

Test-Driven Development and Embedded Systems: An Exploratory Investigation

Michelangelo Esposito, Simone Romano, and Giuseppe Scanniello

University of Salerno, Fisciano, Italy,

{m.esposito253@studenti.unisa.it} and {siromano, gscanniello}@unisa.it

Abstract. We present the results of an exploratory investigation to obtain preliminary evidence on the use of *Test-Driven Development (TDD)*, an incremental approach to software development where tests are written before production code, to develop *Embedded Systems (ESs)*. Specifically, we conducted two experiments in which we compared TDD with a non-TDD approach in terms of the external quality of ESs and developers' productivity. In the experiments, we also gathered qualitative data to better understand the investigated phenomenon. We found that the external quality of the implemented solutions increases when using TDD as compared to a non-TDD approach, while there is not a substantial difference with respect to developers' productivity. However, TDD is perceived as more difficult to apply than a non-TDD approach, and the development task is deemed more challenging with TDD.

Keywords: Embedded Systems, Experiment, Replication, TDD

1 Introduction

An *Embedded System (ES)* is a combination of hardware and software created for a particular purpose. In many cases, ESs operate as part of a larger system (*e.g.*, agricultural-sector equipment, automobiles, medical equipment, airplanes, *etc.*). The global market of ESs is expected to witness notable growth. A recent report evaluated the global market of ESs \$89.1 billion in 2021, and this market is projected to reach \$163.2 billion by 2031; with a compound annual growth rate of 6.5% [1]. This growth is related to an increase in the number of ESs-related research and advanced development projects (*e.g.*, driver-assistance systems) [1].

To date, there has yet to be shown which approach is to be preferred when developing ESs. In his book, Greening [12] asserted that developers can benefit from the application of *Test-Driven Development (TDD)*, an incremental approach to software development in which the developer repeats a short cycle, to incrementally implement functionality, made up of three phases: *Red*, *Green*, and *Refactor* [5]. In the Red phase, the developer writes a unit test for a small chunk of functionality not yet implemented and watches the test fail. In the Green phase, the developer implements that chunk of functionality as quickly as possible and watches all unit tests pass. In the Refactor phase, the developer changes the internal structure of the code while paying attention to keeping the

functionality intact (*i.e.*, all unit tests have to pass). The Red-Green-Refactor cycle is thus repeated until the functionality is completely implemented. When this happens, the developer can tackle new functionality.

TDD has been conceived to develop “traditional” software and it is claimed to improve software quality as well as developers’ productivity [14]. ESs have all challenges of non-Embedded Systems (NoESs), such as poor quality, but add challenges of their own [12]. For example, one of the most cited differences between ESs and NoESs is that the code of ESs depends on the hardware.

A huge amount of empirical investigations has been conducted to study the claimed effects of TDD when developing NoESs [14]. On the other hand, no investigation has been conducted yet to assess the benefits that TDD could bring to the development of ESs although some authors, like Greening [12], claim that the application of this approach produces benefits in the context of ESs development (*e.g.*, high-quality ESs).

To increase the body of knowledge on the application of TDD to ESs development, we investigated the following high-level Research Question (RQ):

Does the use of TDD affect the external quality of ESs, as well as developers’ productivity?

To answer this RQ, we conducted an exploratory study made up of two experiments (*i.e.*, a baseline experiment and a replication) on TDD applied to the development of ESs. The participants in the study were final-year Master’s students in Computer Science (CS) taking an ESs course at the University of Salerno (Italy). To assess the benefits (if any) due to the application of TDD in terms of external (*i.e.*, functional) quality of developed solutions and developers’ productivity, we compared TDD with a non-TDD approach named *YW* (*Your-Way*)—*i.e.*, the approach developers would normally use to produce software when no restriction is imposed to them (except for not using TDD, of course) [2, 3, 11]. In both experiments, whichever the used approach (TDD or YW) was to develop a given ES, the participants were asked to use *mocks* to model and confirm the interactions between the device driver and the hardware. The mocks thus intercepted the commands to and from the device by simulating specific usage scenarios [12]. In the replication, we also asked the participants to replace the mocks with the actual interactions with sensors and actuators by deploying and testing their ES in the target hardware. In both experiments, we also gathered qualitative data to better understand the phenomenon under investigation.

Paper Structure. In Section 2, we outline related work. The design of our investigation is described in Section 3. The results are presented and discussed in Section 4 and Section 5, respectively. We highlight the threats to the validity of our investigation in Section 6. Conclusion and final remarks ends the paper.

