verified, in terms of its root-cause, frequency, and

471 pure-risk.

472     The comparison between PRECOG´s predictions and results of the empirical usability tests

473 will be presented regarding the two profiles (expert and novice) and the three types of conflicts

474 signalized by PRECOG (1 - Under decomposition; 2 - Over decomposition; 3 - No correspondence).

475 In Table 2 we can see that out of a total of 135 potential interaction issues identified by PRECOG, 67

476 (49.6%) occurred during empirical usability tests. Moreover, of all interaction issues verified in the

477 empirical usability tests, only 5 were not predicted by a type of conflict signalized in the PRECOG

478 model. Table 2 summarizes the outcome of applying PRECOG in three confusion matrices,

479 considering the studied profiles both combined and separately.

480 ***Table 2 -*** *Confusion Matrices for a combination of all participants and by participants' profile*

481 *(Novices and Experts).*

| | | All Participants | | |
|---|---|---|---|---|
| | | *Verified* | | |
| | | *Use Error* | *No Use Error* | Total |
| *Predicted* | *Use Error* | 67 | 68 | 135 |
| | *No Use Error* | 5 | 64 | 69 |
| | | | TOTAL | 204 |

482

| EXPERTS | | | | |
|---|---|---|---|---|
| | | *Verified* | | |
| | | *Use Error* | *No Use Error* | Total |
| *Predicted* | *Use Error* | 39 | 42 | 81 |
| | *No Use Error* | 2 | 50 | 52 |
| | | | TOTAL | 133 |

483

| NOVICES | | | | |
|---|---|---|---|---|
| | | *Verified* | | |
| | | *Use Error* | *No Use Error* | Total |
| *Predicted* | *Use Error* | 28 | 26 | 54 |
| | *No Use Error* | 3 | 14 | 17 |
| | | | TOTAL | 71 |

484

485    In order to assess the predictive capability of our DCM we analyzed the values presented in

486    Table 2, where it is possible to see the total number of *true-positives* (i.e., predicted and confirmed

487    use-error), *false-positives* (i.e., predicted but unconfirmed use-error), *true-negatives* (i.e., predicted

488    and confirmed inexistence of use-error), and *false-negatives* (i.e., not predicted but confirmed

489    existence of use-error). An efficient predictive model aims at scoring high in both true-positives and

490    true-negatives and low in both false-positives and false-negatives. From a confusion matrix one can

491    calculate several complementary values to assess a classifier's predictive capability (Powers, 2011;

492    Tharwat, 2018), namely:

493    • *Sensitivity* - or *recall* is the proportion of the positive samples (i.e., verified use-errors) that

494        were correctly classified as so. Thus, Sensitivity depends on true-positives (TP) and false-

495        negatives (FN), which are in the same column of the confusion matrix, and can be calculated

496        as: Sens = TP/(TP+FN)

497    • *Specificity* - or *inverse recall* is the proportion of negative samples (i.e., verified no use-

498        error) that were correctly classified as so. Thus, Specificity depends on true-negatives (TN)

499    and false-positives (FP), which are in the same column of the confusion matrix, and can be

500    calculated as: Spe = TN/(TN+FP)

501    • *Accuracy* - is defined as a ratio between the correctly classified samples to the total number

502    of samples, and can be calculated as follows: Acc = (TP+TN)/(TP+TN+FP+FN)

503    • *F1-score* - also called F1-measure is the harmonic mean of sensitivity and positive

504    predictive value (ppv = TP/(TP+FP)). The value of the F1-score ranges from zero to one,

505    and high values indicate high classification performance. F1-scores are calculated as follow:

506    F1 = (2TP)/(2TP+FP+FN)

507    • *Informedness* - also called *Youden's index* quantifies how informed a predictor is for the

508    specified condition, and specifies the probability that a prediction is informed in relation to

509    the condition (versus chance) (Powers, 2011). The value of the Informedness ranges from

510    zero, or chance-level, to one, representing a perfect predictive capability. Informedness can

511    be calculated as follow: Inf = *sensitivity + specificity* - 1

512    Table 3 shows the PRECOG's values obtained for these different variables, based on the

513    confusion matrices presented in Table 2.

514    **Table 3 -** *PRECOG's predictive metrics for All Participants and divided by user profile*

