

A Longitudinal Cohort Study on the Retainment of Test-Driven Development

Davide Fucci
University of Hamburg
Hamburg, Germany
fucci@informatik.uni-hamburg.de

Simone Romano
University of Basilicata
Potenza, Italy
simone.romano@unibas.it

Maria Teresa Baldassarre
University of Bari
Bari, Italy
mariateresa.baldassarre@uniba.it

Danilo Caivano
University of Bari
Bari, Italy
danilo.caivano@uniba.it

Giuseppe Scanniello
University of Basilicata
Potenza, Italy
giuseppe.scanniello@unibas.it

Burak Thurhan
Brunel University
London, UK
burak.turhan@brunel.ac.uk

Natalia Juristo
Universidad Politecnica de Madrid
Madrid, Spain
natalia@fi.upm.es

ABSTRACT

Background: Test-Driven Development (TDD) is an agile software development practice, which is claimed to boost both external quality of software products and developers' productivity.

Aims: We want to study: (i) the TDD effects on the external quality of software products as well as the developers' productivity; and (ii) the retainment of TDD over a period of five months.

Method: We conducted a (quantitative) longitudinal cohort study with 30 third-year undergraduate students in Computer Science at the University of Bari in Italy.

Results: The use of TDD has a statistically significant effect neither on the external quality of software products nor on the developers' productivity. However, we observed that participants using TDD produced significantly more tests than those applying a non-TDD development process, and that the retainment of TDD is particularly noticeable in the amount of tests written.

Conclusions: Our results should encourage software companies to adopt TDD because who practices TDD tends to write more tests—having more tests can come in handy when testing software systems or localizing faults—and it seems that novice developers retain TDD.

Proceedings of International Symposium on Empirical Software Engineering and Measurement (ESEM). ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 INTRODUCTION

Test-Driven Development (TDD) is a software development technique which leverages unit tests to incrementally deliver small pieces of functionality. Peculiar to TDD is the order in which tests and production code are written—the former are specified first and only in the case of a failure the developer is allowed to write production code to make them pass [2]. An important role in this process is played by refactoring which allows changing the code internal representation (e.g., an algorithm) while preserving its external behavior due to the safety net provided by the test suite [1].

TDD promised to support the delivery of high-quality products, both from a functional (e.g., fewer bugs) and technical perspective (e.g., “cleaner” code), while improving developers' productivity [2]. Consequently, industry has taken an interest in adopting this technique [4] and academia has dedicated large effort to gather empirical evidence to support or disprove its claimed effects. The results, gathered and combined in secondary studies [20, 21, 28], are conflicting and limited conclusions can be drawn. The primary studies, such as controlled experiments and case studies, are often cross-sectional [5] and only capture a “snapshot” of the phenomena at a given time. However, despite being recommended in the literature [10, 19, 26], only few investigations [15, 23] take a longitudinal perspective on the study of TDD. In one of these studies, Latorre [15] followed professional developers of different levels of seniority (but all with no experience in TDD) working on a project for a month while learning to apply the technique. The author studied how the conformance to the TDD process and the participants' productivity evolved during the investigation. The focus was to evaluate how the different subjects' learning curves affect their performance (e.g., in terms of code quality). A long-term case study at IBM by Sanchez *et al.* [23] aimed at understanding whether TDD improves over the process previously adopted in the company. The observation

CCS CONCEPTS

• **Software and its engineering** → **Software development techniques**;

KEYWORDS

Test-driven development, longitudinal cohort study

ACM Reference Format:

Davide Fucci, Simone Romano, Maria Teresa Baldassarre, Danilo Caivano, Giuseppe Scanniello, Burak Thurhan, and Natalia Juristo. 2018. A Longitudinal Cohort Study on the Retainment of Test-Driven Development. In

focused on a team and its sustained use of TDD for a period of five years. However, the investigation was carried out retrospectively—*i.e.*, using existing data gathered during such period but also before TDD was introduced. Conversely, we present a longitudinal cohort study [5] involving two separate observations of the same variables (*i.e.*, functional quality, productivity, and number of tests written), obtained from the same participants (*i.e.*, 30 novice developers), five months apart. Our cohort is composed of software developers of homogeneous experience who attended the same training regarding Agile software development principles, including TDD. Our goal is to understand how well TDD can be applied after the passage of time, giving an indication of its retention.

Thus, the main research question driving our study is:

To what extent can novice software developers *retain* TDD and its effects (if any) over a period of five months?

To establish a baseline, we compared the treatment of interest (*i.e.*, TDD) with the non-TDD development process (*e.g.*, iterative test-last, big-bang testing, or no testing at all) that subjects would normally follow. We refer to the latter as *Your Way* development (*i.e.*, YW).

This paper makes the following main contributions:

- an evidence-based discussion of TDD retainment and its implication for research and practice;
- a longitudinal design methodology that can be applied to other software development processes to distinguish between short-term and long-term phenomena;
- a laboratory package¹ to foster further replication of the presented longitudinal cohort study.

Paper organization. In Section 2, we present background information and related work. In Section 3, we describe the design of our longitudinal cohort study, while the results are presented and discussed in Section 4 and Section 5, respectively. Final remarks conclude the paper in Section 6.

2 BACKGROUND

In this section, we summarize the available evidence supporting (or refuting) the effects of TDD on external quality (or functional *e.g.*, number of defects) and developers' productivity. We also summarize research work on existing longitudinal studies in the context of Software Engineering (SE).