2 Related Work

Despite different authors (*e.g.*, [9, 12, 15, 27]) have promoted the use of TDD in the context of ESs development, nobody has ever investigated the benefits that

TDD could bring to the development of ESs. Conversely, in the context of NoESs development, the claimed benefits of TDD—including increased external quality and developers’ productivity, which we investigated in our study—have been the subject of different primary studies (*e.g.*, [10, 23]), whose results have been also gathered and combined in secondary studies (*e.g.*, [6, 19, 20, 24]). For example, in their Systematic Literature Review (SLR), Turhan *et al.* [24] included 32 primary studies (2000-2009) and found a moderate effect in favor of TDD on external quality while the results about developers’ productivity are inconclusive (*i.e.*, the results do not lead to a firm conclusion in favor or against TDD). Bissi *et al.* [6] conducted an SLR that included 27 primary studies (1999-2014). The authors observed, similar to Turhan *et al.* [24], an improvement of external quality due to TDD, while the results about productivity are inconclusive. Rafique and Misic [20] conducted a meta-analysis of 25 controlled experiments (2000-2011) and observed a small effect in favor of TDD on external quality, while again the results about productivity are inconclusive. Finally, Munir *et al.* [19], in their SLR, classified 41 primary studies (2000-2011) into four categories based on rigor and relevance. They found that in each category different conclusions could be drawn for both external quality and productivity.

Our investigation, although it is preliminary, has the merit to study for the first time the application of TDD in a new development context, the one of ESs, in which different authors have promoted its application (*e.g.*, [9, 12, 15, 27]). Given the preliminary nature of our study, the obtained results are not intended to be conclusive; rather, we aim to gather initial evidence on the application of TDD when developing ESs and pave the way for future research on this matter.

3 Baseline Experiment and Replication

In the following of this section, we detail the design of our exploratory study, which comprises a baseline experiment (**Exp1**) and a replication (**Exp2**) with the same group of participants. To design our study, we took into account the guidelines for experimentation in Software Engineering by Wohlin *et al.* [26].

3.1 Research Questions

Consistent with our high-level RQ, the goal of our study (formulated according to the *GQM* template [4]) is the following:

Analyze the use of TDD **for the purpose of** evaluating its effects in the development of ESs **with respect to** the external quality of the developed solutions and developers’ productivity **from the point of view of** researchers and lecturers **in the context of** an ESs course involving final-year Master’s students in CS who develop ESs in Python.

Accordingly, we defined and investigated the following (low-level) RQs:

RQ1. Does the use of TDD improve the external quality of ESs whose software is written in Python?

RQ2. Does the use of TDD increase developers’ productivity when coding in Python?

We focused our study on Python because it (including its variants like MicroPython) in IEEE Spectrum¹ has ranked as the first programming languages for ESs, Web and Enterprise development.

3.2 Participants

The participants in our study (*i.e.*, Exp1 and Exp2) were final-year Master’s students in CS at the University of Salerno. They were sampled by convenience among the students taking an ES course. Nine students, out of ten, accepted to participate in the study on a voluntary basis. In line with Carver *et al.*’s advice [7], we rewarded the students for their participation with two bonus points on the final mark of the course. The students were aware that: *(i)* they would receive the bonus points independently of their performance in the study; *(ii)* they could drop out of the study at any time, without being negatively judged; *(iii)* they could achieve the highest mark of the course even without participating in the study; and *(iv)* all the gathered data would be confidentially treated and anonymously shared for research purposes only. Our study had both research and educational goals: on one hand, we conceived the study to answer our RQs; on the other hand, the study allowed the students to gain experience with TDD applied to ESs development. As for the educational goal, TDD has attracted great attention from developers [11, 22] and it is claimed to bring benefits to the development of ESs [12].

The participants had all programming experience, and most of them rated such an experience 3 on a 5-point Likert scale (where 1 means “very inexperienced” and 5 means “very experienced”). The participants were mostly not experienced with unit testing since most of them rated 1 or 2 their unit-testing experience (again on a 5-point Likert scale). Three students had heard about TDD but nobody had a practical experience with it.

3.3 Experimental Tasks

The participants had to fulfill three experimental tasks: two in Exp1 and one in Exp2. Each experimental task consisted of developing an ES for *Raspberry Pi* by coding in Python and using *unittest* and *Mock.GPIO* as testing and mocking frameworks, respectively. In Exp1, the experimental tasks were: **Intelligent Office (IO)** and **Cleaning Robot (CB)**. As for Exp2, the experimental task was **Smart Room (SR)**. A brief description of the experimental tasks is reported in Table 1 (further details are available in our online replication package [21]).

For each task, the experimental material provided to the participants included: *(i)* its description, which consisted of an abstract of the ES to be developed, a list of instructions, and a set of user’s stories describing the ES’s

¹ <https://spectrum.ieee.org/top-programming-languages-2022>

Table 1: Experimental task description.

