



Exploring Gender Bias In Remote Pair Programming Among Software Engineering Students: The `twincode` Original Study And First External Replication

Amador Durán Toro^{1,2} · Pablo Fernández^{1,2} · Beatriz Bernárdez^{1,2} · Nathaniel Weinman³ · Aslihan Akalin³ · Armando Fox³

Accepted: 18 October 2023 / Published online: 1 February 2024
© The Author(s) 2024

Abstract

Context Women have historically been underrepresented in Software Engineering, due in part to the stereotyped assumption that women are less technically competent than men. Pair programming is both widely used in industry and has been shown to increase student interest in Software Engineering, particularly among women; but if those same gender biases are also present in pair programming, its potential for attracting women to the field could be thwarted.

Objective We aim to explore the effects of gender bias in pair programming. Specifically, in a remote setting in which students cannot directly observe the gender of their peers, we study whether the perception of the partner, the behavior during programming, or the style of communication of Software Engineering students differ depending on the perceived gender of their remote partner. To our knowledge, this is the first study specifically focusing on the impact of gender stereotypes and bias *within* pairs in pair programming.

Method We have developed an online pair-programming platform (`twincode`) that provides a collaborative editing window and a chat pane, both of which are heavily instrumented. Students in the control group had no information about their partner's gender, whereas students in the treatment group could see a gendered avatar representing the other participant as a man or as a woman. The gender of the avatar was swapped between programming tasks to analyze 45 variables related to the collaborative coding behavior, chat utterances, and questionnaire responses of 46 pairs in the original study at the University of Seville, and 23 pairs in the external replication at the University of California, Berkeley.

Results We did not observe any statistically significant effect of the gender bias treatment, nor any interaction between the perceived partner's gender and subject's gender, in any of the 45 response variables measured in the original study. In the external replication, we observed

Communicated by: Maria Teresa Baldassarre, Jeffrey Carver, Neil Ernst

This article belongs to the Topical Collection: *Special issue on Registered Reports.*

✉ Amador Durán Toro
amador@us.es

Extended author information available on the last page of the article

statistically significant effects with moderate to large sizes in four dependent variables within the experimental group, comparing how subjects acted when their partners were represented as a man or a woman.

Conclusions The results in the original study do not show any clear effect of the treatment in remote pair programming among current Software Engineering students. In the external replication, it seems that students delete more source code characters when they have a woman partner, and communicate using more informal utterances, reflections and yes/no questions when they have a man partner, although these results must be considered inconclusive because of the small number of subjects in the replication, and because when multiple test corrections are applied, only the result about informal utterances remains significant. In any case, more mixed methods replications are needed in order to confirm or refute the results in the same and other Software Engineering students populations.

Keywords Gender bias · Pair programming · Remote pair programming · Distributed pair programming · Software Engineering education · Experiment replication

1 Introduction

Besides being widely used in industry, pair programming is becoming increasingly common in Software Engineering education because of its demonstrated positive influence on grades, class performance, confidence, productivity, and motivation to stay in Software Engineering and Computer Science academic majors (da Silva Estácio and Prikładnicki 2015), especially for women, as reported by Werner et al. (2004).

In pair programming, two partners work closely together to solve a programming task, in which their ability to engage collaboratively with each other is essential. However, these collaborative interactions can be influenced by implicit gender bias (Hofer 2015), which is a widely observed phenomenon even in highly-structured and professional settings, such as those reported by Jarratt et al. (2019) and da Silva Estácio and Prikładnicki (2015), and which is based on the stereotyped assumption that women are less technically competent than men (Martell et al. 1996; Fisher and Cox 2006; Medel and Pournaghshband 2017; Terrell et al. 2017; Allaire-Duquette et al. 2022).

Our study is based on the hypothesis that gender bias will lead to observable differences based on subjects' perceptions of the gender of their pair programming partners, i.e. they will score men and women differently on similar tasks, and they will also behave and communicate differently depending on whether they perceive their partner as a man or as a woman, even though their partner remains the same on all tasks. Specifically, in a non-located, i.e. remote, pair programming setting in which peer gender cannot be directly observed, our goal is to identify the potential effects of gender bias by observing student pairs when the perceived gender of one of the peers changes.

To study our hypothesis, we have applied methodological triangulation (Denzin 2006), using several methods to collect data and approaching a complex phenomenon like human behavior from more than one standpoint (Cohen et al. 2018). In our case, three different data sources have been used: (1) questionnaires to measure changes in subjects' perceptions, (2) data collected automatically during the pair programming tasks to measure behavioral changes, and (3) data produced by several experimenters analyzing the message interchange during the pair programming tasks to measure changes in communication.

Assuming a remote pair programming setting, which has been proved to have similar results than co-located pair programming as reported by Stotts et al. (2003) and Al-Jarrah and Pontelli (2016), our research questions with respect to subjects' perceptions are the following:

RQ₁ Does gender bias affect perceived productivity compared to solo programming? That is, do perceived differences between in-pair and solo productivity depend on the perceived partner's gender?

RQ₂ Does gender bias affect the partner's perceived technical competency compared to one's own technical competency? That is, do perceived differences between one's own and partners' technical competency depend on the perceived partner's gender?

RQ₃ Does gender bias affect the partner's perceived positive and negative aspects? That is, do perceived positive and negative aspects of their partners depend on the perceived partner's gender?¹.

RQ₄ Does gender bias affect how partners' skills are compared? That is, do perceived partners' skills depend on the perceived partner's gender when they are compared?

With respect to the subjects' behavior during remote pair programming, and considering that women are sometimes perceived as less competent in coding because they often adopt less risky approaches (Fisher and Cox 2006; Terrell et al. 2017), we assume that gender bias could cause a subject to be more or less proactive on the programming task, i.e., taking more or less risks, depending on their perception of their self-efficacy and their perception of the competency of their partner (Allaire-Duquette et al. 2022). Thus, our related research question—based on what we can automatically measure—is the following:

RQ₅ Does gender bias affect the frequencies or relative frequencies with which each partner produces source code additions, source code deletions, successful validations, failed validations, and chat utterances? That is, do these frequencies depend on the perceived partner's gender?

Regarding subjects' communication during remote pair programming, we are interested in knowing whether gender bias affects how subjects communicate with their partners, i.e., whether they use a more formal or informal style, and whether they use some types of chat utterances more than others. This interest is motivated by previous research where it is reported that (i) women and men communicate online differently (Hartsell 2005); (ii) the combination of women's lowered perception of self, with the lowered expectations from others can cause them to lower their rates of participation (Medel and Pournaghshband 2017); and (iii) as reported by Oda et al. (2022), the perceived gender of the partner can exert stimulus control over their communication behavior. Thus, our related research questions are the following:

RQ₆ Does gender bias affect the relative frequency of formal and informal chat utterances? That is, does the formality of the messages depend on the perceived partner's gender?

¹ This research question, and its associated variables, were added after the presentation of the related registered report at ESEM'2021 (Durán et al. 2021). We thought that including an open question could improve the data collection process.

RQ₇ Does gender bias affect the frequency or relative frequency of the different types of chat utterances? That is, do the frequencies of the different types of messages depend on the perceived partner's gender?

1.1 The `twincode` Platform

To support our study, we have developed the `twincode` remote pair programming platform (El-Refai et al. 2023), which manages (i) the registration of students collecting demographic data; (ii) the random allocation into experimental and control groups balancing gender proportions, i.e. trying to have the same number of persons of the same gender in both groups; (iii) the random allocation into experimental-control pairs; (iv) the random assignment of programming exercises to individual subjects and pairs; (v) the swapping of gendered avatars between pair programming exercises for those subjects in the experimental group; and (vi) the automatic collection of interaction metrics and chat utterances.

As shown in Fig. 1, `twincode` offers a source code editor where the students concurrently develop the solution to a proposed programming exercise in Javascript and can validate it against several test cases.

Note that, to foster communication, only one partner can validate the source code at the same time and see validation results, which should be communicated to the other partner using the chat window, where they are instructed to collaborate to solve the proposed exercises.

Note also that a gendered avatar is displayed only for the student in the experimental group (see Fig. 1a) but not for the one in the control group (see Fig. 1b).

Experimenters can use `twincode` to create new experimental sessions where they can configure, among other aspects, the type, number, and duration of the programming exercises, and the instructional messages shown to the students. If needed, they can also develop new programming exercises and their corresponding test cases.

The `twincode` platform is in permanent evolution, and several improvements were incorporated for satisfying some emerging requirements during our study, such as allowing the use of Python as an alternative programming language to Javascript for the programming exercises, changing the images used as gendered avatars (see Fig. 9), and improving the user interface with instructions and a gendered message in the chat window (see Fig. 16a and 16b in Appendix B).

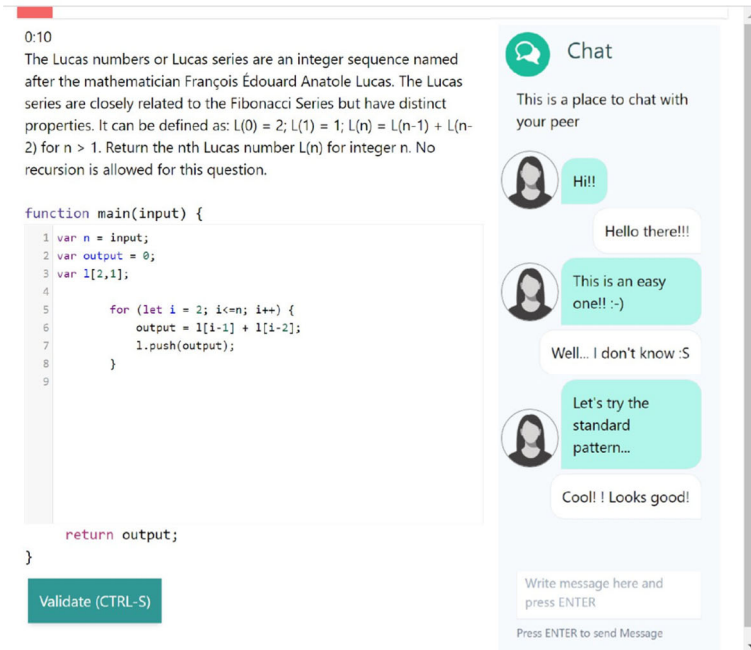
As a companion tool to `twincode`, we have also developed `tag-a-chat`, a tool that help experimenters code chat utterances using different sets of tags, as shown in Fig. 17 in Appendix B.

To assist experimenter s during the training stage of the coding, `tag-a-chat` automatically computes metrics such as Cohen's *kappa* (for two coders) and Fleiss's *kappa* (for three or more coders) in those dialogs that are being coded by several experimenters to achieve inter-coder reliability assessment (O'Connor and Joffe 2020; Syed and Nelson 2015)².

1.2 Pilot Studies

After presenting a very initial approach to our study (Akalin et al. 2021), and to get early feedback on (i) the comprehensibility and internal consistency of the scales used in the

² Although commercial qualitative analysis tools such as MAXQDA (<https://www.maxqda.com>) or Atlas.ti (<https://atlasti.com>) are available, we decided to develop `tag-a-chat` because they are not specifically designed for coding chat utterances, the support for inter-coder reliability metrics is limited, and we prefer to be able to expand its functionality to our future needs and let other researchers use it free of charge.



(a) Experimental group — gendered avatar



(b) Control group — no avatar

Fig. 1 twincode user interface for subjects in the experimental and control groups (original study version)

Table 1 Chat utterance tags by Rodríguez et al. (2017) augmented with orthogonal informal/formal tags

Tag	Description	Examples
I	Informal	<i>LOL! Hahaha!</i>
F	Formal	All messages except informal
S	Statement of information or explanation	<i>We need to create a program for kids to learn math</i>
U	Opinion or indication of uncertainty	<i>Unsure how to add strings together</i>
D	Explicit instruction	<i>Wait put the if back</i>
SU	Polite or indirect instruction	<i>Maybe we can do if user choice = +</i>
ACK	Acknowledgement	<i>Oh ok gotcha</i>
M	Meta-comment or reflection	<i>Hmmm</i>
QYN	Yes/no question	<i>Can the answer be negative?</i>
QWH	Wh- question (who, what, where, when, why, and how)	<i>How do I take in their input?</i>
AYN	Answer to yes/no question	<i>Yea</i>
AWH	Answer to wh-question	<i>The program should be able to generate erroneous questions</i>
FP	Positive task feedback	<i>Oh nice</i>
FNON	Non-positive task feedback	<i>Thats weird</i>
O	Off-task	<i>Wow its sweet in this room</i>

questionnaires; (ii) the usability and performance of the `twincode` platform; and (iii), the applicability of the chat utterance coding based on the one proposed by Rodríguez et al. (2017) and shown in Table 1, two pilot studies with a limited number of students were carried out at the University of Seville and University of California, Berkeley (UC Berkeley) during the 2020–21 academic year.

As a result, the questionnaires were reorganized into three scales that were assessed for internal consistency (see Appendix A), the initial set of chat utterance codes was augmented with formality codes, and the performance and reliability of the `twincode` platform was improved.

1.3 Other Gender Identities

While we recognize that many Software Engineering students may not identify as either men or women, our initial exploration focuses primarily on interactions between students who identify as one of these. The potential biases in interactions involving gender-fluid, gender-nonconforming, and nonbinary students is a complex topic deserving its own subsequent study.

1.4 Structure of the Paper

The rest of the paper is organized as follows. Section 2 reviews related work, although to our knowledge, this is the first study specifically focusing on the impact of gender bias *within* pairs in pair programming. Sections 3 and 4 describe the original study carried out at the

University of Seville (December 2021) and its first external replication performed at UC Berkeley (May 2022) respectively. Section 5 discusses the two studies and the threats to their experimental validity. Finally, Section 6 draws conclusions and proposes future work.

2 Related Work

Several systematic literature reviews (SLR's), which are summarized in Table 2, have compiled the empirical research on pair programming in higher education, including (da Silva Estácio and Prikladnicki 2015), which is focused on distributed pair programming from a teaching perspective.