2.1 Types of longitudinal studies in SE

Longitudinal research in SE is not so common and appears to be mostly associated with the case study methodology. According to Yin [33], in a longitudinal case study data collection happens over an extended period with the goal of investigating “*how certain conditions change over time*” [33]. This is the case when the phenomenon under investigation is a process bounded to its context. Therefore, similarly to ethnography [25], longitudinal case studies require the researchers to be co-located with the case in which the phenomena takes place. For example, in the longitudinal case study reported by McLeod *et al.* [18], researchers spent several hundreds of hours

at the case company site. They attended meetings, observed and interviewed stakeholders within a period of two years to characterize software development as an emergent process. In a similar fashion, Salo and Abrahamson [22] followed how Software Process Improvement (SPI) techniques were introduced in the workflow of five Agile projects. Their investigation lasted for 18 months during which the researchers constantly recorded the output of retrospective meetings, interviews with the developers, as well as the metrics collected from the SPI tool in use at the company.

Longitudinal studies are useful when the observations cover an interesting event—*e.g.*, the introduction of a new practice within a company. Therefore, the researcher is interested in observing the impact of such change while it unfolds. This scenario is similar to interrupted time series in quasi-experimental designs [5]. For example, Li *et al.* [16] studies the changes brought by replacing a waterfall-like approach with Scrum in a small software company. The authors followed the development of a project for more than three years—17 months using waterfall and 20 using Scrum. This approach allowed for a *before-after* comparison of defects density and productivity. The longtime span was necessary to avoid a biased comparison between the established process and an immature one. A similar approach is reported in Vanhanen *et al.* [29] in which the impact of introducing pair programming was assessed over a period of two years with data collected through survey with the developers.

Other examples of longitudinal studies in SE cover a long period of time in retrospect—*e.g.*, by analyzing archival data. Harter *et al.* [11] analyzed the type of defects identified over time by the progressive introduction of SPI techniques in a firm and its subsequent CMMI improvements over a period of 20 years.

In the health science and medicine, longitudinal studies are sometimes realized in the form of cohort studies. A cohort is a sample of subjects (*e.g.*, who undergo a treatment) sharing a specific characteristic of interest (*e.g.*, age). The cohort is tested in several occasions over time to, for example, check for a drug side-effect before releasing it to the market [5].

2.2 Effects of TDD

The effects of TDD on several outcomes, including the ones of interest for this study—*i.e.*, functional quality and productivity—is the topic of several empirical studies, summarized in Systematic Reviews (SR) and Meta-Analysis (MA). Turhan *et al.* [28] SR includes 32 primary studies in which TDD was investigated in different settings. Although the results show a positive effect on quality, the ones regarding productivity are inconclusive. Rafique and Mistic [21] conducted an MA covering 25 primary studies published between 2000 and 2011. When considering participants from academia, TDD seems to improve quality to the loss of productivity. Finally, Munir *et al.* [20] SR classifies the primary studies according to relevance and rigor dimensions. The results show, for both student and professional developers, that TDD increases quality but not productivity. The authors recommend increasing relevance by carrying out long-term studies. One example of such investigation is presented in Marchenko *et al.* [17] which reports a three-year-long case study about the use of TDD at Nokia-Siemens Network. The authors observed and interviewed eight participants (one Scrum master,

¹www2.unibas.it/sromano/downloads/LabPackageUniba.zip

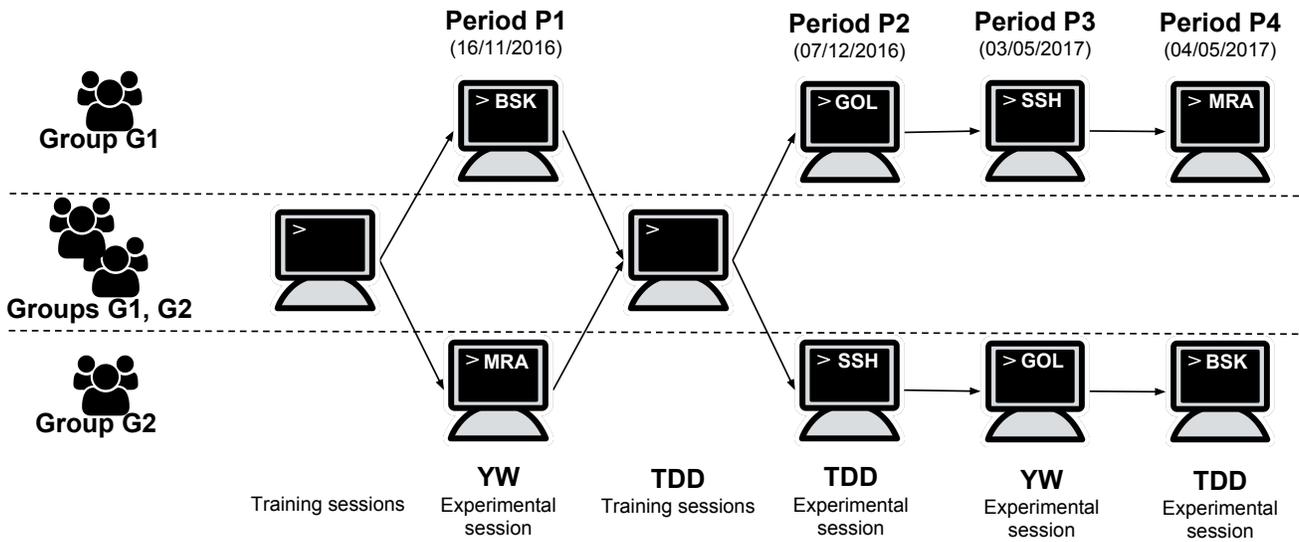


Figure 1: Summary of the study.

one Product owner, and six developers) and extracted themes from the data. The participants perceived TDD as an important driver towards the improvement of their code quality, both from a structural and functional perspective. Moreover, the team confidence with the code base improved, which is associated with improved productivity [2]. The examined team reported that TDD was not suitable for bug fixing, especially for bugs that are difficult to reproduce (e.g., needing a specific environment setup) or for quick “hacks” due to the testing overhead. The authors also report some concerns regarding the lack of a solid architecture when applying TDD.