Name	Description
IO	Developing an intelligent office system (named IO), which allows handling light and CO2 levels inside an office. To handle the light, IO relies on the information gathered from different sensors— <i>i.e.</i> , Infra-Red (IR) distance sensors, a real-time clock, and a photoresistor—and then controls a servo motor (to open/close the blinds) and a light bulb. A carbon dioxide sensor is used to monitor the CO2 levels inside the office and then control the switch of an exhaust fan.
CB	Developing a cleaning robot (named CB) that, while it moves in a room based on the commands it receives from an external system, cleans the dust on the floor along the way (by using rotating brushes). The robot also detects potential obstacles through an IR distance sensor. An intelligent battery sensor is used to check the charge left of the internal battery of the robot, and a LED is turned on to signal the need for a recharge.
SR	Developing a smart room system (named SR), which handles light, temperature, and gas levels inside a room. To handle the light, the system relies on the information gathered from an IR distance sensor and photoresistor, and then turned on/off a light bulb. The information from a pair of temperature sensors is used to control the servo motor of a window. SR is also equipped with an air quality sensor to measure the gas levels inside the room and, if necessary, it triggers an alarm through a buzzer.

functionality to be implemented; and *(ii)* a project template (for *PyCharm* the used IDE), which contained function stubs (*i.e.*, empty functions exposing the expected API signatures), utility functions, and an example unittest test class.

Independently from the used development approach (TDD or YW) and experiment (Exp1 or Exp2), the participants when implementing an ES had to use mocks to model and confirm the interactions between the device driver and the hardware—this is similar to what Grenning suggests in his book [12]. In other words, the mocks allowed the unit tests written by the participants to simulate specific scenarios without depending on the actual sensors and actuators. We deliberately introduced a difference in Exp2 to take into account further steps of the development pipeline for ESs. In particular, we asked the participants in Exp2 (whatever the used development approach was) to apply two additional steps in the development pipeline: *(i)* replace the mocks with the actual interactions with sensors and actuators and then deploy their code in a Raspberry Pi; and *(ii)* test the ES on the field.

3.4 Independent and Dependent Variables

Regardless of the experiment, the participants were asked to fulfill each experimental task by using either TDD or YW as a development approach. Therefore, the independent variable of both Exp1 and Exp2 is **Approach**. It is a nominal variable assuming two possible values: *TDD* and *YW*. The choice of using YW as a baseline for comparison is based on past studies on TDD in the context of NoESs development (*e.g.*, [2, 3, 11]).

In both Exp1 and Exp2, the dependent variables are **QLTY** and **PROD**. The definitions of these variables are borrowed from past studies on TDD, but applied to NoESs development (*e.g.*, [2, 10, 23]). **QLTY** measures the external

quality and, similar to the past studies mentioned above, it is defined as follows:

$$QLTY = \frac{\sum_{i=1}^{\#TUS} QLTY_i}{\#TUS} * 100 \quad (1)$$

where $\#TUS$ is the number of user stories a participant tackled, while $QLTY_i$ is the external quality of the i -th user story tackled. A user story is tackled if at least one assert in the acceptance test suite of that user story passes. The acceptance test suites (one per user story) were pre-defined by the authors of this paper and not delivered to the participants. As for $QLTY_i$, it is computed as the number of asserts passed for the i -th tackled user story out of the total number of asserts for that user story:

$$QLTY_i = \frac{\#ASSERT_i(PASS)}{\#ASSERT_i(ALL)} \quad (2)$$

QLTY assumes values between 0 and 100, where a value close to 100 indicates high external quality of the developed solutions.

As for the PROD variable, it measures the productivity of a participant when developing an ES. Similar to past studies [2, 10, 23], PROD is defined as follows:

$$PROD = \frac{\#ASSERT(PASS)}{\#ASSERT(ALL)} * 100 \quad (3)$$

where $\#ASSERT(PASS)$ is the number of asserts passed in the acceptance test suites, while $\#ASSERT(ALL)$ is the total number of asserts in the acceptance test suites. PROD assumes values between 0 and 100, where a value close to 100 means high productivity of developers.

3.5 Experimental Designs

The experimental design of Exp1 is *ABBA crossover* [25]. It is a kind of *within-participants* design where each participant receives both treatments (*i.e.*, A and B). In ABBA crossover designs, there are two *sequences* (*i.e.*, AB and BA), defined as the order with which the treatments are administered to the participants, and two *periods* (*i.e.*, *Period1* and *Period2*), defined as the times at which each treatment is administered. The experimental groups correspond to the sequences. Also, to mitigate learning effects, each period is paired with a different experimental task. In our case, the experimental groups are *TDD-YW* and *YW-TDD*. As shown in Table 2, the former first experimented with TDD to fulfill IO and then YW to fulfill CR, while the latter first applied YW and then TDD to fulfill IO and CR, respectively. The assignment of the participants to the experimental groups was done randomly.

As for Exp2, the experimental design is *one-factor-two-treatments* [26], a kind of *between-participants* design. In such a design, each participant receives only one treatment (in our case, either TDD or YW) while dealing with the same experimental task (in our case, SR). The experimental groups are paired with the treatments. The assignment of the participants to the groups was done

Table 2: Assignment of the participants.

Approach	Exp1		Exp2 (SR)
	Period1 (IO)	Period2 (CR)	
TDD	P1-P4	P5-P9	P2-P3, P6-P8
YW	P5-P9	P1-P4	P1, P4-P5, P9

randomly: five participants were assigned to TDD and four to YW (see Table 2). We would like to mention that the variation in the design of Exp2, as compared with the design of Exp1, was introduced because of the teaching schedule—*i.e.*, we exhausted the class hours reserved for the practical part of the ESs course.