The SLR by Salleh et al. (2010) reveals that the most important factor under study is solo versus pair programming in terms of effectiveness, quality of code, and satisfaction while students are programming, concluding that pair programming is more effective and satisfactory than solo programming. However, with respect to quality, findings are inconclusive.

Other SLR's, such as the ones by Hanks et al. (2011), Kaur Chahal et al. (2021), and Hawlitschek et al. (2022), show that the focus of the studies is broadened, including factors such as personality, motivation, problem solving, troubleshooting, efficiency, confidence, self-esteem, skill level, gender, or enjoyment but not gender bias. In general, students rate pair programming positively compared to solo programming. Nevertheless, pair programming is effective but not always efficient, as it may take longer.

By means of controlled experiments, remote and co-located pair programming are compared by Stotts et al. (2003) and Al-Jarrah and Pontelli (2016), showing similar results. In most cases, the analyzed variables are related to performance in terms of time, quality, or code tests passed. Students perceptions have also been analyzed in terms of confidence, satisfaction, motivation, or personality by Salleh et al. (2014).

Regarding primary studies, Table 3 summarizes the empirical studies on the influence of gender in pair programming, including findings such as (i) same-gender pairs are more "democratic"; (ii) women working in pairs were more confident than those working solo; and (iii) in mixed-gender pairings, women are less confident compared to same-gender pairings, and report no increase in enjoyment for pair programming compared to solo programming, an effect that is significantly observed in men (Kaur Chahal et al. 2021). Although such studies reveal that gender seems to be a key factor, none of them study gender bias in pair programming.

Many factors other than gender may affect the outcomes of remote programming sessions (Chaparro et al. 2005; Thomas et al. 2003). Previous research on productive pairing looked at factors such as skill levels, autonomy in choosing one's partner (Xinogalos et al. 2017), and different personalities (Hannay et al. 2010). Nevertheless, the work on gender composition of pairs found conflicting results about whether same-gender or mixed-gender pairings are more effective (Choi 2015, 2013; Hofer 2015; Kaur Kuttal et al. 2019). One possible explanation is that gender correlates with other dimensions that may affect the pairs' collaboration, but these correlations may vary between different environments. For example, women in a class may, on average, have higher skill level than men because they had to face more societal barriers to enter the class. On the other hand, they may, on average, have lower skill level if women with no background are more actively recruited.

Table 2 Summary of secondary studies (SMS or SLR) in pair programming in chronological order

Reference	Selected Papers	Main Factors Analyzed	Conclusions
Salleh et al. (2011)	74 papers selected from 1999 - 2000 period; Controlled, longitudinal, observational studies either investigating factors impacting effectiveness of PP (23%) or those measuring effectiveness (90%) via various metrics such as quality (44%)	14 compatibility factors (personality type, actual and perceived skill level, communication skills, self-esteem, gender, ethnicity, learning style, work ethic, time management ability, feel-good factor, confidence level, type of role and type of tasks) and 4 main measures of effectiveness in PP: technical productivity (time spent, knowledge/skill transfer, task performance, code accuracy, number and types of problem, number of solutions to pass test cases), program design and quality (expert opinion, std quality model, code coverage, number of tests passed/failed, LOC, design scores/quality, number of code defects), academic performance (assignment, final, midterm quiz, project and test scores, course grade, course completion rate, retention rate), satisfaction (pair formation, increased knowledge and confidence, positive attitude about collaboration, enjoyment and social interaction)	Paired students report higher satisfaction and achieve productivity similar or better than solo students. Implementing PP in the classroom or lab does not lead to any detrimental effect on students' academic performance. While PP had no significant advantage in improving students' performance in final exams over solo programming (effect size = 0.16) it was effective in helping students get better scores in their assignments (effect size = 0.67). Personality type, actual and perceived skill levels are investigated the most in PP studies but effects of personality were inconclusive. Pair works well when both students have similar abilities and motivation. Students prefer to pair with someone of similar skills, and students' skill level is the most important factor influencing effectiveness of PP. Most popular metric to measure productivity is time spent on completing tasks. Code quality is another common metric to measure productivity that can be measured as internal, external or general categories. When quality was measured according to academic performance and expert opinion (external), students who pair-programmed produced a better quality program compared to students who programmed alone. However, when the quality of the work produced by pair and solo students was measured using metrics at the internal code level, results were contradictory.

Table 2 continued

Reference	Selected Papers	Main Factors Analyzed	Conclusions
Saini et al. (2021)	68 papers selected via thematic analysis	Comparing PP vs. SP settings, student performance, student attitude and enjoyment	PP has mixed effects on students' performance, but almost universally positive effects on students' attitudes (i.e. enjoyment) toward programming. Analysis point out the problems such as small sample size in majority of the studies preventing generalization of the results; the lack of context in reported PP experiments obscures validity of results, and that longitudinal analysis of PP experiments with learning tasks of increasing size and complexity is important to build real-world evidence of the practice.
Korber and Motschnig (2021)	41 papers published after 2000;	Comparing solo vs pair programming, personality, motivation, problem solving, troubleshooting, effectiveness, efficiency, confidence, self esteem, skill level, gender, enjoyment	PP positively impacts motivation, self-esteem and confidence of learners. Students report more fun solving the assignments and think PP helps them solve problems faster. Pair programming students reported more effective programming in both visual and text-based languages but earned more achievement points in the text-based language (python) compared to the visual. PP is effective, but not always efficient. Students who find introductory programming or learning a text-based language especially benefit from the positive effects of pair programming: an appreciative and clear communication regarding possible mistakes and misunderstandings is one of the key factors. Social factors such as gender, personal relationships, effects of successes and failures (attitude), distribution of workload and the influence of partner changes must be considered.

Table 3 Summary of primary studies on gender and pair programming in chronological order

Reference	Object of study	Metrics	Findings
Katira et al. (2005)	Compatibility of student pair programmers	Web-based peer evaluation survey that required the students to evaluate the contributions of their partner and the perceived pair compatibility	Students are compatible with partners whom they perceive of similar skill. Mixed-gender pairs are less likely to report compatibility.
Stetsos et al. (2009)	Effect of personality heterogeneity on PP effectiveness	The Keirsey Temperament Sorter personality test; PP effectiveness is measured by output/performance, communication, velocity, design correctness, passed acceptance tests; pair collaboration-viability is measured by satisfaction, knowledge acquisition and participation.	Heterogeneous personality pairs shows better communication, pair performance and pair collaboration-viability than homogeneous pairs. For heterogeneous pairs, design and code correctness is positively correlated with communication transactions (more communication leads to higher correctness), and satisfaction regarding collaboration, knowledge acquisition and participation was significantly higher.
Salleh et al. (2014)	Personality traits on PP effectiveness	Five Factor Model (FFM2); Conscientiousness, Neuroticism, and Openness to experience	Only openness has a significant role in differentiating paired students' academic performance.
Choi (2015)	PP gender combinations	Productivity, quality of source code, compatibility and communication between pairs	Pair compatibility and communication levels significantly vary between the same gender pair type, woman-woman and man-man.
Gómez et al. (2017)	PP gender combinations	Productivity	Similar productivity rates for the three gender pair combinations. Greater variability of productivity rates with mixed gender pairs (man-woman) was observed.
Jarratt et al. (2019)	PP gender combinations	Weekly attendance, work accomplished during lab and perceived productivity	Students who were randomly assigned a woman partner (rather than a man) attended classes more often, were more confident that the solution was correct, and more confident in the finished product that they developed. However, being assigned a woman partner was also associated with completing a smaller percentage of the assignment.

Table 3 continued

Reference	Object of study	Metrics	Findings
Ying et al. (2021b)	Effect of structured roles in PP; motivation and stress for men and women	Lexical features (number of messages, message length, sentiment), Intrictric Motivation Inventory score (IMI) measuring Interest/Enjoyment, Perceived Competence, Effort/Importance, Pressure/Tension, Perceived Choice, Value/Usefulness, and Relatedness; Self-reported stress, perceived competence, perceived choice, learning gain	No significant differences found between structured vs unstructured PP roles. Women reported significantly higher levels of stress, lower levels of perceived competence in their computing abilities and less perceived choice compared to men during a remote PP activity. Dialogue features significantly correlated with women's reports of stress, perceived competence, or perceived choice. Women tended to feel more relaxed if their partner sent longer messages on average or used more positive language.
Ying et al. (2021a)	Analyzing the differences between women and men's awareness of CS gender gap	Survey (which included six questions related to the gender gap in CS), and some follow-up interviews discussing the experiences and perceptions of CS gender gap.	Men were less aware, had milder beliefs and shallow understanding of the gender disparities in computer science. Women were significantly more aware of the gender gap and felt significantly stronger that efforts should be made to reduce the gender gap. Some participants also expressed discomfort at the idea of opportunities for women within CS because they did not think that those were fair; these students would benefit from understanding the idea of equity over equality.
Galdo et al. (2022)	Young learners in remote compared to co-local PP	Perception and experiences by means of remote collaboration logs, interviews and self evaluation survey	Students felt successful in remote approach, had positive experiences with collaboration, reported remote PP made them have more autonomic and efficient in navigation compared to co-located PP. Furthermore, students recognized who partner with friends become more confident throughout the learning process (vs non friend partners).

3 Original Study (Seville Dec, 2021)

In this section, the original study carried out at the University of Seville in December 2021 is reported, including most of the experimental settings which are in common with the external replication performed at the UC Berkeley in May 2022, reported in Section 4.

3.1 Participants

In the original study carried out at the University of Seville in December 2021, the participants were third-year students of the Degree in Software Engineering enrolled in any of the three groups of the Requirements Engineering course taught in Spanish³. The final number of valid⁴ subjects was 92, arranged in 46 pairs. Only 9 students could not finish the study because of technical problems during the tasks. Considering the 92 valid subjects, 15 identified as *woman* (16.30%), 1 as *non-binary* (1.09%), and the rest as *man* (82.61%) during the registration process.

Note that, although the percentage of women is low, it is above the average percentage in the Degree in Software Engineering at the University of Seville, which unfortunately is close to 11% according to the last academic year official statistics (University of Seville 2021). Note also that, due to the 9 students dropped by technical reasons, the percentage of women could not be kept the same in the control (6 women, 14.29%) and experimental (9 women, 19.57%) groups than in the sample (16.30%), which was our initial intention.

3.2 Experiment Execution

Some weeks before experiment execution, in order to recruit participants, the students enrolled in the three groups of the Requirements Engineering course taught in Spanish were motivated to voluntarily participate in the study as an interesting experience in remote pair programming, but without mentioning neither that the main goal was to study the potential effect of gender bias, nor they were going to be paired with the same classmate during all the study. We also remarked that for the purpose of the study, they must remain anonymous to their partners, so they must neither mention nor ask any personal information, thus not discovering that their partners were always the same person. After providing all that information, including that the participation in the study counted for a 5% bonus on their grades to prevent dropout, the interested students registered in the *twincode* platform providing some demographic data and accepting the participation conditions.

The experiment execution, which is graphically represented in Figs. 2 and 3, took place the same day for the three groups of students of the course during their laboratory sessions, as shown in Fig. 4⁵.

All registered students logged into the *twincode* platform, which automatically allocated them into the control and experimental groups balancing the proportion of women in

³ There is a fourth group of the Requirements Engineering course which is taught in English and in which the enrolled students are approximately 50% Spanish and 50% Erasmus students coming from other countries in the European Union (EU) or from non-UE countries like Israel or Georgia. They were not invited to participate in the study because their command of Spanish was not good enough to chat with a randomly assigned classmate, who would have undoubtedly identified them as foreign students.

⁴ The criteria for considering a subject as valid are strongly dependent on properly performing the experimental tasks, which are described in Section 3.2. The criteria themselves are specified in Section 3.6.

⁵ By the time the experiment was carried out, COVID-19 restrictions in force in Andalucía allowed students to be in the same classroom but wearing masks.

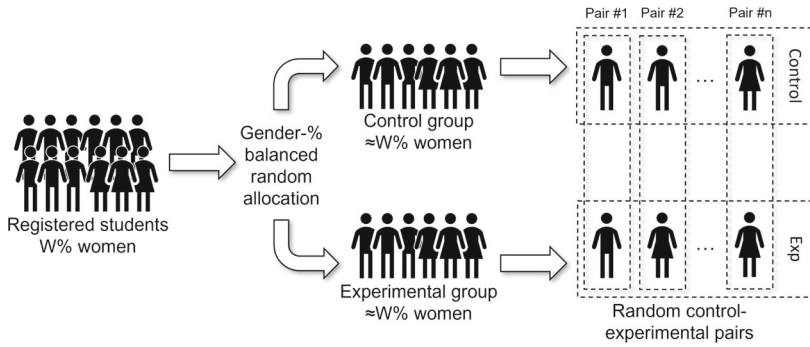


Fig. 2 Experimental process (subject allocation to groups)

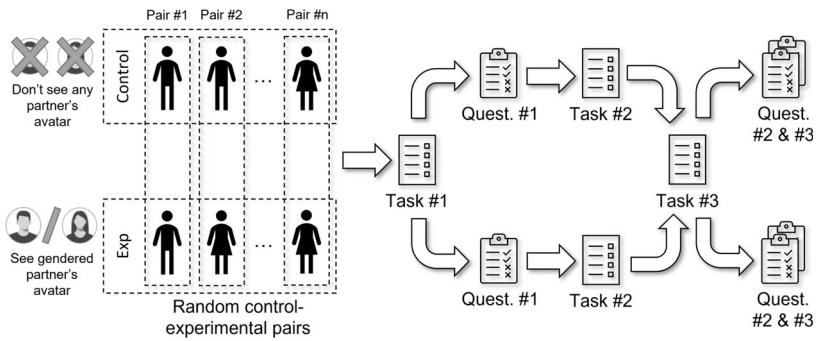


Fig. 3 Experimental process (tasks)



Fig. 4 Experiment execution at University of Seville, Dec 2021

each group as much as possible. Once all the students were allocated to groups, they were randomly allocated into control-experimental pairs by the platform (see Fig. 2).