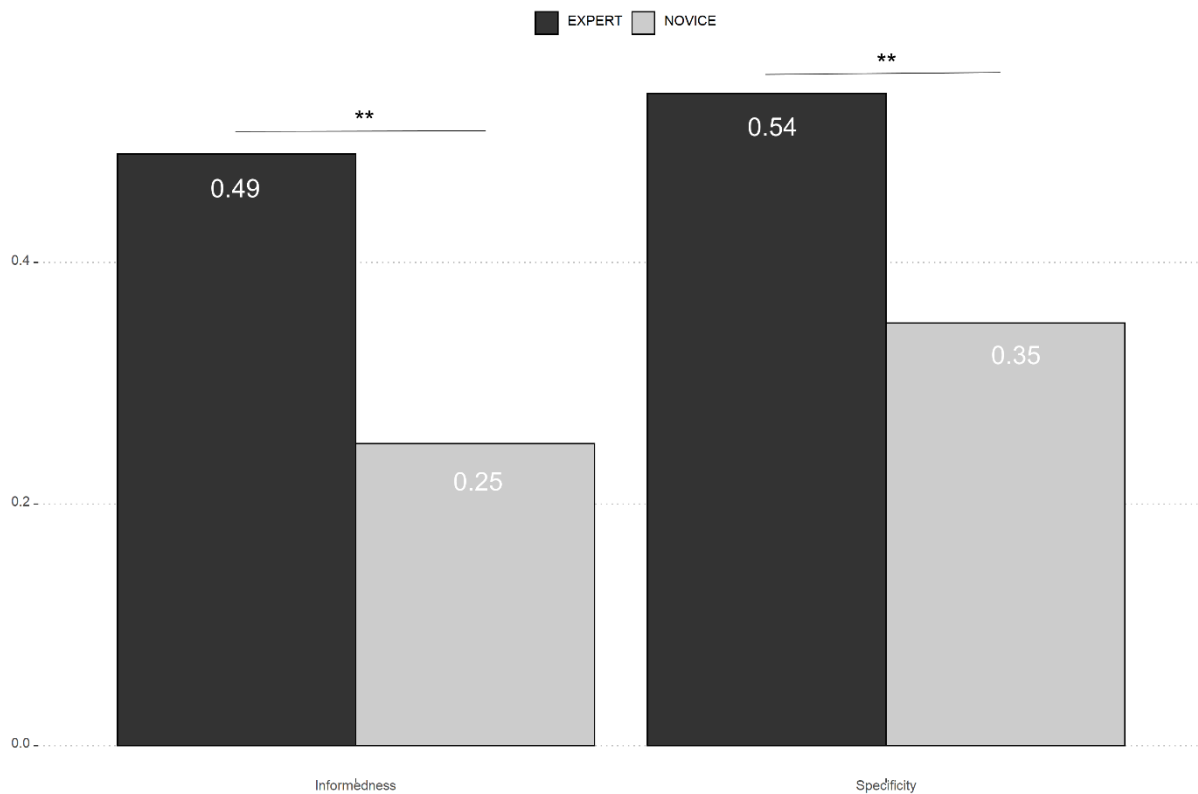|  | *All Participants* | *EXPERTS* | *NOVICES* |
|---|---|---|---|
| Sensitivity | 0.93 | 0.95 | 0.90 |
| Specificity | 0.48 | 0.54 | 0.35 |
| Accuracy | 0.64 | 0.67 | 0.59 |
| $F_1$ Score | 0.65 | 0.64 | 0.66 |
| Informedness | 0.42 | 0.49 | 0.25 |

515

516    Both Powers (2011) and Tharwat (2018) discussed the advantages and limitations of each of

517    these metrics for classification performance. According to Powers, Sensitivity and F1-scores ignore

518    performance in correctly handling negative examples, propagate underlying marginal prevalence and

519    biases, and fail to account for the change level performance. Nevertheless, Tharwat makes the case

520    that all these different metrics, being more focused (such as Sensitivity, Specificity, and F1-score) or

more general (such as Accuracy and Informedness) are useful to understand all the potentialities of a particular classifier. In the case of PRECOG, Sensitivity was generally higher than Specificity and, while Accuracy and F1-scores were fairly similar for both Experts and Novices, Informedness was the measure that varied the most regarding type of participant.

Statistical tests revealed differences between Experts and Novices concerning Specificity and Informedness. An unpaired two-samples Wilcoxon test indicated that Specificity in Experts was significantly higher than in Novices ($W = 73.5$, $p < .01$, $r = -0.59$). Similarly Informedness ($W = 81$, $p < .01$, $r = -0.72$) was also significantly higher in Experts than in Novices.
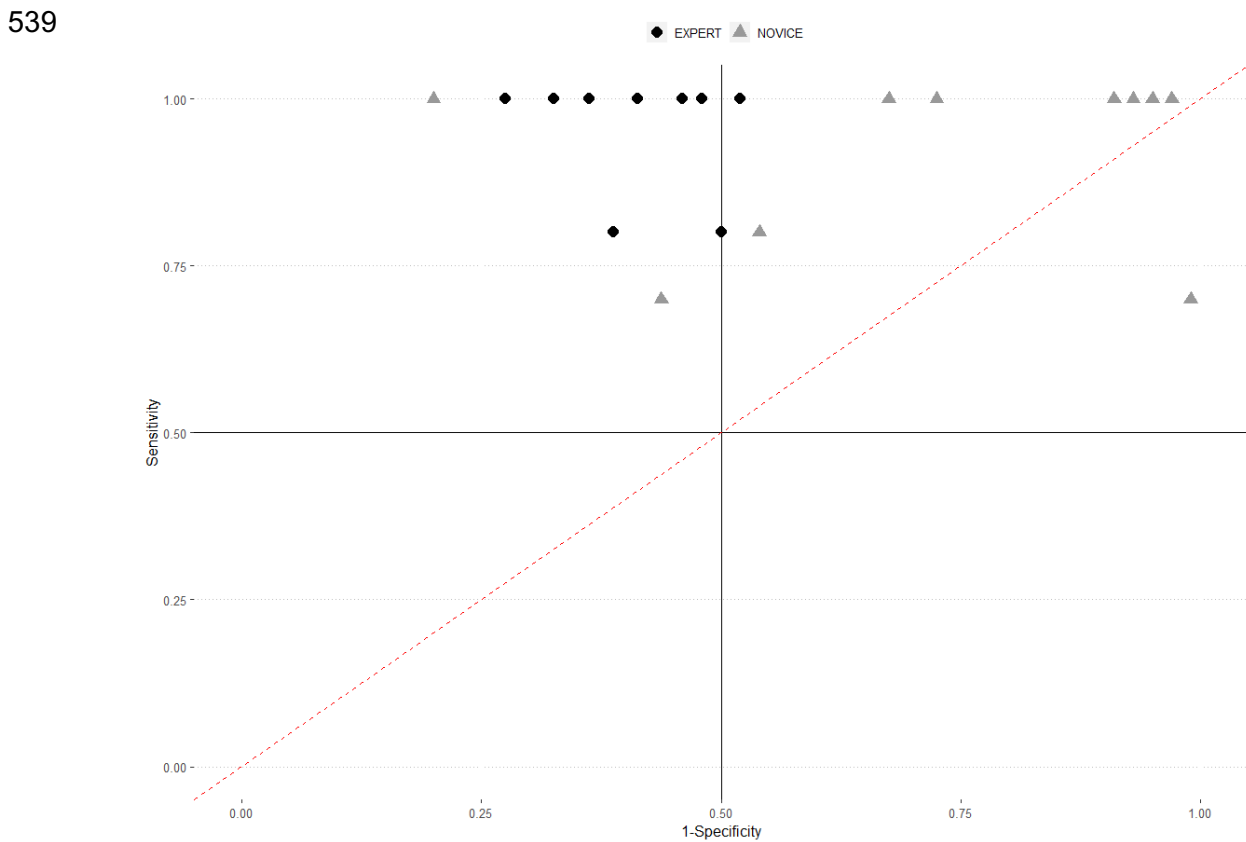