3.6 Experimental Procedure

The experimental procedure of our study consisted of the following steps.

- 1. Recruitment.** We gathered the adhesion of the students to participate in the experiments through an online questionnaire. It was also used to gather demographic information on the participants (*e.g.*, their programming experience).
- 2. Training.** All the participants attended a lesson about unit testing in Python with unittest as a testing framework. The first part of the lesson was theoretical, while the second part was practical—in particular, under the lecturer’s supervision, the students performed the unit testing of a NoES. Then, the participants attended a lesson about TDD, mocking (including the Mock.GPIO framework), and how to program Raspberry Pi. Again, the first part of the lesson was theoretical, while the second one was practical—*i.e.*, the students, with the guide of the lecturer, applied TDD to develop an ES for managing a greenhouse. Finally, as a laboratory activity, the students took part in a *warm-up task* where they all used TDD to develop an ES for managing a parking garage.
- 3. Exp1 Execution.** We (randomly) split the participants into the TDD-YW and YW-TDD groups. As shown in Section 3.5, each participant applied alone each approach only once to fulfill IO and CR. The participants tackled the experimental tasks during two laboratory sessions, each lasting two hours, on two different days. At the beginning of each laboratory session, the participants received the experimental material (see Section 3.3), imported the project template into the PyCharm IDE, and then started implementing the user stories of the experimental task once at a time by using the requested approach. We asked the participants to use mocks. At the end of each laboratory session, the participants filled out an online post-questionnaire by which they delivered their project and shared feedback about both the experimental task and the used approach.
- 4. Exp2 Execution.** We (randomly) split the participants into the TDD and YW groups. The participants tackled alone the experimental task at home and each participant applied either TDD or YW to fulfill SR (see Section 3.5). The experimental material was delivered via email; once received, the participants could import the project template into the PyCharm IDE and then start implementing the user stories of the experimental task once at a time by using the requested approach. To fulfill the experimental task and similar to Exp1, the

participants used mocks without interacting with the actual sensors/actuators. By the delivery date, the participants had to fill out an online post-questionnaire to deliver their PyCharm project and then book a slot for deploying and testing SR in the target hardware (pre-assembled in a research laboratory at the University of Salerno). The deployment of SR required each participant to replace the mocks with the actual interactions with sensors and actuators (see Section 3.3). The participants deployed and tested SR once at a time; in the end, they were (individually) interviewed by one of the authors of this paper. We conducted the interviews to gather feedback from the participants about the entire study. We asked the participants for their consent to audio-record the interviews.

3.7 Data Analysis

We used boxplots to depict the distributions of the values of each dependent variable across the experimental tasks of Exp1 and Exp2 (we also plotted the mean values in these boxplots). Later, we aggregated the data from both experiments by performing a meta-analysis—to perform it, we took into account the guideline by Kitchenham *et al.* [17]. For each dependent variable and experiment, we computed the Standardized Mean Difference (SMD) as *Hedges' (adjusted) g*.² To obtain joint SMDs (one for each dependent variable), we leveraged random-effects meta-analysis models. Finally, we show the results of the meta-analysis through forest plots.³ We did not perform any statistical inference given the number of participants and exploratory nature of our study.

As for the post-questionnaires, we used diverging stacked bar plots to summarize the answers to the closed questions (reported as statements to be rated on 5-point Likert scales) by approach. To analyze the interview data instead, we first transcribed the interview audio recordings and then identified common themes in the interview transcripts by applying a thematic analysis. We took into account the guidelines by King [16] to identify such common themes. Briefly, we exploited the interview script (we used to guide the interviews) to develop some initial themes, and then we revised these themes as the analysis of the interview transcripts progressed until reaching a consensus on the final themes.

Analysis scripts, raw data, and more are available online [21].

4 Results

In this section, we first show the results regarding RQ1 and RQ2, and then those from the post-questionnaires and interviews, respectively.

² We computed any SMD from Exp1 as suggested by Madeyski and Kitchenham [18]; this is to make the SMD from a crossover experiment (Exp1) comparable with that from a between-participants experiment (Exp2).

³ The values of the descriptive statistics are available online [21], as well as the values of the I^2 statistic giving an indication of the between-experiment heterogeneity.