After subject allocation, the pairs were presented a programming exercise that they had to solve collaboratively using `twincode` (labeled as Task#1 in Fig. 3). They were given 10 minutes to solve a first exercise and another 10 minutes to solve a second exercise, thus a total time of 20 minutes. After the first 10-minute period, the second exercise was presented independently of whether the first one was finished successfully or not. Both exercises were randomly selected from a pool of exercises of similar complexity. During this programming exercise in pairs, subjects in the control group received no information about the gender of their partners, whereas subjects in the experimental group could see their partners as having a clearly gendered avatar randomly selected by the platform (see Fig. 1). At the end of the 20-minute period, they were asked to individually fill in a questionnaire (labeled as Quest.#1 in Fig. 3) about the perceived productivity compared to solo programming, the perceived partner's technical competency compared to their own, and about the partner's positive and negative aspects. They were given 10 minutes to fill in the questionnaire.

After filling the first questionnaire, the students were presented another programming exercise to be solved individually in 10 minutes (labeled as Task#2 in Fig. 3). In the case they finished earlier, another exercise of similar complexity was randomly presented. The main purpose of this individual task was to make students forget about their first partners, i.e. their style of writing chat utterances or source code, so they did not recognize them in the second in-pair task.

After the individual task, pairs were presented again a new collaborative programming exercise that they must solve in similar conditions to the exercise in Task#1. In this second in-pair exercise, the gendered avatar was swapped with respect to the first exercise for the subjects in the experimental group. For those in the control group, they continued to receive no information on their partners' genders. Note that pairs were kept the same in order to reduce the variability due to the subjects themselves, which could possibly have had a confounding effect in case of a new pair allocation for Task#3 (see Section 3.5.1 for details).

Once Task#3 was finished, students were asked to fill a questionnaire (labeled as Quest.#2 in Fig. 3) with the same questions than the one they filled after Task#1 but referred to the second partner, and another questionnaire (labeled as Quest.#3 in Fig. 3) comparing the skills of the first and second partners and whether they remembered the gendered avatars of their partners or not. They were given 15 minutes for responding both questionnaires.

Finally, they were informed about the actual purpose of the study. At that point, they were allowed to withdraw their data if they wished, although none of them opted for doing so.

3.3 Factors (Independent Variables)

The four factors, i.e., independent variables, in both the original experiment and the replication are following.

group nominal factor representing the group (EXPERIMENTAL or CONTROL) subjects were randomly allocated to.

time nominal factor representing the moment (T_1 and T_2) in which the first and second in-pair tasks were performed by the subjects.

ipgender nominal factor representing the induced partner's binary gender (MAN or WOMAN for the experimental group, and NONE for the control group) during the in-pair tasks.

Solo programming or pair programming?

All questions should be answered on a numerical scale from 0 to 10, where the midpoint is 5.

2. Regarding the programming exercises you just did, how do you think you would have been **more productive**, programming solo or programming with the partner assigned to you? *

0 - programming solo	5 - the same in both cases	10 - programming in pairs								
0	1	2	3	4	5	6	7	8	9	10

Fig. 5 First response item for pp variable in questionnaires #1 & #2 as presented to the subjects

gender nominal factor representing subject's gender, which may be MAN, WOMAN, or any other option as freely expressed in the demographic form during registration.

3.4 Response Variables (Dependent Variables)

The response variables, i.e., dependent variables, in both studies are described below, organized according to the corresponding three data sources—questionnaires, twincode platform, and chat utterance coding.

3.4.1 Perceived Variables (Questionnaires)

The response variables measuring subjects' perception are mainly scales composed by four or more 0–10 linear numerical response items and they are computed as the average of their corresponding items. Following the recommendations by Hopper (2014), the 0–10 items are labeled not only in the first and last points, but also in the midpoint (see Fig. 5). They are described below.

pp interval variable composed of four 0–10 numerical response items ($pp_{1...4}$) measuring the subject's own *perceived productivity* during each pair programming task compared to solo programming (see RQ₁). Low values correspond to better solo programming productivity whereas high values correspond to better pair programming productivity (see Fig. 5 for an example of a response item and Section A.1 in the Appendix for all the response items in the scale).

pptc interval variable composed of four 0–10 numerical response items ($pptc_{1...4}$) measuring the subject's *partner's perceived technical competency* compared to their own after each in-pair task (see RQ₂). Low values correspond to higher subject's productivity, whereas higher values correspond to higher partner's productivity (see Section A.2 in the Appendix for all the response items).

ppa ratio variable counting the number of *partner's positive aspects* identified by the subject after each in-pair task (see RQ₃)⁶. This variable is automatically computed from an open question item in which subjects are asked to write the most positive and negative aspects of

⁶ According to the four scales of measurement introduced by Stevens (1946), variables ppa and pna are defined as ratio variables because they are numerical variables in which zero represents a lack of the attribute (see Section 2.2 in (Navarro 2018) for an excellent explanation, or (GraphPad 2023) for a graphical representation). Note that this is not the case for the pp, pptc, and cps interval variables, in which zero usually means "the same in both cases" or "both equally".

their partners in the previously performed pair programming exercise (see Section A.3 in the Appendix). They are instructed to prefix positive aspects with a plus sign (+) and negative ones with a minus sign (-). This variable is the result of automatically counting the number of plus signs in the text of the open question.

pna ratio variable counting the number of *p*artner's *n*egative *a*spects identified by the subject after each in-pair task (see RQ₃). In a similar way to the **ppa** variable, this variable is the result of automatically counting the number of minus signs in the text of the aforementioned open question (see also Section A.3 in the Appendix).

ppgender nominal variable measuring the *p*erceived *p*artner's *g*ender during the in-pair tasks. To measure this variable, subjects are asked in questionnaire #3 whether they remember if their partners showed some avatars in chat windows or not. If the answer is NO or *I don't remember* (IDR), this variable is assigned the NONE or IDR levels at T₁ and T₂. If the answer is YES, then the subjects are asked for the avatars of the first and second partner, having MAN, WOMAN, or IDR as options, as shown in Fig. 6.

cps interval variable composed of five 0–10 numerical response items (cps_{1...5}) measuring whether the subject perceived better skills in their first or second partner in the in-pair tasks, i.e., compared *p*artners' skills (see RQ₄). Low values correspond to the first partner, whereas high values correspond to the second partner (see Section A.4 in the Appendix for all the response items).

In the case of the experimental group only, this variable is transformed after collection in such a way that low values correspond to the partner for whom the induced gender was MAN, and high values to the partner for whom the induced gender was WOMAN, in order to analyze whether there is a gender bias in the scoring.

3.4.2 Behavior-Related Variables (twincode Platform)

The response variables automatically collected by the `twincode` platform and related to the behavior during the in-pair programming exercises (see RQ₅) are listed below. Every variable v represents a frequency, i.e., a count, and its associated relative frequency is computed with respect to the the sum of the frequencies of the two subjects in a pair. For example, let us suppose that subjects i and j are the two members of a pair, and v_i and v_j are the corresponding values of the v variable. In this case, the relative frequencies for each subject would be $\frac{v_i}{v_i+v_j}$ and $\frac{v_j}{v_i+v_j}$, respectively.

sca / sca_rf Ratio scale variables representing the count and relative frequency of characters added by a subject to the source code window during an in-pair task (source code *a*dditions).

scd / scd_rf Ratio scale variables representing the count and relative frequency of characters deleted by a subject from the source code window during an in-pair task. (source code *d*eletions).

okv / okv_rf Ratio scale variables representing the count and relative frequency of successful (*ok*) validations of the source code performed by a subject during an in-pair task.

kov / kov_rf Ratio scale variables representing the count and relative frequency of unsuccessful (*ko*) validations of the source code performed by a subject during an in-pair task.

dm / dm_rf Ratio scale variables representing the count and relative frequency of *d*ialog *m*essages (chat utterances) sent by a subject during an in-pair task.

3.4.3 Communication-Related Variables (Utterance Tagging)

The chat utterances registered in the `twincode` platform during the in-pair tasks were manually tagged according to two orthogonal dimensions. The first dimension uses the 13 tags (from S to O in Table 1) proposed by Rodríguez et al. (2017). The second dimension classifies each message as FORMAL or INFORMAL, considering as formal the usual way in which a university student would communicate textually to a professor and informal otherwise.

For the tagging process, we followed a process inspired by the work of O'Connor and Joffe (2020), in which two researchers each tagged 60% of the data, covering all dialogue messages. The overlapping subset of 20%, which was used for the initial training, established the inter-coder reliability using Cohen's *kappa*, which was $\kappa = 0.796$ for the formal/informal tags, and $\kappa = 0.754$ for Rodríguez et al. tags, both indicating *substantial* agreement and sufficient reliability for further coding according to Syed and Nelson (2015).

One last question

Select the option you consider most appropriate

16


Do you remember if your partners showed an avatar in the chat window during the programming exercises?

Yes, they showed an avatar


No, they didn't show any avatar

I don't remember

17



A



B

Do you remember what your first partner's avatar was? *

Selecciona la respuesta
^

Avatar A

Avatar B

I don't remember

Fig. 6 Section in questionnaire #3 for partner's perceived gender (ppgender) variable

The response variables related to the manual tagging of the chat utterances (see RQ₆ and RQ₇) correspond to the tags in Table 1 and are listed below.

Every variable represents a frequency, i.e., a count, and its associated relative frequency is computed with respect to the number of chat utterances generated by the subject during an in-pair task, which is defined by the *dm* variable specified in previous section.

i / i_rf Ratio scale variables representing the absolute and relative frequency of *informal* messages generated by a subject during an in-pair task.

f / f_rf Ratio scale variables representing the absolute and relative frequency of *formal* messages generated by a subject during an in-pair task.

s / s_rf Ratio scale variables representing the absolute and relative frequency of *statement of information or explanation* messages generated by a subject during an in-pair task.

u / u_rf Ratio scale variables representing the absolute and relative frequency of *opinion or indication of uncertainty* messages generated by a subject during an in-pair task.

d / d_rf Ratio scale variables representing the absolute and relative frequency of *explicit or direct instruction* messages generated by a subject during an in-pair task.

su / su_rf Ratio scale variables representing the absolute and relative frequency of *polite or indirect instruction or suggestion* messages generated by a subject during an in-pair task.

ack / ack_rf Ratio scale variables representing the absolute and relative frequency of *acknowledgment* messages generated by a subject during an in-pair task.

m / m_rf Ratio scale variables representing the absolute and relative frequency of *meta-comment or reflection* messages generated by a subject during an in-pair task.

qyn / qyn_rf Ratio scale variables representing the absolute and relative frequency of *yes/no question* messages generated by a subject during an in-pair task.

qwh / qwh_rf Ratio scale variables representing the absolute and relative frequency of *wh-question* (who, what, where, when, why, and how) messages generated by a subject during an in-pair task.

ayn / ayn_rf Ratio scale variables representing the absolute and relative frequency of *answer to yes/no question* messages generated by a subject during an in-pair task.

awh / awh_rf Ratio scale variables representing the absolute and relative frequency of *answer to wh-question* messages generated by a subject during an in-pair task.

fp / fp_rf Ratio scale variables representing the absolute and relative frequency of *positive task feedback* messages generated by a subject during an in-pair task.

fnon / fnon_rf Ratio scale variables representing the absolute and relative frequency of *non-positive task feedback* messages generated by a subject during an in-pair task.

o / o_rf Ratio scale variables representing the absolute and relative frequency of *off-task* messages generated by a subject during an in-pair task.

3.5 Confounding Variables

The confounding variables that were controlled during both studies are described below.

Table 4 Contingency table for induced partner's gender (ipgender) vs. perceived partner's gender (ppgender)

Induced Gender	Perceived Gender			
	MAN	WOMAN	NONE	IDR
MAN	28 (60.87%)	1 (2.17%)	6 (13.05%)	11 (23.91%)
WOMAN	1 (2.17%)	27 (58.70%)	6 (13.05%)	12 (26.09%)
NONE	0 (0.00%)	0 (0.00%)	52 (56.52%)	40 (43.48%)

3.5.1 Subject's technical skills

To control the variability caused by each subject on their partner, pairs were kept the same during the entire experiment, although the subjects were not informed about this fact. Ideally, this would make the conditions of the two in-pair tasks the same except for the programming exercises (see below) and for the induced gender in the case of the experimental group.

3.5.2 Programming exercises

In order to avoid potential differences among the programming exercises used during in-pair tasks, they were all of similar complexity and were randomly assigned.

3.6 Data Analysis

The data analysis was performed only for those subjects considered as valid according to the following criteria: (i) to have filled in both questionnaires; (ii) to have their metrics correctly collected by the `twincode` platform; (iii) to have been paired with another valid subject; and (iv) not to have disclosed their gender or their partner's during the in-pair exercises; This resulted in 46 pairs, i.e. 92 valid subjects, with only 9 subjects dropped because of technical problems with their connections to the `twincode` platform, as previously mentioned in Section 3.1.

3.6.1 Correlation between Induced and Perceived Gender

Before analyzing between and within-group relationships, the correlation of the induced and perceived gender in both groups was analyzed in order to know whether the treatment had been effectively administered to the subjects⁷.

For that purpose, the results of the contingency table in Table 4 were analyzed observing that the percentage of subjects who were induced to think that their partner was a MAN and that effectively remembered they saw a MAN avatar was close to 61%, whereas in the case of WOMAN avatars the percentage was close to 59%. Although Cramer's V for Table 4 showed a *large* effect (0.709) according to Gravetter and Wallnau (2004), we decided to exclude from the remaining analyses those subjects in the experimental group for whom the induced and perceived gender did not match, because we considered that the treatment had not been

⁷ The analysis of the correlation between induced and perceived gender was not included in the registered report originally submitted to ESEM'2021 (Durán et al. 2021). We included it thanks to the reviewers' comments, whose suggestion has definitely improved our analysis.

sufficiently effective in their cases⁸. On the other hand, we kept those subjects in the control group who did not perceived any gendered avatar or did not remember it, discarding the rest. As a result, we kept all the subjects in the control group (39 men, 6 women, 1 non-binary) but only 27 (21 men, 6 women) in the experimental group.