Latorre [15] studied the capability of 30 professional software developers (junior, intermediate, and experts) to develop a real-world system using TDD. The study targeted the *learnability* of TDD, as the participants did not know the technique beforehand. The longitudinal one-month study started after giving the developers, proficient in Java and unit testing, a tutorial about TDD. After only a short practice session, the participants were able to correctly apply TDD (e.g., following the prescribed steps). Although they correctly followed the process between 80% and 90% of the time, their ability to initially apply TDD depended on experience—while seniors needed only few iterations, intermediates and juniors needed more time. Experience had an impact on productivity—only the experts were able to be as productive as they were when applying a traditional development methodology (measured during the initial development of the system). Refactoring and design decision hindered the productivity of intermediates and junior participants. Finally, regarding functional quality, all the participants in the study delivered a correct version of the system regardless of their seniority level.

3 LONGITUDINAL COHORT STUDY

In this section, we describe the planning of our longitudinal cohort study. We summarize the most important steps of this study in Figure 1. In particular, the participants in the study (groups G1 and G2) first took part in training sessions (and accomplished homework) where they practiced unit testing, iterative test-last development, and big-bang testing. Then, the participants in the groups G1 and G2 were asked to perform two implementation tasks on two different experimental objects (i.e., BSK² and MRA,³) in the same period P1—a period is the time during which a treatment is applied [30]. In P1, the participants could apply only YW because they were not aware of TDD yet. Between the periods P1 and P2, all the participants (G1 and G2) practiced TDD during training sessions (and homework). In the second period (i.e., P2), we asked the participants in G1 and G2 to perform other two tasks, GOL⁴ and SSH⁵ respectively. The applied treatment was TDD for both groups. After five months, we asked the same participants in group G1 to implement a new task—i.e., SSH—during P3 using the YW approach. In the same period (i.e., P3), the participants of G2 implemented the GOL task using YW. While, during P4, the participants of G1 and G2 were asked to apply TDD on MRA and BSK, respectively. We considered P3 and P4 to study the effects of TDD and its retainment. We introduced the last periods, P3 and P4, several months apart from the first two to assess whether the initial knowledge of TDD is retained.

The planning of our cohort study is reported according to the template by Jedlitschka *et al.* [13]. When planning and conducting our study, we followed the guidelines by Juristo and Moreno [14] and Wohlin *et al.* [32].

²BSK (Bowling ScoreKeeper) is an API for calculating the score of a bowling game.

³MRA (Mars Rover API) is an API for moving a rover on a planet.

⁴GOL (Game Of Life) is an API for Conway’s game of life.

⁵SSH (SpreadSheet) is an API for a spreadsheet.

3.1 Research Questions

We aimed at investigating the following Research Questions (RQs):

- RQ1.** To what extent do novice software developers retain TDD and how does this affect their performance?
- RQ2.** Are there differences between TDD and YW in the external quality of the implemented solutions, developers' productivity, and number of tests written?

We defined RQ1 to study whether TDD retainment affects the application of YW and whether there are deteriorations (or improvements) in the application of TDD over five months. The considered constructs are, external quality of the implemented solutions, developers' productivity, and the number of tests developers wrote.

Finally, RQ2 aimed at understanding whether the claim that TDD increases both external quality of the software products and developers' productivity is well-founded as well as whether TDD leads developers to write more tests.

3.2 Experimental Units

The participants were third-year undergraduate students in Computer Science. They were sampled by convenience among the students taking the *Integration and Testing* course at the University of Bari in Italy. The course covered the following topics, software quality, unit testing, integration testing, SOLID principles, refactoring, iterative test-last development, big-bang testing, and TDD. The course included frontal lectures, laboratory sessions, and homework. During the laboratory sessions, the students improved their knowledge about how to develop unit tests in Java by using the Eclipse IDE and JUnit, and the refactoring functionality available in Eclipse. During laboratory sessions and by developing homework, the students practiced unit testing, iterative test-last development, big-bang testing, and TDD.

Participation in the cohort study was voluntary. We informed the students that any gathered data would be treated anonymously and used for research purposes only. We also informed them that their performance in the study would not affect their final mark for the Integration and Testing course. To encourage the participation, we rewarded who accepted to take part in the study with a bonus in their final mark. Among the students taking the Integration and Testing course, 30 accepted to participate.

The participants had passed the exams for the courses of Procedural Programming, Object Oriented Programming, Software Engineering, and Databases. During these courses, all participants had acquired programming experience in C and Java. Between the first two periods and the last two, the participants followed the same university curricula courses in which TDD was not used. However, we did not control whether, within such period, the participants practices TDD outside the academic scope (*e.g.*, in personal projects).

3.3 Experimental Materials

The experimental objects were four code katas (*i.e.*, programming exercises used to practice a technique or a programming language).

- **BSK.** It is an API for calculating the score of a bowling game. This API allows adding a frame to a game, as well as bonus

throws; computing the score of a frame; identifying when a frame is spare or strike; and computing the score of a game.

- **MRA.** It is an API for moving a rover on a planet, which is represented as a grid. The cells of this grid can contain obstacles that the rover cannot pass through. MRA allows the initialization of a planet (*i.e.*, defining the grid with the obstacles) and moving the rover on the planet by parsing a list of commands (*i.e.*, turning left/right and moving forward/backward).
- **SSH.** It is an API for a spreadsheet. SSH allows evaluating the content of a cell and thus returning the result of this evaluation. Cells can contain integers, strings, references, and formulas (*e.g.*, concatenation of strings or integer addition).
- **GOL.** It is an API for Conway's game of life. This game takes place on a square grid of cells. Each cell can assume two states: alive or dead. At each step, the current state of the grid is used to determine the next state. GOL allows initializing the grid and determining the next state of each cell (it depends on the current state of the cell and of its neighbors) and then of the grid.