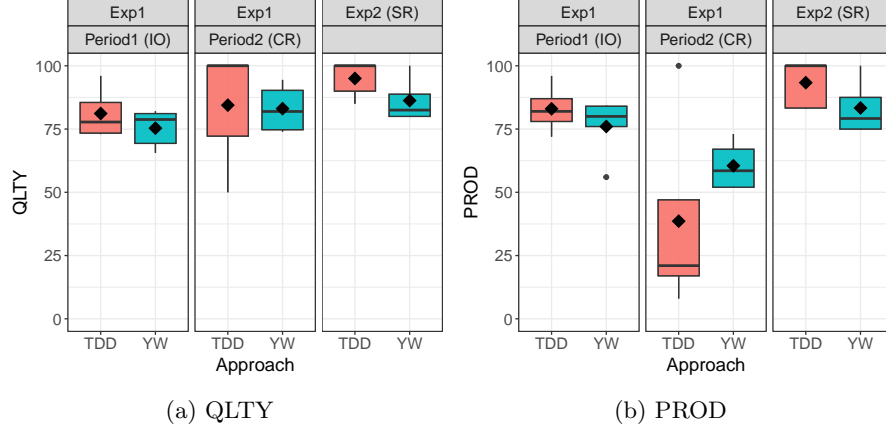


Fig. 1: Boxplots with respect to QLTY and PROD.

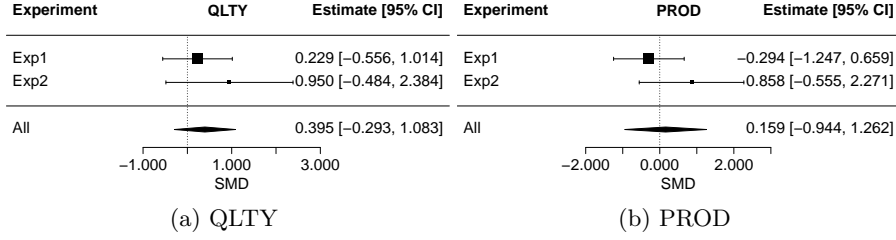


Fig. 2: Forest plots with respect to QLTY and PROD.

4.1 Results on RQ1 and RQ2

In Figure 1, we show the boxplots depicting the distributions of the values of QLTY and PROD for TDD and YW across the experimental tasks (*i.e.*, IO and CR in Exp1, and SR in Exp2). The forest plots depicting, for QLTY and PROD, the SMDs across the experiments and the joint SMD are shown in Figure 2.

QLTY. The comparison with respect to QLTY between TDD and YW seems in favor of the former for any experimental task. Indeed, by looking at Figure 1.a, we can notice that the boxes for TDD are either higher than or comparable to the ones for YW. This trend is confirmed by the mean and median values (represented in the boxplots as diamonds and thick horizontal lines, respectively). We can also notice that the effect of TDD on QLTY is not uniform across the experimental tasks: the difference between TDD and YW is wider for the experimental task in Exp2, while it is more limited for the two experimental tasks in Exp1. This suggests a moderating role of the experimental task. As for the SMDs, we can notice in Figure 2.a that, in the single experiments, the SMDs are both in favor of TDD. In particular, the SMD is small⁴ (0.229) in Exp1 and large (0.95) in Exp2. The joint SMD is in favor of TDD and it is small (0.395).

⁴ Based on Cohen's guidelines [8], the SMD can be interpreted as: *negligible*, if $|\text{SMD}| < 0.2$; *small*, if $0.2 \leq |\text{SMD}| < 0.5$; *medium*, if $0.5 \leq |\text{SMD}| < 0.8$; or *large*, otherwise.

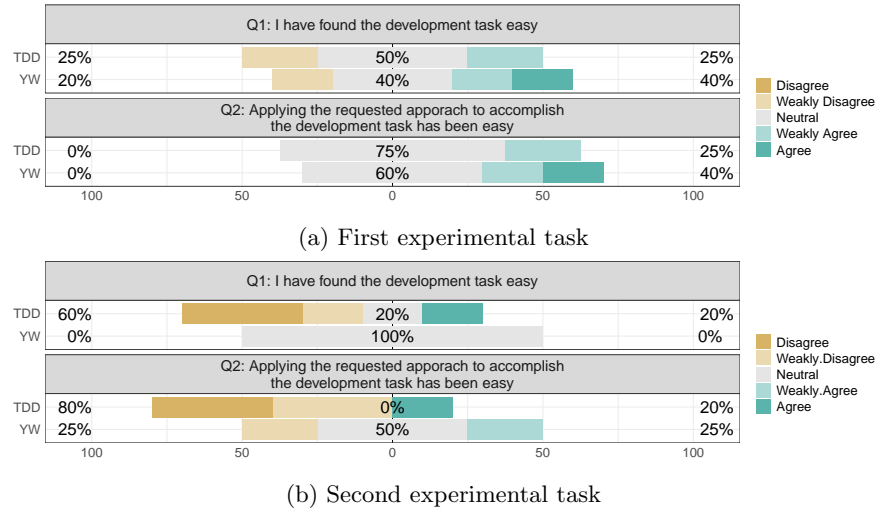


Fig. 3: Diverging stacked bar plots summarizing the answers to the closed questions of the post-questionnaires in Exp1 for first and second experimental tasks.