3.6.2 Between-Groups Analysis

In the analysis between the control and experimental groups, for every response variable v except for cps ⁹, we computed the distance between the two in-pair tasks as the absolute value of the difference, i.e. $|v(t_2) - v(t_1)|$, since the sign of that difference was not relevant in our case. In our research hypothesis, this distance should be smaller for the students in the control group, who received no information about their partners' genders i.e. no treatment, than for those in the experimental group who effectively perceived two different partners' genders at T_1 and T_2 . Therefore, for every response variable except for cps , we performed a one-tailed unpaired mean difference test between groups, applying a t-test or a Mann-Whitney U test (also known as Wilcoxon test), depending on the results of the normality assumption tests.

In the case of the cps variable, for the control group we expected the mean to be closer to the middle point (5) between the first and second partner, as they were unconsciously comparing the skills of the same person. For the experimental group, we expected the mean to be skewed towards 0 (partner perceived as a MAN) or 10 (partner perceived as a WOMAN) due to the effect of the treatment. Therefore, to detect differences between groups for the cps response variable, we performed an unpaired two-tailed t-test because data distribution was not significantly different from normal distribution.

Contrary to our research hypothesis, no significant differences were observed at $\alpha=0.05$ between the control and experimental groups for any of the 45 response variables described in Section 3.4, including cps . The corresponding boxplots are depicted in Fig. 7, where it can be seen that the difference between means—the circles in the boxes—in both groups were very small. Post hoc power analyses using G*Power (Faul et al. 2007) yielded a statistical power ($1 - \beta$) for the applied tests of 0.132 for small effect sizes ($d \leq 0.2$), 0.548 for medium effect sizes ($d \leq 0.5$), and 0.915 for large effect sizes ($d \leq 0.8$).

3.6.3 Within-groups Analysis

Within the experimental group, we wanted to analyze whether there were differences between the response variables when the same subjects perceived their partners as MEN or WOMEN according to our research hypothesis. We also wanted to study the possible interaction between the perceived partner's gender and the subject's gender.

For those purposes, we performed a two-sided paired mean difference test for every response variable except for cps , using the perceived gender ($ppgender$) as a within-subjects variable, and applying a t-test or a Wilcoxon test depending on the results of the normality assumption tests. For studying the interaction, we performed the corresponding mixed-model

⁸ We applied this strict selection of subjects in the experimental group in a manner consistent with the results of the correlation analysis, considerably reducing the number of subjects, especially in the replication reported in Section 4.

⁹ As commented in its description in Section 3.4.1, the cps variable is measured only once at the end of the experimental process, since it compares first and second partners' skills.

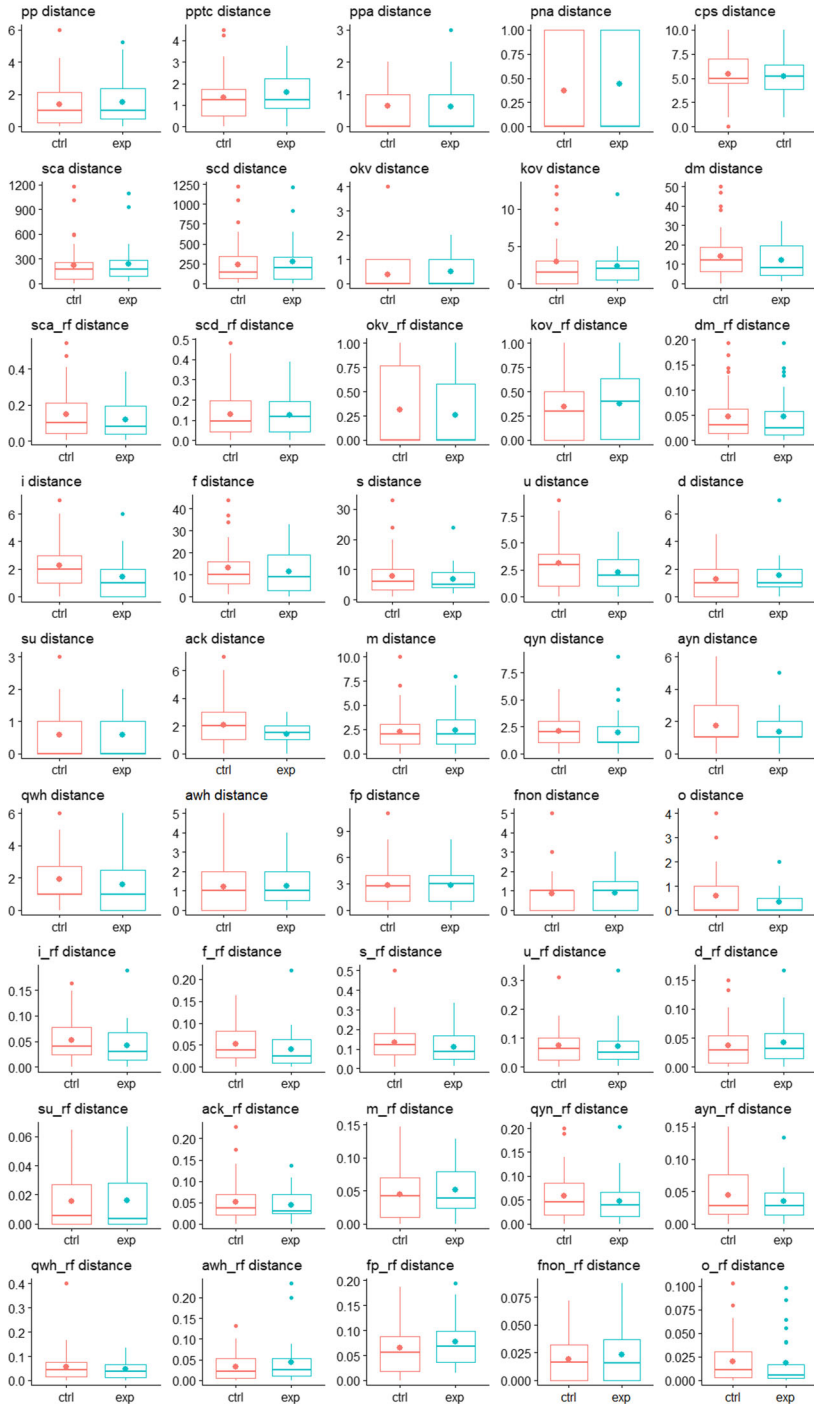


Fig. 7 Boxplots of the 45 response variables for between-groups analysis in the original study

Table 5 Estimated effects on experimental validity of the changes introduced in the replication

Change description	Effect on experimental validity			
	Construct	Internal	External	Conclusion
Third to first year students	–	–	+2	–
Spanish students to U.S. students	–	–	+2	–
Higher percentage of women	–	–	+2	–
Number of subjects reduced	–	–	–	–2
5% grade bonus to \$15 Amazon gift card	–	–	–	–
Remote location of subjects	+2	–1	–	–
Higher number of sessions	–	–1	–	–
Reduced time for tasks	–1	–1	–	–
Different avatars	–1	–	–	–
Blocked exercise assignment	–	+2	–	–
Different programming language	–	–	–	–

Legend: –: it does not affect; –1/ +1: slightly increases/decreases;
 –2/ +2: moderately increases/decreases; –3/ +3: substantially increases/decreases

two-way ANOVA's with the perceived gender (*ppgender*) as a within-subjects variable and the subject's gender (*gender*) as a between-subjects variable.

For the *cps* variable, which passed the Shapiro-Wilk normality tests, we analyzed whether the subject's gender had any effect when comparing partners perceived as MAN or WOMAN by means of a two-tailed unpaired t-test between groups, using *gender* as a between-subjects variable.

Contrary to our research hypothesis, no significant differences were observed at $\alpha=0.05$ between the two levels of the *ppgender* variable for any of the 44 response variables described in Section 3.4. None of the 44 ANOVA tests detected any significant interaction either, and no effect of the subject's gender on the *cps* variable was detected. Post hoc power analyses using G*Power (Faul et al. 2007) yielded a statistical power ($1 - \beta$) for the applied tests of 0.263 for small effect sizes ($d \leq 0.2$), 0.811 for medium effect sizes ($d \leq 0.5$), and 0.991 for large effect sizes ($d \leq 0.8$).

As depicted in Fig. 8, the corresponding boxplots show very small differences between means when partners are perceived as MEN or WOMEN in the experimental group.

4 First Replication (Berkeley May, 2022)

In this section, the first replication carried out at the University of California Berkeley in May 2022 is reported focusing mainly on the changes in the participants and the experiment execution with respect to the original experiment, since the research questions and variables were the same in both studies. For each change, an estimation of their impact on the four types of experimental validity described by Wohlin et al. (2012) is included, following the recommendations by Cruz et al. (2023) about reporting the impact of changes in replications using a 7-point discrete scale from –3 to +3. A summary of the impact of those changes is presented in Table 5, including the labels of the aforementioned scale in its legend.

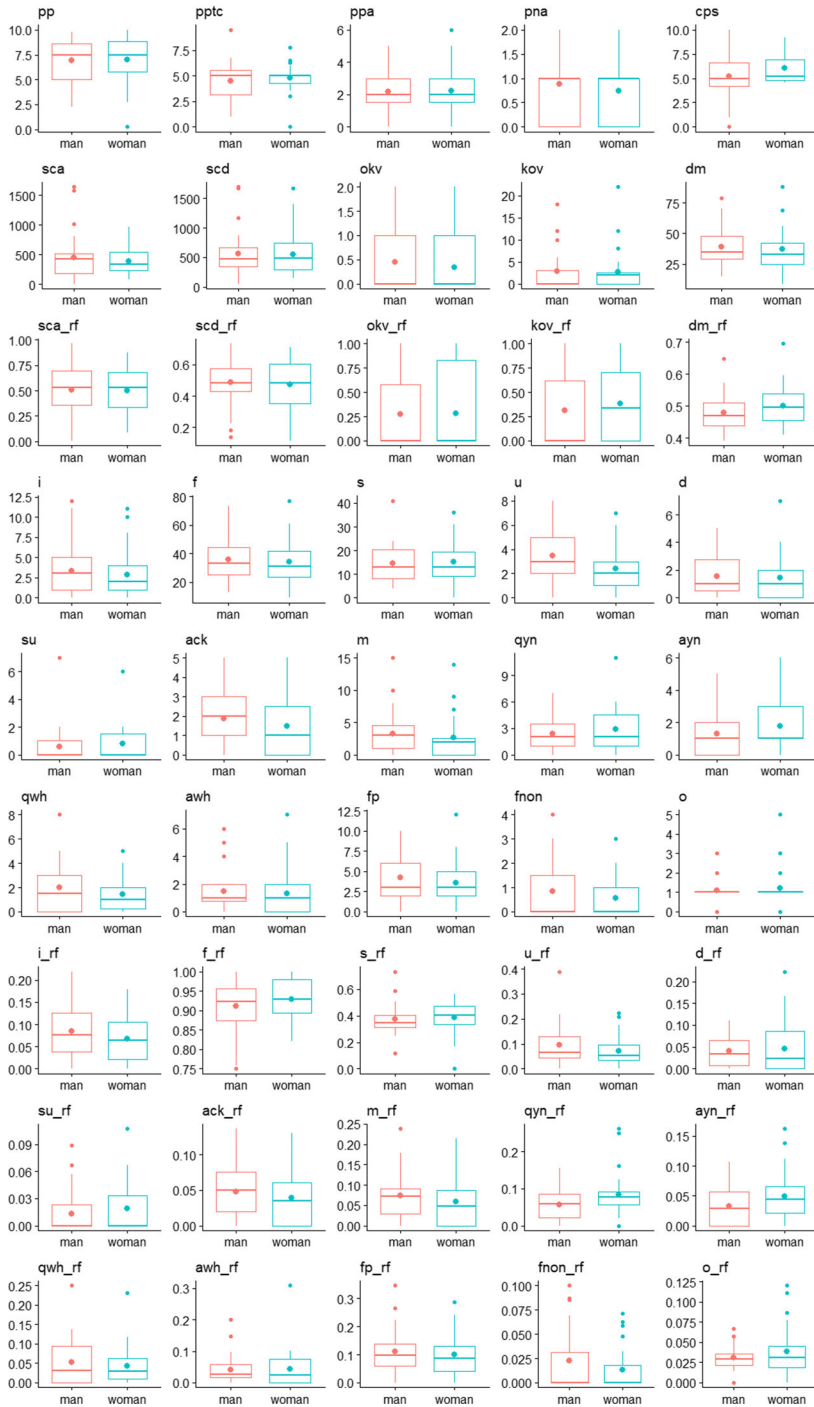


Fig. 8 Boxplots of the 45 response variables for within-groups analysis in the original study

4.1 Participants

In the replication carried out at the University of California, Berkeley, the participants were mainly first year students enrolled in the CS61A (*The Structure and Interpretation of Computer Programs*) and CS88 (*Computational Structures in Data Science*) courses. Applying the same criteria than for the original experiment, the final number of valid subjects was 46, arranged in 23 pairs. Only 6 students, i.e. 3 pairs, were excluded from the initial 52 participants. One pair was dropped due to the disclosure of their identities during the pair programming tasks; another pair was dropped because one of its partners did not actively participate in the experimental tasks; and the third pair was excluded because they lost their connection to the `twincod` platform repeatedly and their metrics could not be properly collected. Among the remaining 46 valid subjects, 26 identified as *woman* (56.52%) and the rest as *man* (43.48%) during the registration process¹⁰.

Note that, contrary to the original experiment, the percentage of women is above that of men because the CS61A and CS88 introductory courses are taken also by students from other majors, usually with a higher presence of women than in Computer Science majors, where is around 25% (University of California, Berkeley 2021). Note also that despite the 6 dropped subjects, the percentage of women in the control (12 women, 52.17%) and experimental (14 women, 60.87%) groups were close to each other.