The implementation of the aforementioned APIs did not require any graphical user interfaces.

Each experimental object was composed of several user stories⁶ to be implemented, as well as a template project (of Eclipse) that contained a stub of the expected API signature and an example JUnit class test. To verify that the user stories were correctly implemented, each experimental object was accompanied by acceptance test suites—an acceptance test suite for each user story. It is worth mentioning that the acceptance test suites were not provided to the participants. That is, these suites were only exploited to quantify the quality of the solutions implemented by the participants and their productivity (see Section 3.5).

The use of code katas in empirical studies on TDD is quite common (*e.g.*, [6–8]). For BSK and MRA, we exploited the materials used in the experiment by Fucci *et al.* [8]. As for SSH and GOL, we created the experimental materials (*e.g.*, template projects).

3.4 Tasks

Each task was coupled to an experimental object (*i.e.*, four tasks, one for each experimental object). A task consisted of implementing a solution for an experimental object (*e.g.*, BSK). To this end, we provided the participants with: (i) the user stories to be implemented for the considered experimental object; and (ii) the template project. Thus, the participants had to use the template project when implementing the user stories for that experimental object.

3.5 Hypotheses and Variables

The participants were asked to carry out each task by using either TDD or the approach they preferred (*i.e.*, YW)—of course, in this latter case, they could not use TDD. Therefore, one of the independent variable (also known as factor) is **Technique**. It is a nominal variable assuming two values, *TDD* and *YW*. Since our study is longitudinal—*i.e.*, we collected data over time—we took into account another independent variable, which represents the

⁶A user story is a description of a feature to be implemented from the perspective of the end user.

Table 1: Design Summary.

		Period			
		P1 (16/11/2016)	P2 (07/12/2016)	P3 (03/05/2017)	P4 (04/05/2017)
Group	G1	YW, BSK	TDD, GOL	YW, SSH	TDD, MRA
	G2	YW, MRA	TDD, SSH	YW, GOL	TDD, BSK

period during which each treatment (*i.e.*, TDD or YW) was applied. We named this variable **Period**. It is a nominal variable and assumes the following values, *P1*, *P2*, *P3*, and *P4*. We also considered the independent variable **Group** representing the two groups of participants. It is a nominal variable assuming two values: *G1* and *G2*.

The dependent variables considered in our study are, **QLTY**, **PROD**, and **TEST**. We choose these dependent variables because they have been previously used in other empirical studies on TDD (*e.g.*, [6–8, 27]). The variable **QLTY** quantifies the external quality of the solution a participant implemented. This variable is defined as follows (*e.g.*, [8]):

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

where $\#TUS$ is the number of user stories a participant tackled, while $QLTY_i$ is the external quality of *i*-th user story. To understand if a user story was tackled or not, we checked the asserts in the corresponding acceptance test suite. Namely, if at least one assert in the test suite (for that story) passed, then the story was tackled. $\#TUS$ is formally defined as:

$$\#TUS = \sum_{i=1}^n \begin{cases} 1 & \#ASSERT_i(PASS) > 0 \\ 0 & otherwise \end{cases} \quad (2)$$

On the other hand, the quality of the *i*-th user story (*i.e.*, $QLTY_i$) is equal to the number of asserts passed for the *i*-th story with respect to the total number of asserts for the same story. More formally:

$$QLTY_i = \frac{\#ASSERT_i(PASS)}{\#ASSERT_i(ALL)} \quad (3)$$

Given the Formulas 1,2, and 3, **QLTY** assumes values between 0 and 100, where a value close to 0 means that the quality of the solution is low, while a value close to 1 indicates high quality of the solution.

The variable **PROD** estimates the productivity of a participant. It is computed as follows (*i.e.*, [27]):

$$PROD = \frac{\#ASSERT(PASS)}{\#ASSERT(ALL)} * 100 \quad (4)$$

where $\#ASSERT(PASS)$ is the number of asserts passed, by considering all the acceptance test suites, with respect to the total number of the asserts in the acceptance test suites. **PROD** assumes values between 0 and 100. A value close to 0 indicates low productivity, while a value close to 1 means high productivity.

The variable **TEST** quantifies the number of unit tests a participant wrote. It is defined as the number of asserts in the test suite written by a participant when tackling a task (*e.g.*, [8]). **TEST** ranges from 0 to ∞ .

We formulated the following parameterized null hypotheses,

HN1_X. There is no significant effect of Period with respect to X (*i.e.*, **QLTY**, **PROD**, or **TEST**).

HN2_X. There is no significant effect of Technique with respect to X (*i.e.*, **QLTY**, **PROD**, or **TEST**).

The alternative hypotheses are two-tailed—*i.e.*, we did not consider the direction of the effect for either independent variable. **HN1_X** was defined to investigate RQ1, while we defined **HN2_X** to investigate RQ2.

3.6 Study Design

The design of the cohort study is depicted in Table 1. The participants were randomly split into two groups—G1 and G2—having 15 participants each. Whatever the group was, the participants were assigned to each treatment (*i.e.*, TDD or YW) twice. In particular, both groups were assigned to: YW in the first period (*i.e.*, P1), TDD in the second period (*i.e.*, P2), YW in third period (*i.e.*, P3), and TDD in the last period (*i.e.*, P4). Therefore, the design of our study can be classified as a *repeated measures, within-subjects* design. In each period, the participants in G1 and G2 dealt with different experimental objects. For instance, in P1, the participants in G1 dealt with BSK, while those in G2 with MRA. At the end of the study, every participant had tackled each experimental object only once.