PROD. The effect of TDD on PROD is contrasting as shown in Figure 1.b. In particular, when considering the first experimental task in Exp1 and the one in Exp2, the boxes for TDD are either higher than the boxes for YW. On the contrary, the box for TDD is lower than the one for YW in the second experimental task of Exp1. Such a contrasting trend is confirmed when looking at the mean and median values, so suggesting a moderating role of the experimental task. As for the SMDs depicted in Figure 2.b, in the single experiments, the SMDs (with respect to PROD) are one in favor of YW and one in favor of TDD (Exp2). In particular, the SMD of Exp1 is in favor of YW and small (-0.294)—this is due to the second experimental task, as shown in Figure 2.b. On the contrary, the SMD of Exp2 is in favor of TDD and it is large (0.858). The joint SMD is still in favor of TDD, but negligible (0.159).

4.2 Results from Post-questionnaires

In Figure 3, we show the diverging stacked bar plots summarizing the answers to the closed questions of the post-questionnaires the participants filled out at the end of the two experimental tasks in Exp1. It seems that the participants who applied TDD found the experimental tasks less easy than those who applied YW. Indeed, the percentages of agreement in the first experimental task are equal to 25% for TDD and 40% for YW (see Figure 3.a). Such a difference seems stronger when considering the second experimental task (see Figure 3.b): the percentage of disagreement for TDD is equal to 60% while the one for YW is equal to 0%.

TDD seems to be considered less easy to be applied than YW in the first experimental task (see Figure 3.a): the percentages of the agreement for TDD and YW are equal to 25% and 40%, respectively. This pattern seems to be

much stronger in the second task (see Figure 3.b), where the percentage of disagreements for TDD and YW are equal to 80% and 25%, respectively.

4.3 Results from Interviews

In the following, we report some of the main themes resulting from the thematic analysis by highlighting them in bold. To bring credibility to our themes, we present them together with some excerpts from the interview transcripts.

1. Appreciations for Overall TDD Experience. The students appreciated the overall experience consisting of theory and practice on TDD and related topics. For example, P4 stated:

| *“I liked the overall experience made of a mix of theory and practical stuff.”*

2. Appreciations for Deployment and Testing in Target Hardware. The students found it interesting and/or pleasant to see their ES running in the target hardware. For example, R2 said:

| *“Nice to see it [RS] run, and test by yourself, with hardware and sensors.”*

3. Intention to Use TDD for ESs Development. The students considered TDD useful to develop ESs. Therefore, they declared they would take into account TDD to develop their ESs in the future; although some of them declared to certainly use TDD, while others on the basis of some factors (*e.g.*, the ES to be developed). An excerpt from P1’s interview follows:

| *“In my opinion, TDD is better [than YW]...I would for sure use it in the future.”*

4. More Practice with TDD. The students found it easier to write their tests after the production code; therefore, they stated to need more time to be comfortable with TDD (where the opposite happens). In this respect, P6 said:

| *“In the beginning, I preferred YW because it is more natural to test the code after; however, after applying TDD in the last tasks I liked using it. However, I still need more practice with TDD.”*

5 Discussion

Below, We highlight the implications of our results for lecturers and researchers.

5.1 Implications for Lecturers

The results suggest that, when using TDD as compared to YW, the external quality of developed solutions increases, while there is not a substantial difference with respect to developers’ productivity (despite in some tasks we observed better productivity with TDD). These initial results should encourage lectures to teach TDD in the context of ESs development. If this happened, companies that deal with the development of ESs would be encouraged to use TDD, rather than other approaches, for the following reasons: (*i*) there would be initial evidence that TDD brings benefits to the development of ESs; and (*ii*) if developers, before

being hired, were already familiar with TDD applied to ESs development, the training costs incurred by software companies would be reduced or even null.

The participants also appreciated the overall experience with TDD applied to the ESs development, especially the deployment and testing in the target hardware. We can postulate that these positive experiences affected participants' intention to use TDD to develop their future ESs. Therefore, we suggest lectures that plan ESs courses with a well-thought mix of theory and practice on TDD applied to ESs development, giving space to deployment and testing in the target hardware.

5.2 Implications for Researchers

Our results (*i.e.*, improved external quality and no difference with respect to productivity) are similar to those of past studies on TDD applied to NoESs development (*e.g.*, [6, 20, 24]). Nevertheless, we advise further studies on TDD in the context of ESs development because: *(i)* the challenges developers tackle to develop ESs are different from those they tackle to develop NoESs [12]; and *(ii)* our results, despite being promising and consistent with those of past studies on TDD in the context of NoESs development, are preliminary. Also, our results suggest that future studies on TDD applied to ESs development should take into account a possible moderating role of the ESs to be implemented.

As for the post-questionnaire results, it emerged that both the development tasks with TDD and application of TDD were perceived as more difficult as compared to YW. To some extent, this outcome is coherent with past studies (*e.g.*, [3]), and a plausible explanation—also supported by the interview results—is that the participants found it more natural to write their tests after the production code rather than doing the opposite as it happens in TDD. Researchers could be interested in conducting longitudinal studies to investigate how developers' perception of TDD applied to ESs development changes over time, and whether changes in such a perception affect external quality and productivity.