From our point of view, this change in the sampled population from third-year Spanish students to first-year U.S. students, and the higher percentage of women, increased external validity, but the reduction in 50% of the number of subjects (46 pairs to 23 pairs) reduced conclusion validity.

4.2 Experiment Execution

The experiment execution at the University of California, Berkeley followed the same process than that performed at the Universidad de Sevilla with some changes, which are described in the following sections.

4.2.1 Bonus for Participating in the Study

As commented in Section 3.2, in the original experiment the participation in the study counted for a 5% bonus on students' grades in the Requirements Engineering course they were enrolled in to prevent dropout. In the replication, considering that the students were enrolled in two different courses with different professors, they were offered a \$15 Amazon gift card for participating actively in the study instead of a grade bonus which would have been difficult to manage. In our opinion, this change did not affect any type of experimental validity.

4.2.2 Location of Students and Number of Sessions

In the original experiment, the experimental execution took place during one of the laboratory sessions of the Requirements Engineering course, as shown in Fig. 4. The three groups of the course had the laboratory sessions the same day at different hours, with 30 students per session on average. In the replication, the students performed the experimental tasks

¹⁰ The only student who reported a non-binary gender described as "who cares" was one of the 6 excluded students.

remotely, coordinated by one of the experimenters using Zoom. There were four sessions that took place during a week with 10 students per session on average.

We think that this change increased construct validity with respect to the original study, since the setting was strictly remote rather than being co-located in a laboratory room, but it also decreased internal validity because of the lack of control of the subject's environment, in which interactions with a third person, interruptions, or distraction could occur. On the other hand, having multiple sessions over a week rather than having three consecutive sessions on the same day also decreased internal validity due to the possibility of some students disclosing the purpose of the study to their peers despite being instructed not to do so.

4.2.3 Timing of the Tasks

In the original experiment, the students were given 20 minutes for the pair programming tasks, 10 minutes for the solo task, 10 minutes for the first questionnaire, and 15 minutes for the second and third questionnaires. In the replication, the students were given 15 minutes for the in-pair tasks, 10 minutes for the solo task, 10 minutes for the first questionnaire, and 10 minutes for the second and third questionnaires, due to the constraints imposed by their busy schedule.

We think that the shortened duration of the in-pair tasks and the second and third questionnaires may have compromised construct validity by reducing the time span for measuring the response variables, the interaction time for assessing the partners' skills, and the reflection time before answering each response item. Moreover, it may have weakened the effect of the treatment over confounding variables, thus decreasing also internal validity.

4.2.4 Gendered Avatars

In the original experiment, the gendered avatars used in the chat windows of the subjects in the experimental group were the silhouettes shown in Fig. 9a, whereas in the replication the avatars were those shown in Fig. 9b, which were generated at <https://getavataaars.com/>. The subjects in the replication were also shown a gendered message at the top of the chat window indicating that their partner was connected, e.g. "Your partner (she/her) is connected" (see Figs. 16a and 16b in Appendix B).

In principle, changing the gendered silhouette avatars by more explicit ones and adding a gendered message in the chat window would have increased construct validity, but the correlation between induced gender and perceived gender in the replication worsened with

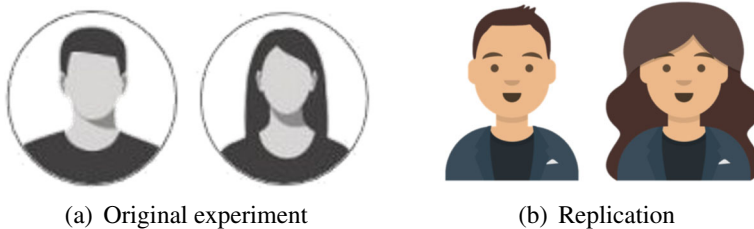


Fig. 9 Gendered avatars used in the original experiment and the replication

Table 6 Contingency table for induced partner's gender (ipgender) vs. perceived partner's gender (ppgender) in the replication

Induced Gender	Perceived Gender			
	MAN	WOMAN	NONE	IDR
MAN	10 (43.48%)	5 (21.74%)	2 (8.70%)	6 (26.09%)
WOMAN	5 (21.74%)	9 (39.13%)	2 (8.70%)	7 (30.43%)
NONE	0 (0.00%)	1 (2.17%)	30 (65.22%)	15 (32.61%)

respect to the original experiment (see Section 4.3.1). As a result, we consider that this change decreased construct validity.

4.2.5 Exercise Assignment

In the original experiment, the programming exercises, which had to be solved using Javascript as the programming language, were randomly assigned to the subjects from a pool of exercises of similar complexity. In the replication, the programming exercises, which had to be solved in Python due to the background of the participants, were organized into two blocks (A and B) that were randomly assigned to the subjects during the experiment.

In our opinion, adapting the programming language to the background of the participants should not have any impact on experimental validity, but using two blocks of exercises instead of a pool of exercises definitely improves the blocking of the related confounding variable (see Section 3.5.2), thus increasing internal validity.

4.3 Data Analysis

The data analysis was performed only for those subjects considered as valid according to the same criteria than in the original experiment. This resulted in 23 pairs, i.e. 46 valid subjects, as previously mentioned in Section 4.1.

4.3.1 Correlation Between Induced and Perceived Gender

As in the original experiment, the correlation of the induced and perceived gender in both groups was analyzed to check treatment effectiveness, especially after having changed the gendered avatars and included a gendered message at the top of the chat window, as described in Section 4.2.4.

As shown in Table 6, the MAN/MAN and WOMAN/WOMAN effectiveness was close to 40% in the replication whereas was close to 60% in the original experiment (see Table 4 in Section 3.6.1). Although Cramer's V for Table 6 showed also a *large* effect (0.530), we applied the same strict criteria than in the original experiment and decided to discard those subjects in the experimental group for whom the induced and perceived gender did not match. For the subjects in the control group, we kept those who did not perceived any gendered avatar or did not remember it. As a result, we kept 22 subjects in the control group (10 men, 12 women) but only 9 (3 men, 6 women) in the experimental group.

4.3.2 Between-Groups Analysis

As in the original experiment, and contrary to our research hypothesis, no significant differences were observed at $\alpha=0.05$ between the control and experimental groups in the replication for any of the 45 response variables¹¹ described in Section 3.4, including *cps*. The corresponding boxplots are depicted in Fig. 10. Post hoc power analyses using G*Power (Faul et al. 2007) yielded a statistical power ($1 - \beta$) for the applied tests of 0.081 for small effect sizes ($d \leq 0.2$), 0.249 for medium effect sizes ($d \leq 0.5$), and 0.536 for large effect sizes ($d \leq 0.8$).

4.3.3 Within-Groups Analysis

Within the experimental group (see Fig. 11 for the corresponding boxplots), we performed the same analysis than in the original experiment, finding statistically significant differences at $\alpha=0.05$ in the following four response variables when using the perceived partner's gender (*ppgender*) as a within-subjects variable. The four variables passed the Shapiro-Wilk normality test and were therefore analyzed using a two-sided paired t-test. Their effect sizes were computed using Cohen's *d*.

- **scd** (source code deletions): the test detected ($p = 0.0485$) that subjects deleted more source characters when they perceived their partners as a WOMAN, with a *medium* effect size ($d = -0.775$).
- **i_rf** (relative frequency of informal messages): the test detected ($p = 0.0138$) that subjects increased the relative frequency of informal messages when they perceived their partners as a MAN, with a *large* effect size ($d = 1.050$).
- **m_rf** (relative frequency of meta-comments or reflections): the test detected ($p = 0.0377$) that subjects increased the relative frequency of meta-comments or reflections when they perceived their partners as a MAN, with a *large* effect size ($d = 0.829$).
- **qyn_rf** (relative frequency of yes/no questions): the test detected ($p = 0.0297$) that subjects increased the relative frequency of yes/no questions when they perceived their partners as a MAN, with a *large* effect size ($d = 0.880$).

Note that these results must be considered carefully because of the small number of selected subjects ($n=9$), and because when *multiple test corrections*¹² are applied, only the hypothesis test corresponding to the *i_rf* variable remains significant. Post hoc power analyses using G*Power (Faul et al. 2007) yielded a statistical power ($1 - \beta$) for the applied tests of 0.137 for small effect sizes ($d \leq 0.2$), 0.393 for medium effect sizes ($d \leq 0.5$), and 0.707 for large effect sizes ($d \leq 0.8$).

No significant interactions were detected between the perceived partner's gender and the subject's gender for the same response variables than in the original study.

¹¹ Actually, only 41 variables were analyzed in the replication due to technical problems with the Python server used by the *twincod* platform. As a result, *okv*, *okv_rf*, *kov*, and *kov_rf* could not be measured and, therefore, analyzed. As can be seen in Fig. 10 and 11, their means are 0 in all cases.

¹² As recommended by de Oliveira Neto et al. (2019), multiple test corrections must be applied in Empirical Software Engineering to correct for Type I and Type II error rates. In our case, we have applied all the corrections available in the `p.adjust` function in the R language (RDocumentation 2023) such as Bonferroni-Holm, Hochberg, Hommel, Benjamini-Hochberg (also known as *false rate discovery*), or Benjamini-Yekutieli, obtaining the same results in every case.

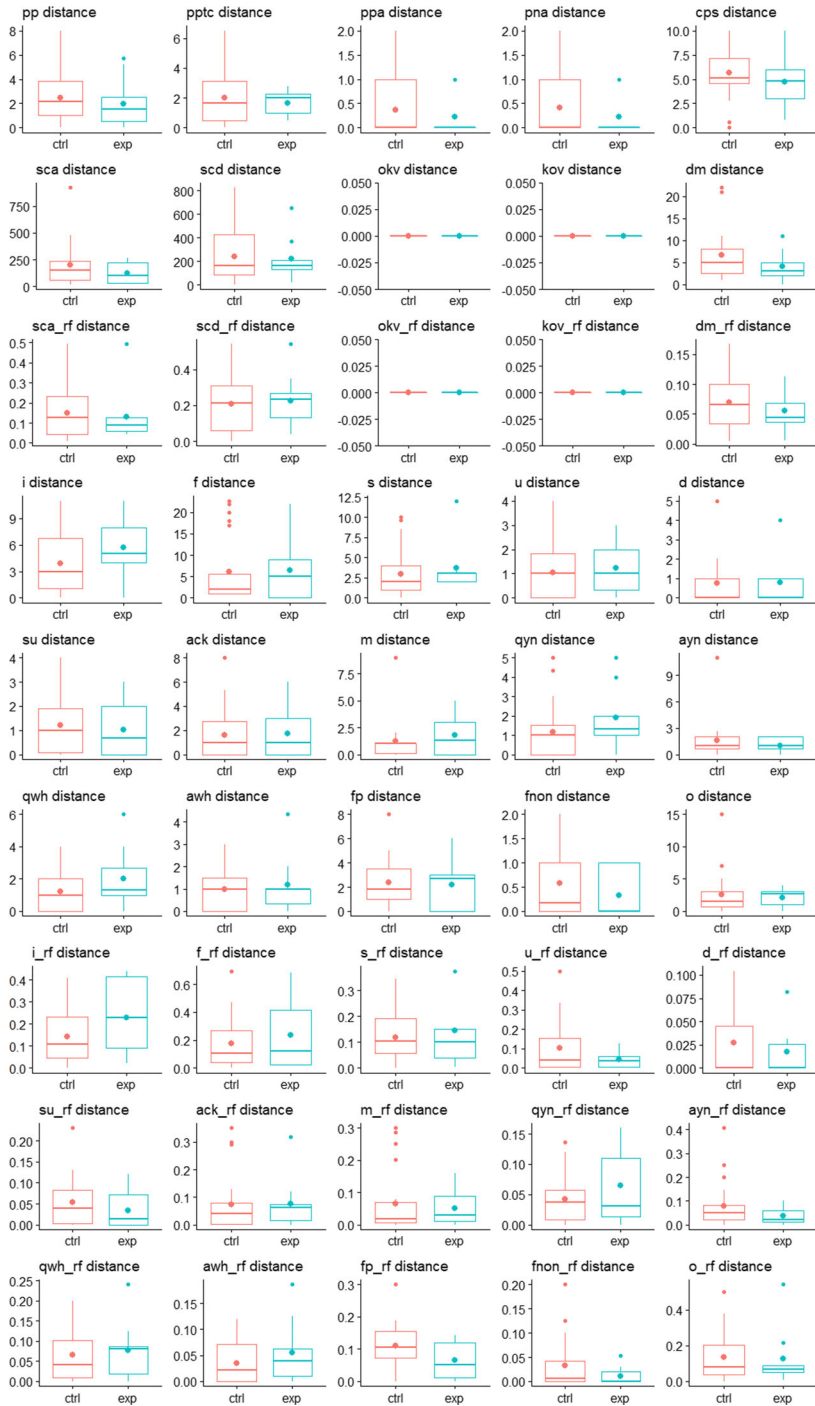


Fig. 10 Boxplots of the 45 response variables for between-groups analysis in the replication