3.7 Procedure

Before our study took place, we collected some demographic information on the participants. To this end, the participants filled out an on-line pre-questionnaire (created by means of Google Forms).

The Integration and Testing course—in which the cohort study was conducted—started in October, 2016. The first application of the YW treatment (*i.e.*, P1) took place on November 16th, 2016 (see Table 1). Between the beginning of the course and P1, the participants had never dealt with TDD, while they knew unit testing, iterative test-last development, and big-bang testing. On these techniques, the participants had taken part in two training sessions and carried out some homework. TDD was introduced to the participants between P1 and P2. They also had taken part in three training sessions on TDD and completed some homework by using this development practice. Given the previous considerations, we can exclude that the knowledge of TDD affected in anyway the YW treatment in P1. The YW treatment was applied again on May 3rd, 2017 (in P3), while TDD was applied the day after in P4. From P2 to P3 five months passed. Since the participants knew TDD in P3, we cannot exclude that the knowledge of TDD would have affected the treatment YW in P3 in some way. That is, if the TDD retainment had affected the application of YW or not. On the other hand, we assessed the retainment of TDD by asking the participants to use TDD (once again) in P4.

Table 2: Some descriptive statistics for each dependent variable grouped by Period and by Technique.

Variable	Statistic	Period (Technique)				Technique	
		P1 (YW)	P2 (TDD)	P3 (YW)	P4 (TDD)	YW	TDD
QLTY	mean	59.3989	63.1002	63.0505	58.535	61.2247	60.8176
	median	76.7614	69.7251	71.2867	74.7614	72.9702	71.9962
	SD	37.8509	31.989	30.7322	34.5895	34.232	33.1112
PROD	mean	34.1145	32.4793	30.991	37.9692	32.5527	35.2243
	median	27.5281	29.0698	27.907	42.8571	27.907	34.8837
	SD	32.182	29.039	28.9798	29.194	30.403	29.0012
TEST	mean	4.9333	7.8333	7.9333	10.1	6.4333	8.9667
	median	4	6.5	5	8.5	5	7
	SD	4.0508	5.5216	7.5198	7.2462	6.1764	6.4885

The execution of the study tasks took place under controlled conditions in a laboratory at the University of Bari. In each period, the participants in G1 and G2 were randomly assigned to the laboratory PCs. We alternated participants in G1 in G2 to avoid that participants of the same group assigned to the same experimental object were close to each other. This setup limited interactions among the participants. In addition, we monitored them during the execution of tasks.

All the PCs in the laboratory were equipped with the same hardware and software. Furthermore, they were set up with all the experiment materials necessary for carrying out the task, *i.e.*, the template project (of Eclipse) corresponding to the assigned experimental object. Each subject provided a solution for the assigned task by using the template project. The participants implemented the tasks in Java and used JUnit as testing framework. At the end of each task, the participants uploaded their solution through GitHub and then filled out a post-questionnaire. This questionnaire collected feedback on how the participant perceived the execution of each task.

3.8 Analysis Procedure

The gathered experimental data were analyzed according to the following procedure:

- (1) **Descriptive Statistics.** We computed descriptive statistics, *i.e.*, mean, median and standard deviation (SD), to summarize the distributions of the dependent variable values. We also used boxplots to graphically summarize these distributions.
- (2) **Inferential Statistics.** We used Linear Mixed Model (LMM) analysis methods to test the defined null hypotheses. LMM is a popular method for analyzing data from longitudinal studies [31]. For each dependent variable, we built an LMM that included the following terms: Period, Group, and Period:Group (*i.e.*, the interaction between Period and Group) as fixed effects, while the participant represents the random effect (this is customary in SE experiments [30]). It is worth noting that the periods P1 and P3 correspond to the application of the YW treatment, while TDD was applied in the periods P2 and P4. This means that, if the LMM does not indicate a statistically significant effect of Period, then there is no statistically significant effect of Technique. To build

LMMs, we considered Group because, based on the study design, it also represents the sequence (*i.e.*, the order in which the treatments are applied in combination with the experimental objects). In repeated measures designs the effect of sequence on the dependent variables must be analyzed [30]. LMM analysis methods have two assumptions: (i) the residuals of LMM have to be normally distributed and (ii) their mean has to be equal to zero [30]. If these two assumptions are not verified, transforming the data of the dependent variable is an option (*e.g.*, by using log or square-root transformation) [30]. To check if the residuals were normally distributed, we used the Shapiro-Wilk test (Shapiro test, from here onwards) [24]. As it is customary with tests of statistical significance, we accepted a probability of 5% of committing Type-I error (*i.e.*, $\alpha = 0.05$).

4 RESULTS

In this section, we first report the results from the descriptive statistics followed by those pertaining the inferential statistics.

4.1 Descriptive Statistics

In Table 2, we report the values of mean, median, and SD for each depended variable. These values are grouped by both Period and Technique. We also show the boxplots for the dependent variables in Figure 2.

QLTY. As shown in Table 2 and Figure 2.a, there are no noticeable differences in the QLTY values between the periods. In particular, by comparing the boxplots for P1 and P3—*i.e.*, same treatment (YW) but different experimental objects—we can observe that these boxplots overlap and the median level in P1 is higher than that in P3 (see also the median values in Table 2).

Similarly, we can observe that the boxplots for P2 and P4—*i.e.*, same TDD treatment but different experimental objects—overlap and the median level is higher in P4. Therefore, such slight differences in the QLTY values seem to be due to the experimental objects rather than the retainment of TDD. Namely, when the experimental objects are BSK and MRA (*i.e.*, in P1 and P4), the median levels are higher.