*Figure 7 - Differences PRECOG's classification of Experts and Novices concerning Specificity and Informedness.*

Having calculated the Sensitivity and Specificity of each participant's classification, it was possible to map the performance of PRECOG in a Receiver Operating Characteristic (ROC) plot (Figure 8). Participants are mainly mapped in the upper left-hand of the ROC space, meaning that, for the generality of the participants, PRECOG's classification was predictive of actual behavior.
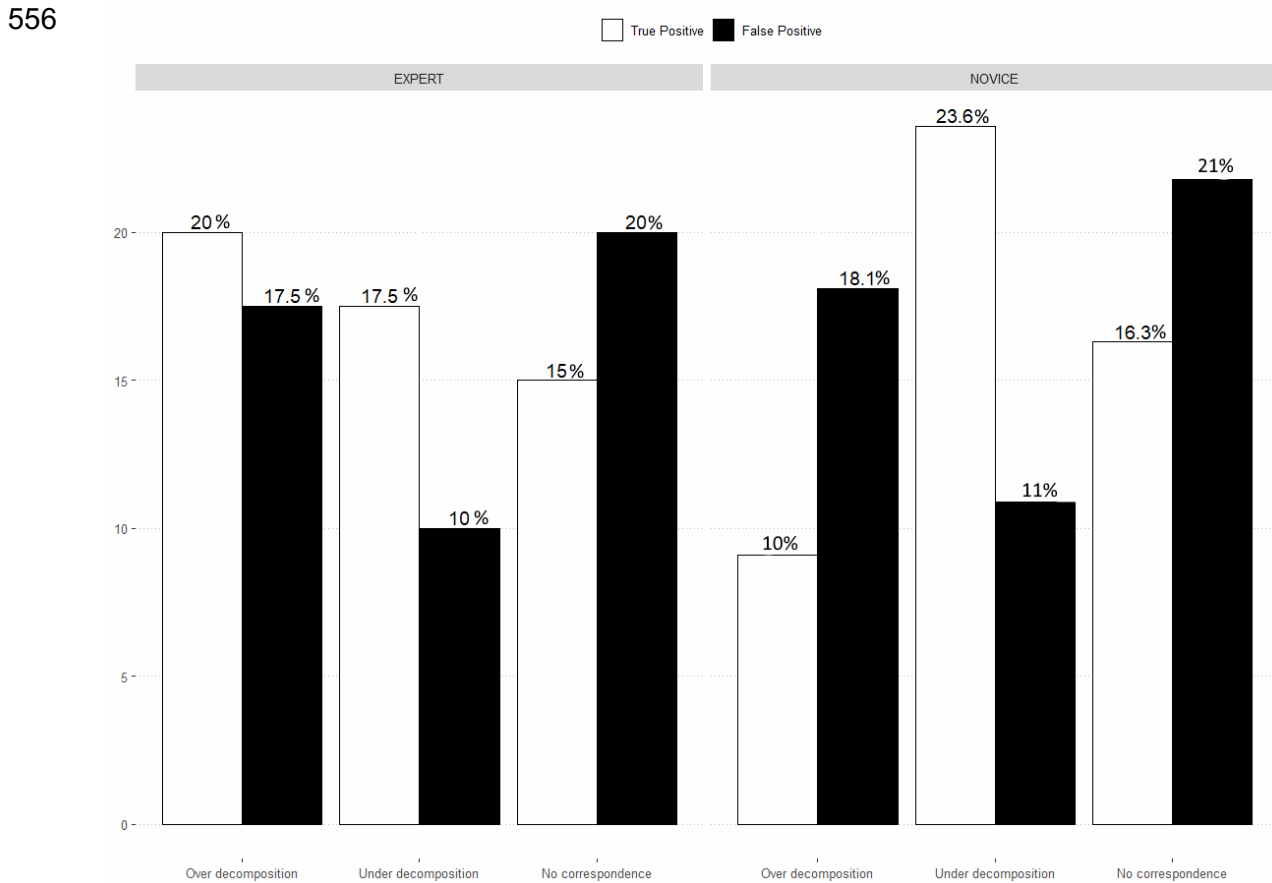
24

537    Again, it is possible to verify that data points of Expert participants are further from chance level

538    (x=y) than data points of Novice participants.