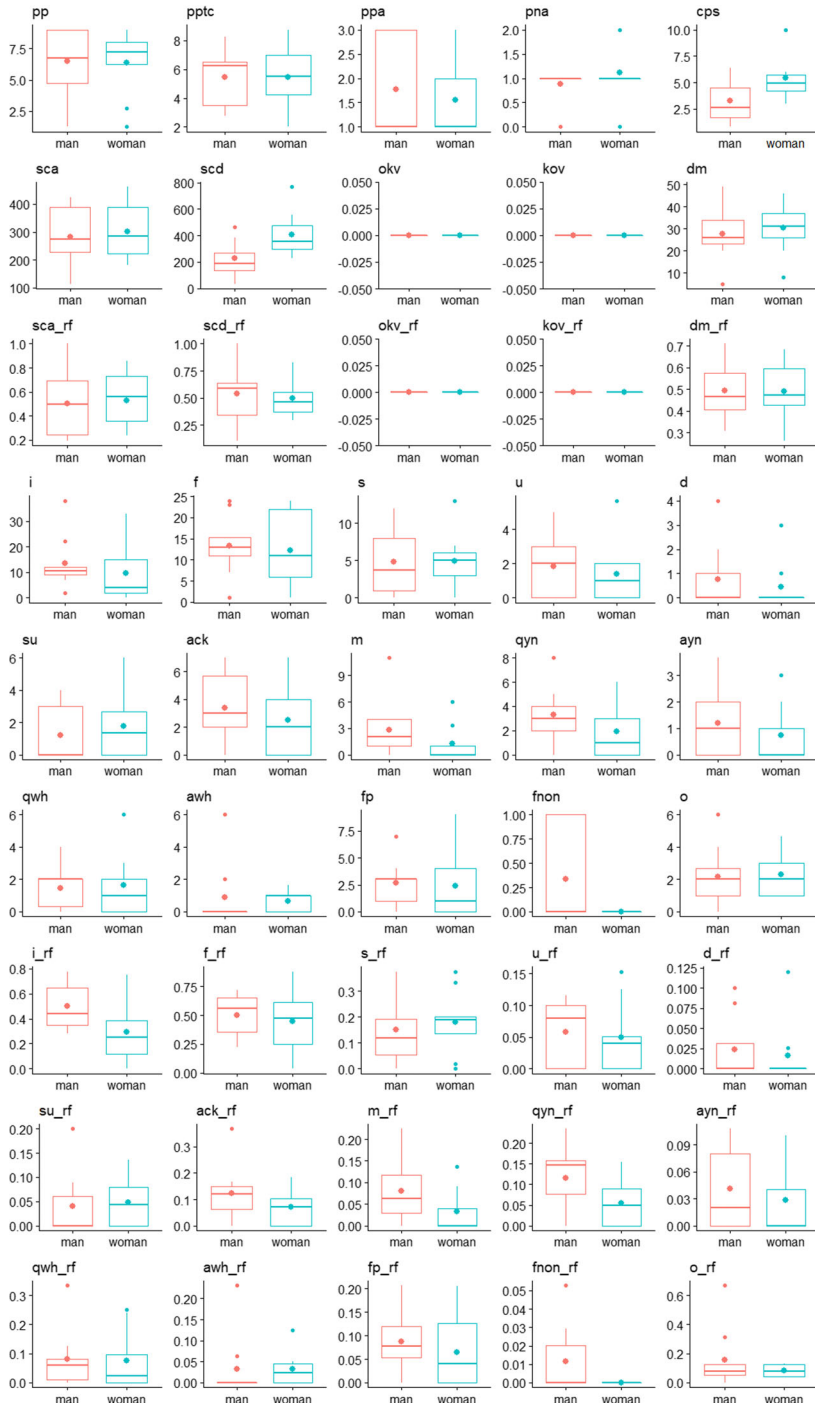


Fig. 11 Boxplots of the 45 response variables for within-groups analysis in the replication

5 Discussion and Threats to Validity

In this section, the original study and its external replication are discussed. Since the main concerns are about their threats to the experimental validity regarding operationalization and sampling, the discussion is organized around these type of threats, especially those that were not previously discussed in the description of the replication changes in Sections 4.1 and 4.2.

5.1 Operationalization of the Cause Construct — Treatment

The operationalization of gender bias into a treatment is not a trivial task and, according to the obtained results, we may not have designed our treatment as adequately as we intended, thus threatening construct validity.

Considering our experimental design, telling the subjects that they were going to collaborate with a man or a woman more explicitly could have caused in many of them the suspicion of being observed about that fact, behave unnaturally and, probably, having mentioned it unintentionally during the chat messaging, thus discovering that they were being deceived about their partner's gender and invalidating the study.

However, although the silhouetted avatars in the original experiment (see Fig. 9a) had an effectiveness close to 60% (see Table 4), when they were changed in the replication into what we thought were more explicitly gendered avatars (see Fig. 9b), their effectiveness dropped under 40% (see Table 6). Apart from the change of the avatars, this decrease in treatment effectiveness could have been probably affected by other factors, such as the remote setting, which increased the likelihood of distractions compared to a controlled environment such as a laboratory session, as commented in Section 4.2.2. Other factors could have been the reduced duration of the in-pair tasks and the second and third questionnaires, as previously discussed in Section 4.2.3, and the so-called *Zoom burnout* (Samara and Monzon 2021), i.e., the fatigue and exhaustion caused by prolonged use of video conferencing platforms during the COVID-19 pandemic, which may have influenced the motivation and performance of students at UC Berkeley, who are also exposed to very high levels of stress (Study International 2016; Newser 2023).

As commented in Section 6.2, we are evaluating the use of chatbots together with a within-subjects design in future replications to improve the treatment and thus mitigate this threat to construct validity.

5.2 Operationalization of the Effect Construct — Metrics

The main goal of our work is exploring the effects of gender bias in remote pair programming. Due to this exploratory nature, we have applied methodological triangulation (Denzin 2006), observing the phenomenon from as many points of view as possible, with an operationalization based in 45 response variables of different types which were measured during a reasonable interaction time. It is possible that some of the metrics used were not sensitive enough to treatment effects. For example, during the coding of the chat utterances in the original experiment, some of the authors, who were in their late forties and early fifties at that time, perceived strong differences between their communication and that of the significantly younger (Generation Z, Dimock (2019)) experimental subjects. These intergenerational communicative differences might have led to some noise in the labeling of chat utterances in the original experiment, although this was not the case in the replication at Berkeley, where the coders were close in age to the participants.

With all due caution, and taking into account the strong socio-political environment in Spain and the U.S. against any type of gender discrimination, we think that apart from the aforementioned limitations of the operationalization of the cause and effect constructs, another possible explanation of the obtained results is that the presence of gender bias in current Software Engineering students might not be as strong as in previous generations, although we do not have enough evidence to affirm it. In addition, if we consider that gender bias still persists in current generations—as reported by Medel and Pournaghshband (2017); Terrell et al. (2017); Allaire-Duquette et al. (2022); Oda et al. (2022)—it is also possible that most subjects self-censor, thus hindering the detection of its effects.

To improve this situation, we are currently evolving the `twincode` platform to include more metrics, and we are also considering the inclusion of qualitative research that might lead to new findings in future replications by widening the spectrum of collected information.

5.3 Sampling the Population — Participants

5.3.1 Low Percentage of Women in the Original Study

Unfortunately, the small proportion of women in STEM studies is a common issue in most higher education institutions (AAUW 2020; STEM Women 2021). The low number of women participants in the original study was an obstacle to study whether gender bias was mainly a masculine trait or if it was also present in women in any way. Nevertheless, the percentage of women increased substantially in the first replication without significant findings on the interaction of subject's gender with other factors.

5.3.2 Small Size of the Sample in the Replication

The small size of the sample in the replication and the low effectiveness of the treatment supposed a clear threat to conclusion validity that can only be mitigated by taking the outcomes as provisional and performing more replications with bigger samples and alternative experimental designs in the future.

5.3.3 Using Students as Subjects

Although in other empirical studies in which subjects are Software Engineering students, findings can be reasonably generalized to a wider community because the experimental tasks do not usually require high levels of industrial experience (Porter et al. 1999), and the students, who are the next generation of professionals, are close to the population under study (Kitchenham et al. 2002; Runeson 2003; Falessi et al. 2018), the intergenerational differences commented in Section 5.2 and the lack of conclusive results makes that very difficult in our case.

6 Conclusions and Future Work

After performing the original study and an external replication, we can conclude that we did not observe any effect of the gender bias treatment, nor any interaction between the perceived partner's gender and subject's gender, in any of the 45 response variables in the original study.

With respect to the external replication, we only observed statistically significant effects within the experimental group, i.e. comparing how subjects acted when they thought their partner was a man or a woman, in four dependent variables. One variable was related with changes in the behavior (source code deletions), and the other three were related with the relative frequency of different type of chat utterances (informal messages, reflections, and yes/no questions). In the case of the source code deletions, subjects deleted more characters when they perceived their partners as a woman, but the relative frequency of informal messages, reflections, and yes/no questions was higher when they perceived their partners as a man. We also observed a lower effectiveness of the treatment in the replication, that could be caused by the changes in the gendered avatars but also for having used a remote setting instead of a controlled environment like a laboratory session, free of distractions and interruptions.

That lower effectiveness of the treatment led to a small number of selected subjects in the experimental group, thus leading to consider the replication results inconclusive because of the small sample they are based on, and because when multiple test corrections are applied, only the result of the relative frequency of informal messages remains significant.

These outcomes have raised a number of potential research questions that we plan to address in the future and that are briefly described in the next subsections.

Take away messages

- No effect of the gender bias treatment was observed in the original study.
- Differences in four variables were observed in the replication. Only one (relative frequency of informal messages) presented statistically significant differences after multiple test corrections.
- Statistical power was low due to the reduced number of participants.
- Possible non-exclusive explanations for the results are:
 - Low treatment effectiveness: 60% in the original study (performed in a laboratory) but only 40% in the replication (performed remotely).
 - Metrics not sufficiently sensitive to the treatment effect.
 - Intergenerational communication differences in the chat utterance coding process in the original experiment.
 - Not so strong gender bias among current Software Engineering students.
 - Self-censorship of participants due to socio-political pressures.

6.1 Replication in Different Cultural Background

The cultural differences between Spanish and U.S. students could have also influenced the outcomes of both studies, so we would like to replicate it other countries and analyze those potential differences caused by cultural backgrounds.

6.2 Using Chatbots as Partners and AI-based Utterance Coding

Another two research lines we would like to explore in the future are the use of chatbots as pair programming partners and the use of deep learning to automatically code chat utterances, thus reducing the manual effort of carrying out a replication.

Inspired by current trends in Psychology (Bendig et al. 2019; Greer et al. 2019) and taking into account not only the absence of significant differences between groups in the original study and the replication, but also the difficulties in recruiting a relevant number of subjects, we are considering the possibility of changing from a between-groups design to a within-subject design in which each subject performs the pair programming tasks with a chatbot simulating being a man or a woman instead of with another human subject. Obviously, developing such a chatbot is not a trivial task, but current advances in the area, such as LaMDA (Collins and Ghahramani 2021), BERT (Devlin et al. 2019), or GPT-3 (Lim et al. 2021), make this approach a technical challenge worth exploring. A very relevant aspect in the development of such a chatbot is avoiding gender bias in the training data, as recently studied by McAuliffe et al. (2022).

On the other hand, now that we have a relevant number of coded chat utterances in Spanish and English, we could use that labeled dataset to fine train a large language model system similar to those used in chatbots to classify user intents and apply it for the automatic coding of chat utterances, which is one of the most time-consuming tasks we have had to perform as experimenters in our exploratory study. If the results of such a fine trained system were accurate, future replications would required much less effort than the two presented in this article and experimenter bias would be considerably mitigated.

Appendix

A Questionnaire #1 and #2 Response Items

In this section, the response items of the scales used in questionnaires #1 and #2 are enumerated. Those scales were analyzed for internal consistency using the data collected during the pilot studies, and the results of those analysis consisting in the Pearson's correlations, Cronbach's α , and principal components scree plot are also reported (UCLA: Statistical Consulting Group 2022), indicating whether some response items were dropped or not according to the obtained results.

A.1 Response Items for Perceived Productivity Scale (pp)

All the items in this questionnaire section, entitled as "Solo programming or pair programming?", are 0–10 numerical response items in which 0 means "programming solo", 5 means "the same in both cases", 10 means "programming in pairs".

pp₁ Regarding the programming exercises you just did, how do you think you would have been **more productive**, programming solo or programming with the partner assigned to you?

pp₂ Regarding the programming exercises you just did, how do you think you would have achieved a **better program quality**, programming solo or programming with the partner assigned to you?

pp₃ Regarding the programming exercises you just did, how do you think you would have developed a **more reliable** program, i.e., a program more likely to run without failures, programming solo or programming with the partner assigned to you?

pp₄ Regarding the programming exercises you just did, how do you think you would have **enjoyed more**, programming solo or programming with the partner assigned to you?

As shown in Fig. 12, all the items presented high Pearson correlations with Cronbach's $\alpha = 0.83$, and the scree plot confirmed they were unidimensional according to the Kaiser criterion. As a result, all of them were kept after the reliability analysis on the data from the pilot studies.

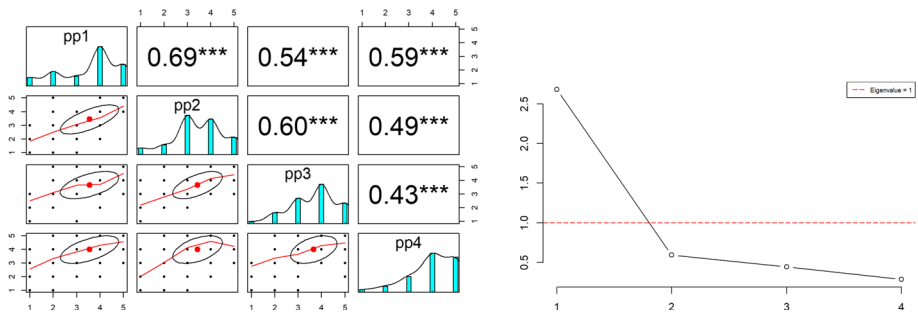


Fig. 12 Pearson correlations and scree plot of pp scale items

A.2 Response Items for Partner's Perceived Technical Competency (pptc)

All the items in this questionnaire section, entitled as “My partner or me?”, are 0–10 numerical response items in which 0 means “me”, 5 means “both equally”, 10 means “my partner”.

pptc₁ During the programming exercises you just did, who do you think had more **knowledge and technical skills**, you or the partner assigned to you?

pptc₂ During the programming exercises you just did, who do you think has been more **cooperative**, you or the partner assigned to you?

pptc₃ During the programming exercises you just did, who do you think has had a **faster pace at solving the exercises**, you or the partner assigned to you?

pptc₄ During the programming exercises you just did, who do you think has **led more to the solutions**, you or the partner assigned to you?