Finally, we believe that our results could encourage other researchers to contribute to increasing the body of knowledge on TDD applied to ESs development, especially by conducting laboratory or field studies with software professionals—after all, it is easier to involve software companies and professionals when promising preliminary results are available. We also made our replication package, including the raw data, available online to allow researchers to replicate our study and ease the execution of secondary studies (*e.g.*, meta-analyses) on TDD applied to ESs development [21].

6 Threats to Validity

To determine the threats that might affect our results, we followed Wohlin *et al.*'s guidelines [26]. Although we tried to mitigate/avoid as many threats to validity as possible, some of them are unavoidable. This is because mitigating/avoiding a kind of threat (*e.g.*, internal validity) might intensify/introduce another kind of

threat (*e.g.*, external validity) [26]. Since we conducted the first study (comprising two experiments where both quantitative and qualitative data are gathered) on the application of TDD in the development of ESs, we preferred to reduce threats to internal validity (*i.e.*, making sure that the cause–effect relationships were correctly identified), rather than being in favor of external validity.

Construct validity. We measured each construct with a single dependent variable (*e.g.*, external quality with QLTY). Thus, in the case of measurement bias, the study results would be affected (threat of *mono-method bias*). Although we did not disclose the research goals of our study to the participants, they might have guessed them and changed their behavior based on their guess (threat of *hypotheses guessing*). To mitigate a threat of *evaluation apprehension*, we informed the participants that they would obtain the bonus points on the final exam mark regardless of their performance in both Exp1 and Exp2. Finally, there might be a threat of *restricted generalizability across constructs: i.e.*, TDD might have influenced some non-measured constructs.

Conclusion validity. We mitigated a threat of *random heterogeneity of participants* through two countermeasures: (*i*) we involved students taking the same course, allowing us to have participants with a similar background, skills, and experience; (*ii*) the participants underwent a training period to make them as more homogeneous as possible within each experimental group. A threat of *reliability of treatment implementation* might have occurred (*i.e.*, some participants might have followed TDD more strictly than others). To mitigate this threat, during the experiment, we reminded the participants to use the development approach we assigned them. Finally, we did not perform any statistical inference given the number of participants and exploratory nature of our study. That is, if we had performed statistical inference, we would have suffered from a threat of *low statistical power* (*i.e.*, even if there had been a true effect, we would have been probably unable to reject the corresponding null hypothesis).

Internal validity. The participation in the experiments was voluntary. Since volunteers are generally more motivated to carry out new tasks than the whole population, this might affect the obtained results (*selection* threat). Another potential threat is *resentful demoralization—i.e.*, participants receiving a less desirable treatment might not behave as they normally would. To mitigate a possible threat of *diffusion or imitation of treatments*, we monitored the participants during the execution of the experimental tasks in Exp1 and checked any delivered project for plagiarism—no hint of plagiarism emerged.

External validity. The participants were Master’s students, this might pose a threat of *interaction of selection and treatment—i.e.*, the results might not be generalized to other populations (*e.g.*, software professionals). However, the use of students has the advantage that they have more homogeneous backgrounds and skills, and allow obtaining initial empirical evidence [7, 13]. The implementation tasks might represent another threat: *interaction of setting and treatment*. However, we opted for ESs that can be considered representative of the domain of interest for our study.

7 Conclusion and Final Remarks

In this paper, we present the results of an exploratory investigation, comprising two experiments, to study the external quality of ESs when applying TDD, as well as developers' productivity. Specifically, the participants in the two experiments carried out ESs development tasks in Python by applying TDD or a non-TDD approach (YW). Our results suggest that the external quality of the implemented solutions increases when using TDD as compared to YW, while there is not a substantial difference with respect to developers' productivity. However, TDD is perceived as more difficult to apply than YW, and the development task is deemed more challenging with TDD. These results have implications for lecturers and researchers. For example, it is worth teaching TDD in the context of ESs development, with a well-thought mix of theory and practice, including the deployment and testing in the target hardware. Also, despite we gather preliminary evidence that TDD can be successfully applied to the development of ESs, we foster replications of our experiments, especially by involving software companies and professionals. Our exploratory investigation has the merit to justify these replications since it is easier to involve software companies and professionals when promising preliminary results are available.