539



540

541    ***Figure 8 -*** *Receiver Operating Characteristic (ROC) data point cloud representing Experts (points)*

542    *and Novices (triangles). The diagonal line represents the chance level.*

543

544         Looking at the distribution of True Positives and False Positives according to profile and

545    considering the three types of conflicts predicted by PRECOG, it is possible to observe that Under

546    decomposition conflicts were the type of mapping conflict where PRECOG performed better.

547    PRECOG obtained the highest difference between True Positives and False Positives in this type of

548    conflict, having Under decomposition conflicts in Experts accounting for 17,5% of the overall correct

549    predictions and in Novices accounting for 23.6% of the overall correct predictions. In the case of

550    Experts, Over decomposition got the highest rate of True Positives (20%), however, this type of

551    conflict also accounted for 17.5% of False Positives. For the Over decomposition conflicts in Novices,

552    the reverse pattern was observed, with a higher number of False Positives (18.1%) than of True

553 Positives (10%). Finally, the No correspondence type of conflicts were the type of mapping conflict

554 where PRECOG had a worse performance, with a higher number of False Positives than True

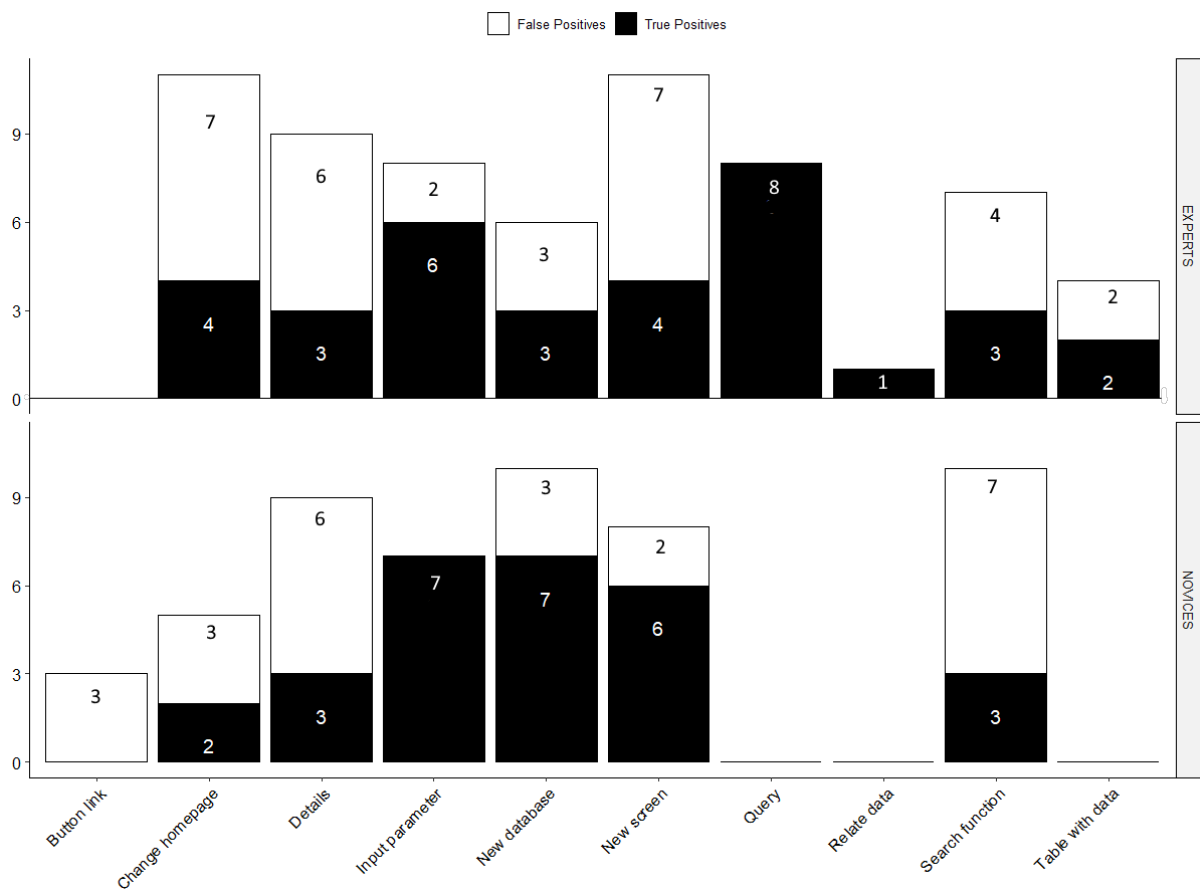555 Positives both in Experts (20%) and in Novices (21%).

556



557 **Figure 9 -** *Distribution of True Positives and False Positives according to profile, considering the*

558 *three types of conflicts predicted by PRECOG*

559 Focusing on the data coming from the usability empirical studies, and considering the

560 Frequency of the observed interaction issues, we were able to identify 10 types of issues: Input

561 parameter; New database; Table with data; Button link; New screen; Query; Search function;  Relate

562 data; Change homepage; Details. Figure 10 depicts the Frequency with which each of these issues

563 occurred in both profiles. It is possible to observe that the same issues did not occur for both profiles.