As shown in Fig. 13, in the initial version of the scale used in the pilot studies, the pptc₅ item, which asked whether the assigned partner had been condescending, presented low correlations with the rest of the items in the scale and the scree plot indicated two factors. After removing that uncorrelated item, the Cronbach's α increased from 0.73 to 0.85, and the scree plot indicated only one factor, as shown in Fig. 14.

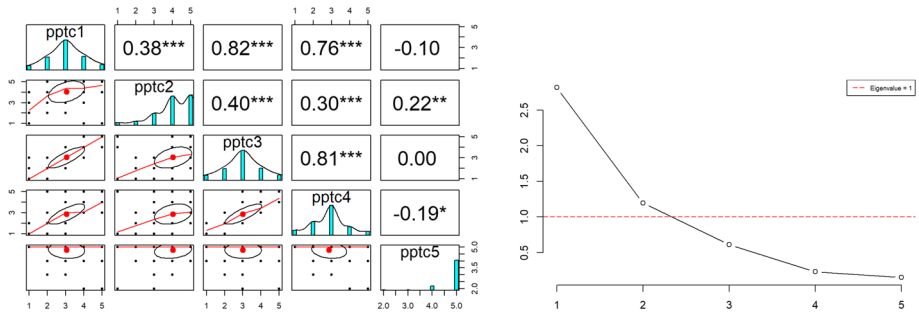


Fig. 13 Pearson correlations and scree plot of the initial version of pptc scale items

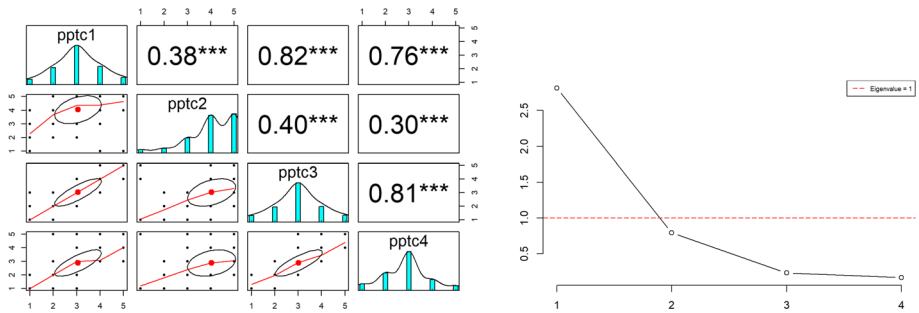


Fig. 14 Pearson correlations and scree plot after dropping pptc₅ from pptc scale

A.3 Response Item for Partner’s Perceived Positive and Negative Aspects (ppa and pna)

The only item in this questionnaire section, entitled as “Describe your partner”, is a free text field in which subjects are instructed to describe the most positive and most negative aspects of the partner assigned to them in the programming exercises they just did, indicating the positive ones with a “+” sign and the negative ones with a “-” sign in front of each aspect.

A.4 Response Items for Compared Partners’ Skills (cps)

All the items in this questionnaire section, entitled as “First or second partner?”, are 0–10 numerical response items in which 0 means “first partner”, 5 means “both equally”, 10 means “second partner”.

cps₁ Comparing your assigned partners in sessions 1 and 3, who do you think provided **more clear and constructive feedback**, your first partner or your second partner?

cps₂ Comparing your assigned partners in sessions 1 and 3, who do you think was **easier to communicate with**, your first partner or your second partner?

cps₃ Comparing your assigned partners in sessions 1 and 3, who do you think who do you think was **more knowledgeable about the subject material**, your first partner or your second partner?

cps₄ Comparing your assigned partners in sessions 1 and 3, who do you think would be a **better project partner**, your first partner or your second partner?

cps₅ Comparing your assigned partners in sessions 1 and 3, who do you think would be a **better teaching assistant**, your first partner or your second partner?

As shown in Fig. 15, all the items presented high Pearson correlations with Cronbach's $\alpha = 0.88$, and the scree plot confirmed they were unidimensional according to the Kaiser criterion. As a result, all of them were kept after the reliability analysis on the data from the pilot studies.

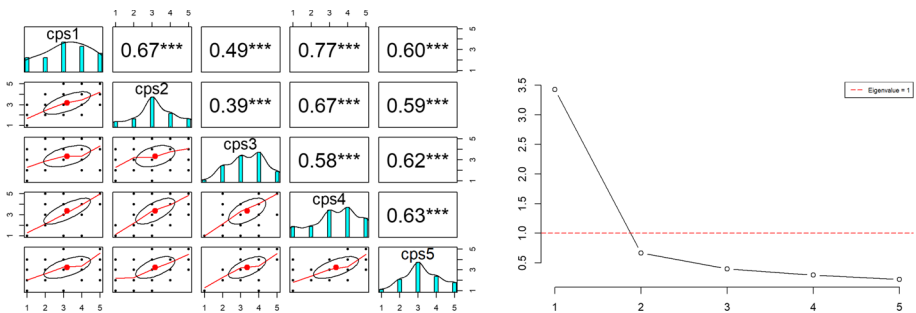
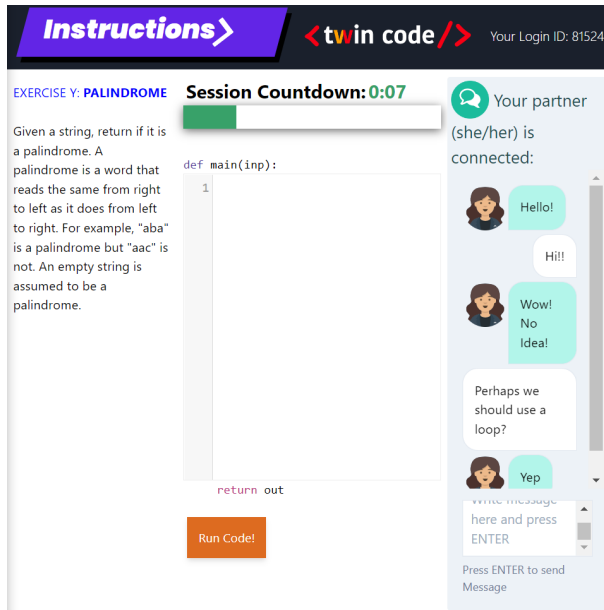


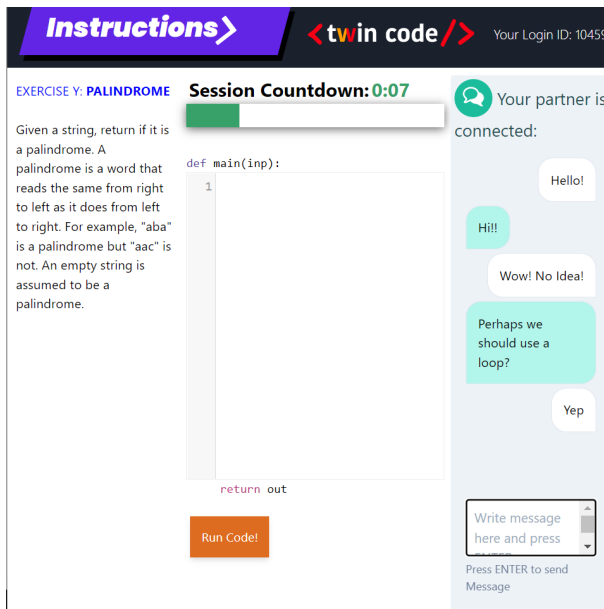
Fig. 15 Pearson correlations and scree plot of cps scale items

B Evolution of the twincode User Interface

The twincode user interface used in the external replication at UC Berkeley is shown in Fig. 16a and 16b.



(a) Experimental group — gendered avatar



(b) Control group — no avatar

Fig. 16 twincode user interface for subjects in the experimental and control groups (replication version)

C User Interface of tag-a-chat

The user interface of the tag-a-chat tool used for collaboratively coding chat utterances is shown in Fig. 17.

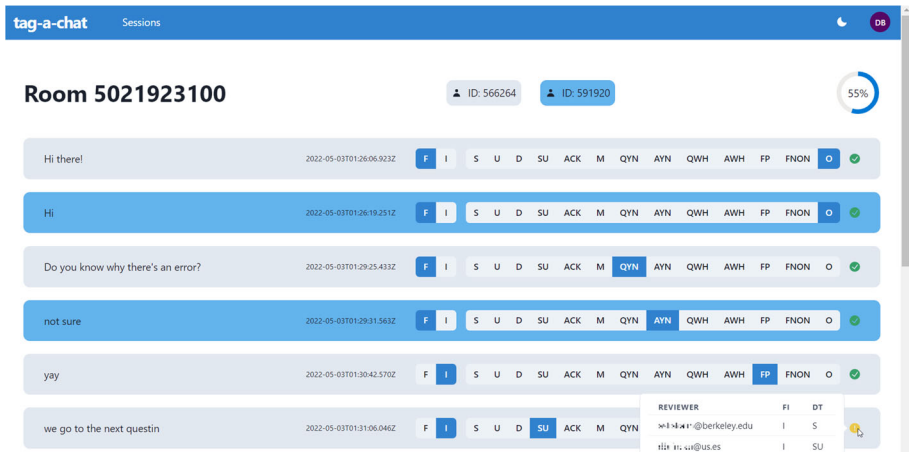


Fig. 17 User interface of the tag-a-chat tool

Acknowledgements We would like to thank the students who volunteered to participate in the pilot studies, the original experiment and the first replication at the Universities of Seville (US) and California Berkeley (UCB). We also want to thank David Brincau (undergraduate student at US) for their support in the development of the `twincod` platform; José Sandoval (Master's student at US) for developing tag-a-chat, the collaborative tool for tagging chat utterances; and Daewon Kwon and Karim el Refai (undergraduate students at UCB) for their support in the evolutive changes to the `twincod` platform and in the experiment execution at UCB. We particularly acknowledge Vron Vance (UCB alumnus, Data Analyst at Google) for their assistance regarding inclusive language around gender identity. Last but not least, we would like to thank the anonymous reviewers for their valuable comments and suggestions that helped us improve the quality and clarity of this article.

Funding Funding for open access publishing: Universidad de Sevilla/CBUA. This work has been partially supported by grants PID2021-126227NB-C21, PID2021-126227NBC22 funded by MCIN/AEI/10.13039/501100011033/FEDER and European Union "ERDF a way of making Europe"; TED2021-131023B-C21, TED2021-131023B-C22 funded by MCIN/AEI/10.13039/501100011033 and European Union "NextGenerationEU"/PRTR; EKIPMENT-PLUS (P18-FR-2895), MEMENTO (US-1381595) funded by Junta de Andalucía/ERDF, European Union; and Universidad de Sevilla under the 2021 and 2023 Grants for the Exchange Mobility of Professors, Researchers, and PhD Students between the University of Seville and the University of California.

Data Availability The datasets generated and analyzed during the current study are available in the Zenodo repository, <https://doi.org/10.5281/zenodo.6783717>.

Declarations

Conflicts of interest The authors declared that they have no conflict of interest with any aspect of the reported studies.

Ethical standard The experiment protocols were approved by the Institutional Review Board (IRB) at UC Berkeley. At the University of Seville, only studies involving experimentation animals or biomedical experiments involving humans need to be approved by the Ethics Committee on Experimentation, so no approval was required in this case.

Open Access This article is licensed under a Creative Commons Attribution 4.0 International License, which permits use, sharing, adaptation, distribution and reproduction in any medium or format, as long as you give appropriate credit to the original author(s) and the source, provide a link to the Creative Commons licence, and indicate if changes were made. The images or other third party material in this article are included in the article's Creative Commons licence, unless indicated otherwise in a credit line to the material. If material is not included in the article's Creative Commons licence and your intended use is not permitted by statutory regulation or exceeds the permitted use, you will need to obtain permission directly from the copyright holder. To view a copy of this licence, visit <http://creativecommons.org/licenses/by/4.0/>.

References

- AAUW (2002) The STEM gap: Women and girls in science, technology, engineering and mathematics. American Association of University Women. <https://www.aauw.org/resources/research/the-stem-gap/>
- Akalın A, Weinman N, Stasaski K, Fox A (2021) Exploring the impact of gender bias on pair programming. In: Proceedings of the 17th ACM conference on international computing education research, p 435–437
- Al-Jarrah A, Pontelli E (2016) On the effectiveness of a collaborative virtual pair-programming environment. In: International conference on learning and collaboration technologies, p 583–595
- Allaire-Duquette G, Chastenay P, Bouffard T, Bélanger SA, Hernandez O, Mahhou MA, Giroux P, McMullin S, Desjarlais E (2022) Gender differences in self-efficacy for programming narrowed after a 2-h science museum workshop. *Can J Sci Math Techn Educ* 22:87–100
- Bendig E, Erb B, Schulze-Thuesing L, Baumeister HH (2019) The next generation: Chatbots in clinical psychology and psychotherapy to foster mental health - a scoping review. *Verhaltenstherapie*. <https://doi.org/10.1159/000501812>
- Chaparro EA, Yuksel A, Romero P, Bryant S (2005) Factors affecting the perceived effectiveness of pair programming in higher education. In: Proceedings of the 17th workshop of the psychology of programming interest group
- Choi KS (2013) Evaluating gender significance within a pair programming context. In Proceedings of the hawaii international conference on system sciences, p 4817–4825
- Choi KS (2015) A comparative analysis of different gender pair combinations in pair programming. *Behav Inf Technol* 34(8):825–837
- Cohen L, Manion L, Morrison K (2018) *Research Methods in Education*, 8th edn. Routledge
- Collins E, Ghahramani Z (2021) LaMDA: our breakthrough conversation technology. Google Research. <https://blog.google/technology/ai/lamda/>
- Cruz M, Bernárdez B, Durán A, Guevara-Vega C, Ruiz-Cortés A (2023) A model-based approach for specifying changes in replications of empirical studies in computer science. *Computing* 105:1189–1213
- da Silva Estácio BJ, Prikladnicki R (2015) Distributed pair programming: A systematic literature review. *Inf Softw Technol* 63:1–10
- de Oliveira Neto FG, Torkar R, Feldt R, Gren L, Furia CA, Huang Z (2019) Evolution of statistical analysis in empirical software engineering research: Current state and steps forward. *J Syst Softw* 156:246–267
- Denzin NK (2006) *Sociological Methods: A Sourcebook*. 5th ed. Aldine Transaction
- Devlin J, Chang M-W, Lee K, Toutanova K (2019) BERT: Pre-training of deep bidirectional transformers for language understanding. In Proceedings of the 2019 conference of the north american chapter of the association for computational linguistics: human language technologies, Volume 1 (Long and Short Papers), p 4171–4186
- Dimock M (2019) Defining generations: Where millennials end and generation z begins. <https://pewrsr.ch/2szqtJz>
- Durán A, Fernández P, Bernárdez B, Weinman N, Akalın A, Fox A (2021) Gender bias in remote pair programming among software engineering students: The twincode exploratory study. In Proceedings of ESEM 2021 – registered report track. [arXiv:2110.01962](https://arxiv.org/abs/2110.01962)
- Eckles D, Kizilcec R, Bakshy E (2016) Estimating peer effects in networks with peer encouragement designs. *Proc Natl Acad Sci* 113(27):7316–7322