References

1. Embedded systems market by component, by application: global opportunity analysis and industry forecast, 2021-2031. Tech. rep., Allied Market Research (2022), <https://www.alliedmarketresearch.com/embedded-systems-market-A08516>
2. Baldassarre, M.T., Caivano, D., Fucci, D., Juristo, N., Romano, S., Scanniello, G., Turhan, B.: Studying test-driven development and its retainment over a six-month time span. *J. Syst. Softw.* 176, 110937 (2021)
3. Baldassarre, M.T., Caivano, D., Fucci, D., Romano, S., Scanniello, G.: Affective reactions and test-driven development: results from three experiments and a survey. *J. Syst. Softw.* 185, 111154 (2022)
4. Basili, V., Caldiera, G., Rombach, H.D.: Goal Question Metric (GQM) Approach, pp. 528–532. John Wiley & Sons, Ltd (2002)
5. Beck, K.: Test-driven development: by example. Addison-Wesley (2003)
6. Bissi, W., Neto, A.G.S.S., Emer, M.C.F.P.: The effects of test driven development on internal quality, external quality and productivity: a systematic review. *Inf. Softw. Technol.* 74, 45–54 (2016)
7. Carver, J., Jaccheri, L., Morasca, S., Shull, F.: Issues in using students in empirical studies in software engineering education. In: *Proc. of International Symposium on Software Metrics*. pp. 239–249. IEEE (2003)
8. Cohen, J.: A power primer. *Psychol. Bull.* 112, 155–9 (1992)
9. Cordemans, P., Van Landschoot, S., Boydens, J., Steegmans, E.: Test-Driven Development as a Reliable Embedded Software Engineering Practice, pp. 91–130. Springer (2014)
10. Fucci, D., Scanniello, G., Romano, S., Shepperd, M., Sigweni, B., Uyaguari, F., Turhan, B., Juristo, N., Oivo, M.: An external replication on the effects of test-driven development using a multi-site blind analysis approach. In: *Proc. of International Symposium on Empirical Software Engineering and Measurement*. pp. 3:1–3:10. ACM (2016)

11. Ghafari, M., Gross, T., Fucci, D., Felderer, M.: Why research on test-driven development is inconclusive? In: Proc. of International Symposium on Empirical Software Engineering and Measurement. pp. 25:1–25:10. ACM (2020)
12. Grenning, J.: Test Driven Development for Embedded C. Pragmatic Bookshelf Series, Pragmatic Bookshelf (2011)
13. Höst, M., Regnell, B., Wohlin, C.: Using students as subjects—a comparative study of students and professionals in lead-time impact assessment. *Empir. Softw. Eng.* 5(3), 201–214 (2000)
14. Karac, I., Turhan, B.: What do we (really) know about test-driven development? *IEEE Software* 35(4), 81–85 (2018)
15. Karlesky, M.J., Bereza, W.I., Erickson, C.B.: Effective test driven development for embedded software. In: Proc. of International Conference on Electro/Information Technology. pp. 382–387. IEEE (2006)
16. King, N.: Using templates in the thematic analysis of text. In: Cassell, C., Symon, G. (eds.) *Essential Guide to Qualitative Methods in Organizational Research*, pp. 256–270. Sage (2004)
17. Kitchenham, B., Madeyski, L., Brereton, P.: Meta-analysis for families of experiments in software engineering: a systematic review and reproducibility and validity assessment. *Empir. Softw. Eng.* 25(1), 353–401 (2020)
18. Madeyski, L., Kitchenham, B.: Effect sizes and their variance for ab/ba crossover design studies. *Empir. Softw. Eng.* 23(4), 1982–2017 (2018)
19. Munir, H., Moayyed, M., Petersen, K.: Considering rigor and relevance when evaluating test driven development: a systematic review. *Inf. Softw. Technol.* 56(4), 375–394 (2014)
20. Rafique, Y., Mišić, V.B.: The effects of test-driven development on external quality and productivity: A meta-analysis. *IEEE Trans. Softw. Eng.* 39(6), 835–856 (2013)
21. Romano, S.: Replication package of “test-driven development and embedded systems: an exploratory investigation” (2023), <https://figshare.com/s/6322beb3a23373149862>
22. Romano, S., Zampetti, F., Baldassarre, M.T., Di Penta, M., Scanniello, G.: Do static analysis tools affect software quality when using test-driven development? In: Proc. of International Symposium on Empirical Software Engineering and Measurement. pp. 80–91. ACM (2022)
23. Tosun, A., Dieste, O., Fucci, D., Vegas, S., Turhan, B., Erdogmus, H., Santos, A., Oivo, M., Toro, K., Jarvinen, J., Juristo, N.: An industry experiment on the effects of test-driven development on external quality and productivity. *Empir. Softw. Eng.* 22(6), 2763–2805 (2017)
24. Turhan, B., Layman, L., Diep, M., Erdogmus, H., Shull, F.: How effective is test-driven development. In: *Making Software: What Really Works, and Why We Believe It*, pp. 207–217. O’Reilly Media (2010)
25. Vegas, S., Apa, C., Juristo, N.: Crossover designs in software engineering experiments: Benefits and perils. *IEEE Trans. Softw. Eng.* 42(2), 120–135 (2016)
26. Wohlin, C., Runeson, P., Hst, M., Ohlsson, M.C., Regnell, B., Wessln, A.: *Experimentation in Software Engineering*. Springer (2012)
27. Xiao, J.: Application of agile methods on embedded system development. In: Proc. of International Conference on Mechatronics, Materials, Biotechnology and Environment. pp. 667–670. Atlantis Press (2016)