564 As a first characterization, it can be observed that Experts had a more diverse typology of issues,

565 having experienced nine while Novices had six different types of issues. No participant had any

566 "Button Link" related issues. On the Novice side, they did not have issues related to "Table with

567  data", "Query" and "Relate data".  For the Experts, the issue with the highest True Positives was

568  "Query" with 8 correctly identified issues. The highest number of False Positives among Experts

569  were found in the "New Screen" and "Change homepage" issues, where 7 identified issues were

570  considered False Positives. Regarding Novices, "Input Parameter" and "New database" issues had 7

571  True Positives followed by "New Screen" with 6 True Positives. Concerning False Positives, "Search

572  function" with 7 False positives was followed by "Details" with 6.

573



574  *Figure 10 - Frequency of True and False positive issues identified in the empirical tests after KBD-*

575  *SBD mapping*

576          Besides Frequency, another component for the Pure Risk analysis is the Criticality of the issues.

577  Table 4 presents the final Criticality attributed to the issues that emerged during our analysis.

578  **Table 4 -** *Criticality of the identified issues by a pool of eight evaluators*

| Issue | Criticality | Description |
|---|---|---|
| Input parameter | 8 | Relevant |
| New database | 6 | Marginal |
| Table with data | 8 | Relevant |
| Button link | 8 | Relevant |
| New screen | 8 | Relevant |
| Query | 8 | Relevant |
| Search function | 8 | Relevant |
| Relate data | 8 | Relevant |
| Change homepage | 6 | Marginal |
| Details | 8 | Relevant |

579

580      Figure 11 depicts all ten different issues observed during the interaction with the LCDP. The

581 issues are presented in terms of the Pure Risk evaluation (min=0, max=50). Besides Criticality, the

582 Pure Risk evaluation considers the Frequency of occurrence of each issue, hence the difference in the

583 value of the same issue for both profiles.

584



585 ***Figure 11 -*** *Pure Risk of each observed issue per profile. Issues found by the LCDP Descriptive*

586 *Cognitive Model according to their Pure Risk evaluation (min=0, max= 50)*

587    This data indicates that, for instance, the "Query" issue should be solved, as almost all Experts

588    who interacted with the LCDP had issues with creating a query. It also indicates that the weight of

589    each issue is different depending on the LCDP user. For Novices, the highest Pure Risk value is in

590    the "Input Parameter", followed by "New Database" and "New Screen" issues. The only Pure Risk

591    value similar to both profiles occurred in the "Search Function" issue, indicating similar difficulties

592    in this task for both profiles.

593    # DISCUSSION

594    We have successfully applied our descriptive cognitive model (PRECOG) to a relevant use

595    case, which allowed us to validate the viability of the proposed approach.

596    **Predictive power**

597    Although time-consuming, the methodology proved to correctly identify several high-

598    criticality issues from both user profiles. Indeed, we were able to use the model to predict a relevant

599    number of problems, prior to the user studies that confirmed them. The confirmation came with a

600    relatively high number of False Positives. Such is due to the conservative nature of the model, as

601    highlighted by the fact that results show a high Sensitivity, with a comparatively lower Specificity.

602    We opted for an exhaustive approach, and no relevant error was disregarded in the KBD-SBD

603    mapping analysis, since such a stance allowed a more extensive list of expected problems. Despite

604    this conservative approach, the majority of the predicted issues or root causes occurred during

605    empirical user tests (all were observed except for the "Button Link" issue).  This is a relevant output

606    as very little unpredicted issues occurred during the user studies and these were mostly related to

607    navigational (interaction with the platform) difficulties and not conceptual mismatches.

608    Another explanation for the high number of False Positives predicted by PRECOG concerns

609    "No correspondence" issues, where one item from the KBD or the SBD finds no correspondent on

610    the other side. Some issues of this type may have arisen because the participant did not mention a

611    certain development step or minor detail - which ended up being evident when interacting with the

612    LCDP, thus not resulting in a real issue. This is particularly evident in the case of Novice

613     programmers, leading to a significantly lower Specificity in their case, and can be attributed to the

614     fact that their mental models of the programming tasks differ from the platform's model more than

615     the Experts' mental model. That some of the problems were not observed is, in effect, a positive

616     indicator towards the ability of the LCDP to guide users during the learning stage. Conversely, "Over

617     decomposition" was the type of issue with more false positives, which seems to indicate a need to

618     raise the abstraction level of the platform. Doing this without compromising the ability of advanced

619     users to fine-tune applications that are more complex, implies exploring strategies for adaptive user

620     support (Oppermann, 1994; Gajos, Czerwinski, Tan & Weld, 2006). An important observation should

621     be made regarding moderation of the Knowledge-Based Description stage. The moderator should be

622     very familiar with the platform and the tasks under study in order to prompt the necessary information

623     to complete the mapping. It is extremely important to obtain all the information that will allow

624     illustrating the participant's mental model before interaction with the platform.

625     Another significant result concerns the Criticality of the issues observed. All issues which

626     occurred had a Criticality evaluation of more than 6, that is, were either Marginal (*Mistakes due to*

627     *unmatched expectations, eventually solved through exploration/help*), or Relevant (*Continuously*

628     *affects user's understanding of the development platform and actions*). This is an indicator that our

629     approach is useful in detecting high-criticality issues. Unfortunately, and this is a limitation of the

630     current study, no Criticality evaluation was performed in the non-observed issues in order to compare

631     the criticality between the observed and unobserved issues.