- El-Refai K, Kwon D, Brincau D, Akalın A, Fox A, Fernández P, Durán A (2023) Twincode: An instrumented platform for pair programming research. In Proceedings of the 54th ACM technical symposium on computer science education v. 2, p 1264
- Falessi D, Juristo N, Wohlin C, Turhan B, Münch J, Jedlitschka A, Oivo M (2018) Empirical software engineering experts on the use of students and professionals in experiments. *Empirical Softw Eng* 23(1):452–489
- Faul F, Erdfelder E, Lang A-G, Buchner A (2007) G*Power 3: A flexible statistical power analysis program for the social, behavioral, and biomedical sciences. *Behav Res Methods* 39:175–191
- Fisher M, Cox A (2006) Gender and programming contests: Mitigating exclusionary practices. *Inf Educ* 5(1):47–62
- Galdo AC, Celepkolu M, Lytle N, Boyer KE (2022) Pair programming in a pandemic: Understanding middle school students' remote collaboration experiences. In Proceedings of the 53rd ACM technical symposium on computer science education V. 1, p 335–341
- Gómez O, Solari M, Calvache C, Ledezma-Carrizalez A (2017) A controlled experiment on productivity of pair programming gender combinations: Preliminary results. In Proceedings of the XX Ibero-American conference on software engineering, p 197–210
- GraphPad (2023) What is the difference between ordinal, interval and ratio variables? Why should I care?. <https://t.ly/rxCW>
- Gravetter FJ, Wallnau LB (2004) *Statistics for the Behavioural Sciences*. 6 edn. Wadsworth/Thompson Learning
- Greer S, Ramo D, Chang Y-J, Fu M, Moskowitz J, Haritatos J (2019) Use of the chatbot “vivibot” to deliver positive psychology skills and promote well-being among young people after cancer treatment: Randomized controlled feasibility trial. *JMIR Mhealth Uhealth*, 7(10)
- Hanks B, Fitzgerald S, McCauley R, Murphy L, Zander C (2011) Pair programming in education: A literature review. *Comput Sci Educ* 21(2):135–173
- Hannay JE, Arisholm E, Engvik H, Sjöberg DIK (2010) Effects of personality on pair programming. *IEEE Trans Softw Eng* 36(1):61–80. <https://doi.org/10.1109/TSE.2009.41>
- Hartsell T (2005) Who's talking online? a descriptive analysis of gender & online communication. *Int J Inf Commun Technol Educ* 1(1):42–54
- Hawlitsek A, Berndt S, Schulz S (2022) Empirical research on pair programming in higher education: a literature review. *Computer science education*, p 1–29
- Hofer SI (2015) Studying gender bias in physics grading: The role of teaching experience and country. *Int J Sci Educ* 37(17):2879–2905
- Hopper J (2014) How to label your 10-point scale. Versta Research. <https://verstaresearch.com/blog/how-to-label-your-10-point-scale/>
- Jarratt L, Bowman NA, Culver KC, Segre AM (2019) A large-scale experimental study of gender and pair composition in pair programming. In Proceedings of the ACM conference on innovation and technology in computer science education, p 176–181
- Katira N, Williams L, Osborne J (2005) Towards increasing the compatibility of student pair programmers. In: International conference on software engineering, p 625–626. <https://doi.org/10.1109/ICSE.2005.1553618>
- Kaur Chahal K, Kaur A, Saini M (2021) Research and evidence in software engineering: from empirical studies to open source artifacts, chapter empirical studies on using pair programming as a pedagogical tool in higher education courses: A systematic literature review, p 251–287. Taylor & Francis Group
- Kaur Kuttal S, Gerstner K, Bejarano A (2019) Remote pair programming in online cs education: Investigating through a gender lens. In 2019 IEEE symposium on visual languages and human-centric computing (VL/HCC), p 75–85. <https://doi.org/10.1109/VLHCC.2019.8818790>
- Kitchenham BA, Pleegeer SL, Hoaglin DC, Emam KE, Rosenberg J (2002) Preliminary Guidelines for Empirical Research in Software Engineering. *IEEE Trans Softw Eng* 28(8):721–734
- Korber P, Motschnig R (2021) The effects of pair-programming in introductory programming courses with visual and text-based languages. In IEEE frontiers in education conference, p 1-9
- Lim R, Wu M, Miller L (2021) Customizing GPT-3 for your application. OpenAI. <https://openai.com/blog/customized-gpt-3/>
- Martell RF, Lane DM, Emrich C (1996) Male-female differences: A computer simulation. *Am Psychol* 51(2):157–158
- McAuliffe A, Hart J, Kuttal SK (2022) Evaluating gender bias in pair programming conversations with an agent. In 2022 IEEE symposium on visual languages and human-centric computing (VL/HCC), p 1–4. <https://doi.org/10.1109/VLHCC53370.2022.9833146>
- Medel P, Pournaghshband V (2017) Eliminating gender bias in computer science education materials. In Proceedings of the 2017 ACM SIGCSE technical symposium on computer science education, p 411–416

- Navarro D (2018) Learning statistics with R: A tutorial for psychology students and other beginners (version 0.6). <https://learningstatisticswithr.com/>
- Newser (2023) This university has the most stressed-out students. <https://www.newser.com/story/330315/10-most-least-stressed-college-towns.html>
- O'Connor C, Joffe H (2020) Intercooder reliability in qualitative research: Debates and practical guidelines. *Int J Qual Methods* 19:1–13
- Oda F, Lechago SA, da Silva BE, Hunt JC (2022) An experimental analysis of gender-biased verbal behavior and self-editing using an online chat analog. *J Exp Anal Behav* 118(1):24–45
- Per Runeson (2003) Using students as experiment subjects - an analysis on graduate and freshmen student data. In Proceedings 7th International conference on empirical assessment & evaluation in software engineering, p 95–102
- Porter AA, Votta LG, Basili VR (1999) Building Knowledge through Families of Experiments. *IEEE Trans Softw Eng* 25(4):456–473
- RDocumentation (2023) p.adjust: Adjust p-values for multiple comparisons. <https://www.rdocumentation.org/packages/stats/versions/3.6.2/topics/p.adjust>
- Rodríguez FJ, Price KM, Boyer KE (2017) Exploring the pair programming process: Characteristics of effective collaboration. In Proceedings of the 2017 ACM SIGCSE technical symposium on computer science education, p 507–512
- Saini M, Chahal KK, Kaur A (2021) Empirical studies on using pair programming as a pedagogical tool in higher education courses: A systematic literature review. Auerbach Publications
- Salleh N, Mendes E, Grundy J (2011) Empirical studies of pair programming for cs/se teaching in higher education: A systematic literature review. *IEEE Trans Software Eng* 37:509–525. <https://doi.org/10.1109/TSE.2010.59>
- Salleh N, Mendes E, Grundy J (2014) Investigating the effects of personality traits on pair programming in a higher education setting through a family of experiments. *Empirical Soft Eng* 19(3):714–752
- Salleh N, Mendes E, Grundy J, Burch G (2010) The effects of neuroticism on pair programming: an empirical study in the higher education context. In Proceedings of the 2010 ACM-IEEE international symposium on empirical software engineering and measurement, p 1–10
- Samara O, Monzon A (2021) Zoom burnout amidst a pandemic: Perspective from a medical student and learner. *Therapeutic Advances in Infectious Disease*, 8
- Sfetsos P, Stamelos I, Angelis L, Deligiannis I (2009) An experimental investigation of personality types impact on pair effectiveness in pair programming. *Empirical Softw Eng* 14(2):187–226
- STEM Women (2021) Percentages of women in STEM statistics. *STEM Women*. <https://www.stemwomen.com/women-in-stem-percentages-of-women-in-stem-statistics>
- Stevens SS (1946) On the theory of scales of measurement. *Science* 103(2684):677–680
- Stotts D, Williams L, Nagappan N, Baheti P, Jen D, Jackson A (2003) Virtual teaming: Experiments and experiences with distributed pair programming. In: Conference on extreme programming and agile methods, p 129–141
- Study International (2016) Students at these U.S. universities are under the most stress. <https://www.studyinternational.com/news/students-mental-health-us-universities-stress/>
- Syed M, Nelson SC (2015) Guidelines for establishing reliability when coding narrative data. *Emerging Adulthood* 3(6):375–387
- Terrell J, Kofink A, Middleton J, Rainear C, Murphy-Hill E, Parnin C, Stallings J (2017) Gender differences and bias in open source: pull request acceptance of women versus men. *PeerJ Computer Science*, 3(e111)
- Thomas L, Ratcliffe M, Robertson A (2003) Code warriors and code-a-phobes: A study in attitude and pair programming. In Proceedings of SIGCSE, p 363–367
- UCLA: Statistical Consulting Group (2022) What does cronbach's alpha mean?. Accessed 29-June-2022. <https://stats.oarc.ucla.edu/spss/faq/what-does-cronbachs-alpha-mean/>
- University of California, Berkeley (2021) Demographic information (restricted access). <https://calanswers.berkeley.edu/home>
- University of Seville (2021) Statistical yearbook 2020–2021. <https://servicio.us.es/splanestu/WS/Anuario2021/AESY20-21.html>. English version starts at page 400
- Werner LL, Hanks B, McDowell C (2004) Pair-programming helps female computer science students. *J Educ Resour Comput*, 4(1)
- Wohlin C, Runeson P, Höst M, Ohlsson MC, Regnell B, Wesslén A (2012) Experimentation in Software Engineering: an Introduction. Springer
- Xinogalos S, Satratzemi M, Chatzigeorgiou A, Tsompanoudi D (2017) Student perceptions on the benefits and shortcomings of distributed pair programming assignments. 2017 IEEE global engineering education conference (EDUCON), p 1513–152

- Ying KM, Martin AC, Rodríguez FJ, Boyer KE (2021a) Cs1 students' perspectives on the computer science gender gap: Achieving equity requires awareness. In 2021 Conference on research in equitable and sustained participation in engineering, computing, and technology (RESPECT), p 1–9. IEEE
- Ying KM, Rodríguez FJ, Dibble AL, Boyer KE (2021) Understanding women's remote collaborative programming experiences: The relationship between dialogue features and reported perceptions. *Proc ACM Hum -Comput Interact* 4(CSCW3):1–29

Publisher's Note Springer Nature remains neutral with regard to jurisdictional claims in published maps and institutional affiliations.



Amador Durán Toro is an associate professor at the University of Seville, Spain. His current research is focused on empirical software engineering, requirements engineering, software testing, business process modeling, service engineering, and applications of artificial intelligence. He is the author of the [requirements management tool REM](#), used by universities and companies in several countries. He also serves regularly as a reviewer for relevant journals and conferences.



Pablo Fernández is an associate professor at the University of Seville, Spain. His current research is focused on the automated governance of organizations based on service level agreements and commitments. He has been the lead architect for several projects in public administrations and major firms.



Beatriz Bernárdez is an assistant professor with the University of Seville, Spain. Her current research focuses on requirements engineering, pair programming and the application of mindfulness to software engineering from an empirical perspective, having developed a number of experiments in these areas. She has the Search Inside Yourself Certification of the personal growth program developed by Google. She serves as a reviewer for relevant journals and has collaborated in the organization of international conferences such as ICSE'2021.



Nathaniel Weinman received his Ph.D. degree in Computer Science from UC Berkeley in 2023. His research lies at the intersection of human-computer interaction, education, and software engineering. He has helped develop materials for multiple computer science courses at Berkeley and received the Outstanding Graduate Student Instructor award in 2022.







Aslihan Akalın received her Masters degree in Computer Science from UC Berkeley in 2023. Her academic interests center around CS education and learning (humans and machines) with a special interest in natural language processing. She has served as a student instructor for many Computer Science courses, and received the Outstanding Graduate Student Instructor award for her work in Summer 2022.



Armando Fox is Professor of Computer Science, Faculty Advisor for Digital Learning Strategy, and Equity and Diversity Advisor UC Berkeley. He and colleague David Patterson co-created Berkeley's first Massive Open Online Course, "Engineering Software as a Service", offered through edX, and the award-winning accompanying textbook of the same name, for which he received the ACM Karl V. Karlstrom Outstanding Educator Award in 2015. His CS education research focuses on creating novel technologies to help students learn advanced programming concepts at scale.

Authors and Affiliations

Amador Durán Toro^{1,2}  · **Pablo Fernández**^{1,2}  · **Beatriz Bernárdez**^{1,2}  ·
Nathaniel Weinman³ · **Aslıhan Akalın**³ · **Armando Fox**³ 

Pablo Fernández
pablofm@us.es

Beatriz Bernárdez
beat@us.es

Nathaniel Weinman
nweinman@berkeley.edu

Aslıhan Akalın
asliakalin@berkeley.edu

Armando Fox
fox@berkeley.edu

¹ I3US Institute, Universidad de Sevilla, Seville, Spain

² SCORE Lab, Universidad de Sevilla, Seville, Spain

³ Computer Science Division, University of California, Berkeley, CA, USA