When comparing TDD and YW, the results in Table 2 do not suggest differences in QLTY values (*e.g.*, on average, QLTY is equal

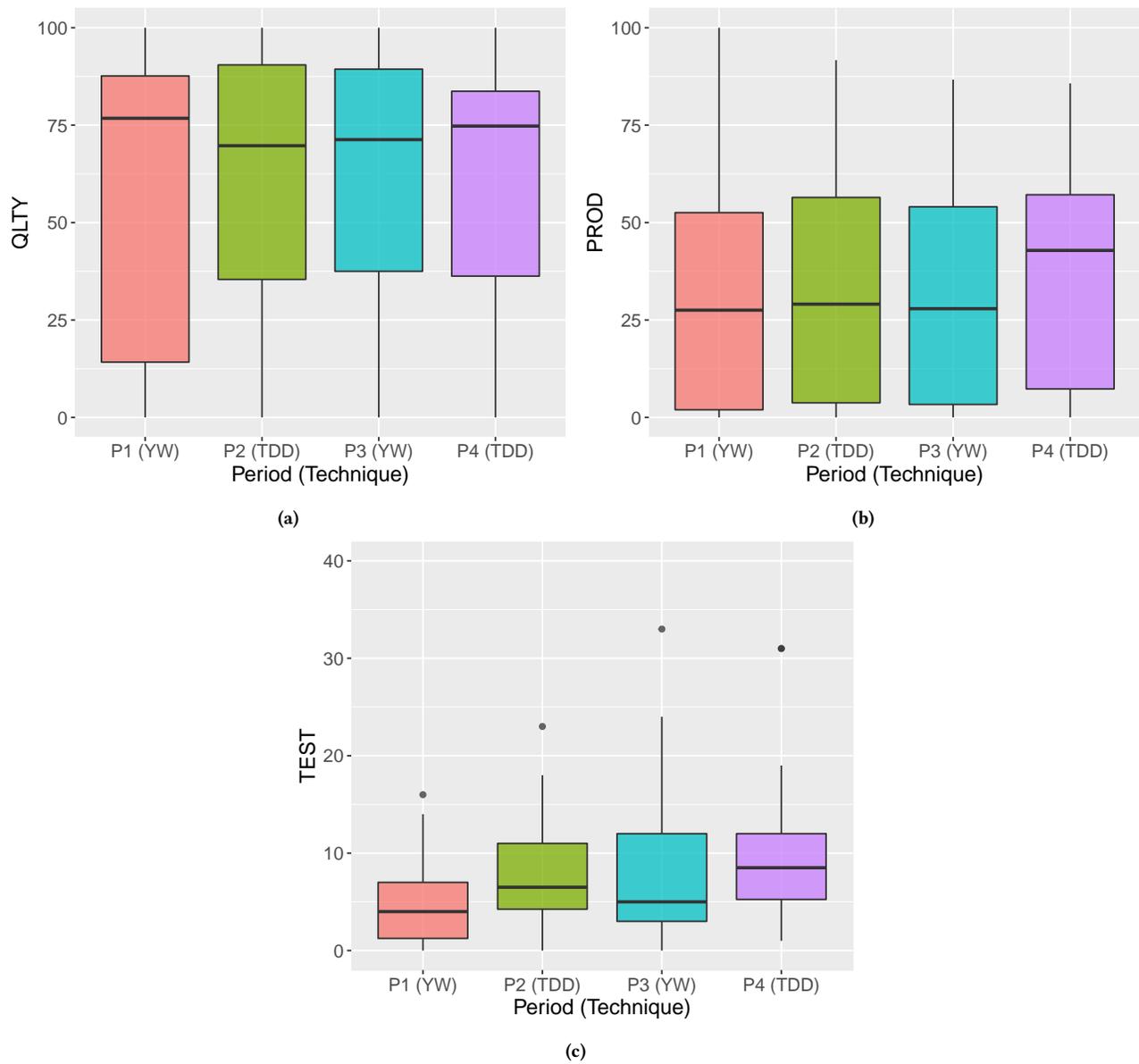


Figure 2: Boxplots for QLTY (a), PROD (b), and TEST (c) for each period.

to 60.8176 for TDD, while it is equal to 61.2247 for YW). This outcome is confirmed when we compare P4 (TDD) with P1 (YW) and P2 (TDD) with P3 (YW)—same experimental objects. For instance, the participants in P4 and P1 achieved, on average, similar values for QLTY (58.535 vs. 59.3989) although, when dealing with the same experimental objects, they applied either TDD or YW. The comparison between P2 and P3 lead to a similar observation.

PROD. The boxplots in Figure 2.b do not indicate noticeable differences in the PROD values among the periods. Indeed, when passing from P2 to P4, we can observe that the boxplots overlap, but the median level for P4 is higher than for P2. In particular, the medians of the PROD values are equal to 42.8571 and 29.0698 for P4 and P2,

respectively. This improvement in the PROD values might be due to the TDD retainment. As for the comparison between P1 and P3, the boxplots are very similar to each other. Therefore, it seems that the knowledge the participants had of TDD (*i.e.*, its retainment) in P3, with respect to P1, did not affect QLTY.

The results in Table 2 seem to suggest that there is a slight difference in the PROD values between TDD and YW in favor of TDD (*e.g.*, PROD for TDD is equal to 35.2243 on average, while it is equal to 32.5527 for YW). By comparing pairs of periods in which the same experimental objects are used, we can observe that the PROD values in P4 (TDD) are better than those in P1 (YW). Namely, it seems that the participants who applied TDD on BSK and MRA

Table 3: Results (i.e., p-values) from the LMM analysis methods for the dependent variables QLTY, PROD, and TEST.