632     Regarding the number and type of issues, the applied model allowed us to distinguish two

633     patterns between Novice and Expert programmers. Expert programmers had more types of issues (9

634     in total), but each happening less frequently, depicting a more exploratory behaviour, Novice

635     programmers had less types of issues (6 in total) but each with more repetitions, meaning that the

636     issues were effectively problematic for this profile, which explored less and whose participants had

637     very similar performances. Another potential explanation for this difference in number of observed

638     issues results from the speed with which each profile performed the tasks. Expert programmers

639     managed to complete more tasks than Novice programmers, hence there was more opportunity for

640     issues to arise. These different patterns affected the Pure Risk evaluation of the issues according to

641 profiles. Even though the Criticality was the same for all issues independently of the profile, the

642 Frequency of the issue affected the Pure Risk value. The issues found in Novice programmers were

643 consistent and robust, and this is even more interesting if we consider their heterogeneous background

644 education. The Pure Risk evaluation should work as an order of priorities in order to improve the

645 LCDP under study. Indeed, the LCDP development team has since addressed some of the issues

646 found.

647 In global terms, the combined results show PRECOG is somewhat a conservative model, which

648 despite eliciting some False Positives also identified correctly the majority of issues that effectively

649 impaired the participants' progression. It should be noted that we are considering first-time users of

650 the platform without any formal training.

651 As it is the case for safety-critical systems, when analyzing applications such as LCDPs, it is

652 safer to overestimate the occurrence of potential for use-error, than fail to identify serious use-errors

653 that occurred. In this regard, PRECOG appears to be a promising approach in the sense that it was

654 able to identify almost all issues faced by Expert and Novices programmers.

655 From the application of this LCDP Descriptive Cognitive Model, and as highlighted by the

656 ROC analysis, we conclude that the method has predictive power (i.e., most of the identified

657 knowledge-system conflicts were resulting in use issues during the performance of the task), and that

658 this methodology could be used as an effective tool to predict, understand, and mitigate use errors

659 and faulty interactions in an LCDP platform.

660 **Model applicability**

661 Regarding the applicability of PRECOG, the descriptive cognitive model itself was tailored to

662 a specific set of applications, and for the LCDP it proved useful, granular and precise. In theory, the

663 model is not limited to LCDP platforms. We address applications where the user tasks can be

664 decomposed through an HTA, so that the KBD-SBD mapping is possible. That said, we are only able

665 to support our claims in the low-code development area, due to the performed user studies. As for the

666 type of participants, the model assumes naïve participants or first-time users, but maybe the observed

667 granularity and detail of the predictions would allow its application with participants with some

668 knowledge of the platform under study (i.e., testing tasks that the participants never performed in that

669 particular platform). This would be an interesting application for the future, as we expect the model

670 to identify the missing knowledge independently of the proficiency of the participant.

671 **Threats to validity**

672 Having performed a user study, results are susceptible to threats to its validity. Specifically, we

673 acknowledge the relatively limited number of participants (a total of 20). While the achieved results

674 regarding the percentage of effective issues vs. identified ones are positive, the impact of a larger

675 study group should be analyzed.

676 Would this model work without empirical studies? We believe it would, and in the section

677 where PRECOG is presented, we provide the necessary steps to do so. We believe its predictive

678 capability would be improved if the criticality evaluation had been performed after the KBD-SBD

679 mapping for all identified potential issues (including false positives). Having all the criticalities of all

680 potential issues would, for instance, allow the analysts to choose the issues evaluated with a criticality

681 of 6 or more. These issues, according to our results, have a higher probability of being captured by

682 PRECOG. In the future, it is our intention to validate this second approach by performing cognitive

683 walk-throughs in the platform and performing criticality evaluations in replacement of the empirical

684 studies with real participants.

685 **Value**

686 The approach used in the current study allowed us to identify problems, difficulties and issues

687 participants faced during the interaction with the LCDP. Although not detailed in the present study,

688 a root-cause analysis of each issue allowed us to understand that these arose mostly due to two types

689 of lacking concepts: LCDP-related concepts and development-related concepts. Whereas Expert

690 programmers knew the development concepts but had difficulty in translating them into the LCDP

691 terms, Novice programmers lacked basic development-related concepts, which largely affected their

692 performance.

693 Although these results are based on the study of one development use-case, and we cannot

694 generalize to the entire LCDP, PRECOG allowed the identification, prioritization and root-cause

695 analysis of several issues. This is valuable information for the LCDP platform developers as they

696 have as main goal to place both types of profiles (Experts and Novices) within the Optimal Flow (cf.

697     Csikszentmihalyi, 1990; Repenning & Ioannidou, 2006), with just the right amount of challenges and

698     just the right amount of skill-acquisition - each at their own pace. The nature of the issues also

699     provides valuable inputs to support adjusting the LCDP learning process, according to the Optimal

700     Flow. Specifically, the nature of the errors should be taken into account: Novice users, due to the lack

701     of software development skills, will fall into anxiety, as they are not able to develop the desired

702     features; Expert users, while lacking knowledge of the platform, will perform the tasks resorting to

703     the previously acquired knowledge, which might result in repetitive and monotonous tasks, leading

704     to boredom.