Variable	Period	Group	Period:Group
QLTY	0.8837	0.6108	<0.0001*
PROD	0.7973	0.8225	<0.0001*
TEST	0.0002*	0.0617	0.4632

* Statistically significant effect.

achieved PROD values better than the participants who applied YW. This trend is not observed when comparing P2 (TDD) and P3 (YW). For instance, the boxplot for P2 is very similar to that for P3, so suggesting that there is no difference with respect to the dependent variable PROD.

TEST. By looking at the boxplots in Figure 2.c, we can observe differences in the TEST values among the periods. In particular, if we compare the YW treatments in P1 and P3, it appears that the boxplot for P3 is higher than that for P1. The TEST values for P3 are also better on average (7.9333 for P3, 4.9333 for P1). This difference might suggest a positive effect of the TDD retainment when participants had to apply YW in P3. On the other hand, the boxplots for TDD in P2 and P4 suggest a less pronounced difference in TEST values. The boxplots for P2 and P4 overlap, although, as for P4, the boxplot is shorter and the median value is higher (8.5 in P4 vs. 6.5 in P2).

The comparison between TDD and YW seems to suggest that the TEST values for TDD are higher than those for YW. For instance, the mean values are equal to 8.9667 and 6.4333 for TDD and YW, respectively. By considering only P4 and P1, we can observe a clear improvement in the TEST values in P4 (see the boxplots). Namely, the participants who applied TDD in P4 seem to achieve values for TEST higher than those who applied YW in P1 on the same experimental objects (e.g., the mean values are 10.1 in P4, while 4.9333 in P1). Interestingly, the comparison between P2 (TDD) and P3 (YW) does not confirm the trend previously observed. Namely, it seems that the distributions of the TEST values for P2 (TDD) and P3 (YW) are quite similar (see both boxplots and descriptive statistics), despite the application of either TDD (in P2) or YW (in P3) on the same experimental objects. This outcome can indicate that the TDD retainment influenced the participants who applied YW in P3.

4.2 Inferential Statistics

The results (i.e., p-values) from the LMM analysis methods are reported in Table 3. When a p-values is less than α , we highlighted it with the \star symbol.

QLTY. The assumptions of the LMM analysis method for QLTY were both verified, so we did not perform any data transformation. As shown in Table 3, the LMM analysis method does not allow us to reject $HN1_{QLTY}$, the p-value for Period is 0.8837, namely the effect of Period is not statistically significant. This means that either there is no deterioration nor improvement in the observed time period (i.e., no effect of time period) with respect to QLTY, or that the test does not have enough statistical power to show differences, would they exists. The built LMM also suggests that the

effect of Group is not statistically significant, while the interaction between Group and Period is. This interaction is due to the effect of the experimental objects (e.g., whatever the treatment is, the distributions for BSK are higher than those for GOL).

Since LMM analysis method for QLTY does not indicate an effect of Period, the effect of Treatment is not statistically significant either. Therefore, we cannot reject $HN2_{QLTY}$.

PROD. The LMM analysis method for PROD needed data transformation since the method assumptions were not satisfied. In particular, we applied a square-root transformation to meet these assumptions. The results in Table 3 show that the effect of Period is not statistically significant. Therefore, we can neither reject $HN1_{PROD}$ nor $HN2_{PROD}$, indicating that the participants may retain TDD with respect to PROD. Moreover, applying either TDD or YW seems to not affect the PROD values. The LMM also includes a significant effect, namely Group:Period. Again, this significant interaction is due to the effect of the experimental objects.

TEST. To apply the LMM analysis method for TEST, we had to transform the data of the dependent variable. In particular, we performed a log transformation so that the assumptions were verified. The LMM analysis suggests that the effect of Period is statistically significant (the p-value is equal to 0.0002). Therefore, we can reject $HN1_{TEST}$. There is evidence that *a significant effect of Period on the number of tests the participants wrote* exists. According to the boxplots in Figure 2.c, the significant difference in Period is not due to a deterioration of TEST values for TDD over time—the worst distribution can be observed in P1—therefore, we can conclude that the ability of writing unit tests is retained by developers using TDD.

Since we found a significant effect of Period and in accordance with the results from the descriptive statistics (i.e., there is a clear difference in favor of TDD in P4 with respect to YW in P1 on the same experimental objects), we reject $HN2_{TEST}$. Therefore, we can conclude that *there is a significant effect of Technique on the number of tests the participants wrote*.

5 DISCUSSION

In this section, we discuss the results obtained according to the RQs and present possible practical implications from our research. Finally, we delineate threats that could have affected the validity of our study.

5.1 Answers to Research Questions

The data analysis gives some indication that developers retain TDD. In particular, we observed neither deteriorations in the external quality of the solutions developed by the participants nor in their productivity. Moreover, it seems that, with time, there is an improvement in the number of tests written when using TDD.

Our results do not suggest differences between TDD and YW with respect to the quality of the implemented solutions, as well as the developer’s productivity. However, who practices TDD tends to write more tests.