## CONCLUSION

706     Low-code development platforms have the potential to dramatically change how software is

707     developed, making it possible, at least for particular domains, for someone without a formal education

708     in computer science to develop quality software, and for experienced developers to significantly speed

709     up the development process. Understanding how programmers and non-programmers approach this

710     type of platform, is key to support their design and evolution. By developing and applying PRECOG,

711     a new Descriptive Cognitive Model (DCM), aimed at identifying interaction issues with the learning

712     of low-code platforms, we were able to gain insights into potential problems with a specific low-code

713     platform's use. The proposed DCM was validated, using empirical techniques. Twenty participants

714     were observed interacting with the LCDP, of which 10 were expert programmers and the other 10

715     were novice programmers. All performed the same tasks and all interactions were analyzed according

716     to the proposed model.

717     Although a high number of False Positives were identified after a first mapping between the

718     user's mental model and the system's requirements, it is relevant to notice that all issues but one

719     (Button link), which occurred during users' interaction with the LCDP, were predicted by this

720     mapping. Expert programmers had a higher number of observed issues, although each occurring less

721     frequently. This was due to expert programmers performing the tasks more quickly and with a more

722     explorative behaviour, giving room for more issues to occur. On the other hand, Novice programmers

723     faced fewer issues, although each occurred more frequently. These results allowed us to successfully

724 identify high criticality use errors through the analysis of the users' mental model and, importantly,

725 the results allowed us to identify the root causes of each issue. One of the future goals of the current

726 research is to validate PRECOG as a predictive model without recurring to user studies.

727 PRECOG revealed itself quite valuable in the search for more usable LCDP and effective EUD

728 platforms. As Maeda (2006) points out, "*observing what fails to make sense to the non-expert, and*

729 *then following that trail successively to the very end of the knowledge chain is the critical path to*

730 *success* [i.e., in developing simple and easy to learn systems]". Our proposed method allows the

731 systematic and effective exploration of the conflict between users' knowledge and system

732 requirements/challenges, thus providing important insights for system developers that aim at creating

733 a broadly accessible development platform. Moreover, this method can be applied in other contexts

734 where learnability might be an issue for it allows to identify possible sources of faulty interaction and

735 sub-tasks where the users' background knowledge will be insufficient to guarantee a successful

736 performance of the task at hand.

## KEY POINTS

738 ● An effective Low-Code Development Platform (LCDP) requires an understanding of the

739 distance between the LCDP end-users' conceptualization of programming, and the actions required

740 in the platform.

741 ● We propose and evaluate a Descriptive Cognitive Model (DCM) for the identification of initial

742 use issues in a low-code development platform.

743 ● We propose three mapping rules for the identification of knowledge-system conflicts: over

744 decomposition, under decomposition and no correspondence conflicts.

745 ● Applying the proposed DCM we were able to predict the interaction problems felt by first time

746 users of the LCDP.

## REFERENCES

748 Amir-Heidari, P., Farahani, H., & Ebrahemzadih, M. (2015). Risk assessment of oil and gas well

749 drilling activities in Iran – A case study: Human factors. *International Journal of Occupational*

750 *Safety and Ergonomics*, *21*(3), 276–283.

751 Barricelli, B. R., Cassano, F., Fogli, D., & Piccinno, A. (2019). End-user development, end-user

752 programming and end-user software engineering: A systematic mapping study. *Journal of Systems*

753 *and Software*, *149*, 101–137.

754 Blackwell, A. F. (2017). End-user developers–What Are They Like?. *New Perspectives in End-*

755 *User Development* (pp. 121-135). Cham, Springer International Publishing.

756 Blackwell, A. F., Petre, M. & Church, L. (2019). Fifty years of the psychology of programming.

757 *International Journal of Human-Computer Studies, 131*, 52-63.

758 Brooks, F. P. (1987). No silver bullet: Essence and accidents of software engineering, *Computer*,

759 20(4), 10-19.

760 Ceruzzi, P. E. (2012). *Computing: a concise history*. Cambridge, Massachusetts: MIT Press.

761 Cicchetti, D. V. (1994). Guidelines, criteria, and rules of thumb for evaluating normed and

762 standardized assessment instruments in psychology. *Psychological assessment*, *6*(4), 284-290.

763 Csikszentmihalyi, M. (1990). *Flow: The Psychology of Optimal Experience*. New York: Harper &

764 Row.

765 Dijkstra, E. W. (1982). How do we tell truths that might hurt? In *Selected Writings on Computing:*

766 *A personal Perspective* (pp. 129–131). New York: Springer New York.

767 Dix, A., Finlay, J., Abowd, G. D., & Beale, R. (2004). *Human-Computer Interaction*. Harlow,

768 England: Pearson Education.

769  Embrey, D. E. (1986). SHERPA: A systematic human error reduction and prediction approach.

770  Paper presented at the International Meeting on Advances in Nuclear Power Systems, Knoxville,

771  TN.

772  Fischer, G., Giaccardi, E., Ye, Y., Sutcliffe, A. G., & Mehandjiev, N. (2004). Meta-design: A

773  manifesto for End-user development. *Communications of the ACM*, *47*(9), 33–37.

774  Fowler, M., & Kobryn, C. (2004). *UML distilled: a brief guide to the standard object modeling*

775  *language*. Addison-Wesley Professional.

776  Gajos, K. Z., Czerwinski, M., Tan, D. S., & Weld, D. S. (2006, May). Exploring the design space

777  for adaptive graphical user interfaces. *Proceedings of the working conference on Advanced visual*

778  *interfaces* (AVI '06) (pp. 201-208). New York: ACM.

779  Gonçalves, R. C., Maciel, R. H., Maia, L. M., Nascimento, A. P. T., & Canuto, K. M. (2014).

780  Electric system control room operators: Cognitive task analysis and Human error. *Proceedings of*

781  *the 5th International Conference on Applied Human Factors and Ergonomics (AHFE 2014)*

782  (pp.621-630)*.

783  Huddlestone, J. A., & Stanton, N. A. (2016). New graphical and text-based notations for

784  representing task decomposition hierarchies: Towards improving the usability of an Ergonomics

785  method. *Theoretical Issues in Ergonomics Science*, *17*(5–6), 588–606.

786  ISO - International Organization for Standardization. (2010). Ergonomics of human-system

787  interaction - Part 210: Human-centred design for interactive systems. ISO 9241-210.

788  Lieberman, H., Paternò, F., Klann, M., & Wulf, V. (2006). *End-user development*. Dordrecht:

789  Springer.

790  Maeda, J. (2006) *The laws of simplicity*. MIT press.

791     MacLean, A., Carter, K., Lövstrand, L., & Moran, T. (1990). User-tailorable systems: Pressing the

792     issues with buttons. *Proceedings of the SIGCHI conference on Human factors in computing systems*

793     (pp. 175-182). ACM.

794     Marshall, A., Stanton, N., Young, M., Salmon, P., Harris, D., Demagalski, J., Dekker, S. (2003).

795     *Development of the Human Error template: A new methodology for assessing design induced errors*

796     *on aircraft flight decks. Final Report of the ERRORPRED Project.* London.

797     Marvin, R. (2018). The best low-code development platforms for 2019. Retrieved April 30, 2019,

798     from https://www.pcmag.com/roundup/353252/the-best-low-code-development-platforms

799     Mayer, R. E. (1981). The psychology of how novices learn computer programming. *ACM*

800     *Computing Surveys (CSUR)*, *13*(1), 121-141.

801     Moran, T. P. (1983). Getting into a system: External-Internal task mapping analysis. *Proceedings of*

802     *the SIGCHI Conference on Human Factors in Computing Systems*, (December), 45–49.

803     Myers, B. A., Pane, J. F., & Ko, A. (2004). Natural programming languages and environments.

804     *Communications of the ACM*, *47*(9), 47.

805     Nielsen, J. (1994). *Usability engineering*. Elsevier.

806     Oppermann, R. (1994). *Adaptive user support: Ergonomic design of manually and automatically*

807     *adaptable software* (pp. 1-13). CRC Press.

808     Paternò, F., & Wulf, V. (2017). *New Perspectives in End-User Development*. Springer.

809     Powers, D. (2011). Evaluation: from Precision, Recall and F-measure to ROC, Informedness,

810     Markedness and Correlation. *Journal of Machine Learning Technologies*, 2(1), 37-63.

811     Repenning, A., & Ioannidou, A. (2006). What makes end-user development tick? 13 Design

812     guidelines. *End user development* (pp. 51–85). Dordrecht: Springer.

813 Sajaniemi, J. (2008). Guest Editor's Introduction: Psychology of Programming: Looking into

814 programmers heads. *Human Technology: An Interdisciplinary Journal on Humans in ICT*

815 *Environments 4*(1), 4–8.

816 Silva, C. C. (2013). Audiovisual perception in a virtual world: An application of human-computer

817 interaction evaluation to the development of immersive environments. *Proceedings of the 5th ACM*

818 *SIGCHI symposium on Engineering Interactive Computing Systems (EICS 2013)* (pp. 175-178).

819 Stanton, N. A., Hedge, A., Brookhuis, K., Salas, E. & Hendrick, H.W. (2004). *Handbook of human*

820 *factors and ergonomics methods*. Florida: CRC Press.

821 Tharwat, A. (2018). Classification assessment methods. *Applied Computing and Informatics*.

822 Young, R. M. (1981). The machine inside the machine: Users' models of pocket calculators.

823 *International Journal of Man-Machine Studies*, *15*(1), 51–85.

824

825 **Carlos César Loureiro Silva** is the Research and Development Coordinator of the Perception,

826 Interaction and Usability group at CCG - *Centro de Computação Gráfica*. He holds an MSc in

827 Experimental Psychology from the University of Minho (Portugal, 2011) and a PhD degree in

828 Informatics from the University of Minho (Portugal, 2019).

829 **Joana Catarina Fernandes Vieira** is a Usability Analyst and Researcher at CCG - *Centro de*

830 *Computação Gráfica*. She holds an MSc in Experimental Psychology from the University of Minho

831 (Portugal, 2008) and is concluding a PhD in Ergonomics from the University of Lisbon.

832 **José Francisco Creissac Freitas de Campos** is an Auxiliary Professor at the Department of

833 Informatics of the University of Minho and a Senior researcher at HASLab/INESC TEC. He holds a

834 PhD. degree in Computer Science from the University of York (UK, 2001).

835 **Rui Miguel Silva Couto** is a Senior Researcher at the HASLab/INESC TEC & University of

836 Minho. He holds a PhD. degree in Informatics from the University of Minho (Portugal, 2017).

837 **António Manuel Nestor Ribeiro** is an Auxiliary Professor at the Department of Informatics of the

838 University of Minho and a Senior researcher at HASLab/INESC TEC. He holds a PhD. degree in

839 Informatics from the University of Minho (Portugal, 2008).