5.2 Implications

The participants retention of TDD is particularly noticeable in the amount of unit tests written. This is in line with the findings of a cross-sectional study by Erdogmus *et al.* [6] which pointed out that

the number of tests correlates with the ability of novice developers to follow TDD. Our study extends that notion to a longitudinal perspective—TDD helps retain TEST over a period of five months. The motivations for such effect are to be considered for further studies. However, we believe that TDD raises the participants awareness about the importance of writing several (fine-grained) unit tests. Nevertheless, this does not translate in improved QLT_Y, nor PROD. The latter result contrasts the ones of Latorre [15] which saw, over a period of a month of *constant* observation, a steady and significant improvement in performance measure similar to our QLT_Y.⁷ We conjecture that this can be the case due to the better experience of Latorre’s study participants (*i.e.*, professional developers), furthering the thesis that TDD alone is not a silver bullet but pre-existing skills play a crucial role [4, 10]

Based on our findings, software companies that value unit testing (*e.g.*, for creating a regression for continuous integration) should encourage the use of TDD as developers are likely to produce more tests when using such technique. We showed that a small initial investment in training results in retainment of this particular feature on the long term. Likewise, computer science educators should include TDD early in their curricula to install a long-term unit-testing mentality in the students. Finally, “Experience with TDD” is a characteristic that researchers should foster when building a sample for studies requiring (novice) participants familiar with unit testing. Likewise, when designing experiments where unit testing is desirable, researchers can avoid or limit time spent on training as, at least in our time-frame, such skill is retained by (novice) participants who already have minimal TDD experience.

Our results did not show any improvement of TDD over YW, contributing to the null results in TDD research [8, 9, 21]. However, differently from previous attempts, we showed that no effects are observable also when the same subjects are tested again several months later, under similar conditions. Time did not drastically decrement the novices performance when TDD was applied, hinting at the fact that they soon regain familiarity with technique similarly to what the study of Latorre [15] reports for the junior developers in the sample. Although carrying out cohort longitudinal studies—in particular, with several observations over a long time span—is difficult in SE (*e.g.*, controlling for maturation or learning effects), we put forward the idea that we might not be looking long enough (rather than hard enough) for the claimed effects of TDD to become apparent. As a starting point towards this direction, we recommend longitudinal studies in academia, which allow to follow the “career” of students over several years and can achieve a good amount of control (*e.g.*, based on grades, attendance)

5.3 Threats to Validity

We discuss the threats that could have affected the validity of the obtained results according to the guidelines by Wohlin *et al.* [32]. Accordingly, we ranked our threats from the most sensible for the goal of this study to the least one. In particular, being this the first test for a theory of TDD retainment, we prefer to limit threats to internal validity (*i.e.*, make sure that the cause-effect relationships are correctly identified), rather than being in favor of generalization.

5.3.1 Threats to Internal Validity. This kind of threat concerns internal factors of the study that could have affected the results. The effect of letting volunteers take part in the study may influence the results because volunteers are generally more motivated [32] (*i.e.*, *selection threat*). To prevent participants exchanging information during the tasks (*i.e.*, *threat of diffusion or treatments imitations*), at least two researchers monitored them. We also prevented the diffusion of experimental materials by gathering it at the end of each task. A threat of *resentful demoralization* might exist. For instance, a participant, who was given a less desirable treatment or task, might not perform as good as they generally would. This threat to validity might have equally affected both TDD and YW. Finally, control over subject *maturation* was checked by making sure that the students attended the same courses between the first observation and the last one.

5.3.2 Threats to Construct Validity. These threats concern the relationship between theory and observation. The investigated constructs were quantified by means of one dependent variable each. This might affect the results (*i.e.*, *threat of mono-method bias*). However, we used well-known and widely adopted dependent variables in TDD experiments (*e.g.*, [8]). Although the participants were not informed about the goals of the study, they might guess them and change their behavior accordingly (*i.e.*, *threat of hypotheses guessing*). To mitigate an *evaluation apprehension threat*, we told the participants that they would not be evaluated on the basis of their performances in the study.

5.3.3 Threats to Conclusion Validity. This kind of threat concerns the relationship between the dependent and independent variables. To mitigate a *threat of random heterogeneity of participants*, our sample included students with similar backgrounds—*i.e.*, students taking the same courses in the same university with similar development experience. In empirical studies like this one, a *threat of reliability of treatment implementation* might also exist. For example, a participant might follow TDD more strictly than another one. We did not control for such effect in this study, however, we explicitly reminded the participants to follow the treatment they were assigned to. The treatments might have impacted other constructs which were not observed (*e.g.*, number of refactoring, code complexity). Nevertheless, we focused on the most salient dependent variables according to the literature. Finally, our sample was limited due to difficulty of recruiting participants available for the period of the entire study.

5.3.4 Threats to External Validity. These threats concern the generalizability of the results. The participants in our longitudinal study were students, thus generalizing the obtained results to the population of professional developers poses a *threat of interaction of selection and treatment*. However, the use of students as participants also implies a number of advantages [3], such as participants with homogeneous background, the possibility to obtain preliminary evidence. In addition, the tasks to be performed did not require a high level of industrial experience, we believe that the use of students as participants could be considered appropriate, as suggested in the literature [3, 12]. The experimental objects might also affect the external validity of the results (*i.e.*, *threat of interaction of setting and treatment*) as they are not representative of real-world settings.

⁷In Latorre [15], all the subjects completed the task—*e.g.*, achieved PROD of 100%.

On the other hand, simpler tasks, which can be completed in a single exercise session (approximately three hours), allow a better control over the participants. The latter was our preferred trade-off due to the theory-testing nature of this study.

6 CONCLUSION

In this paper, we present a quantitative longitudinal cohort study to investigate: (i) the TDD effects on the external (*i.e.*, functional) quality of software products as well as the developers' productivity; and (ii) the retainment of TDD over a period of five months. The results indicate that the use of TDD has a statistically significant effect neither on the external quality of software products nor on the developers' productivity. However, we observed that participants using TDD produced significantly more tests than those applying a non-TDD development process, and that this capacity was retained over time. On the basis of these findings, we speculate that software companies can be encouraged to adopt TDD because (i) it compels developers to write more unit tests, which can be leveraged for localizing faults through regression and (ii) it requires minimal initial effort to retain its effect.

There are several future directions for the research presented in this paper. First, we will replicate this study in academic context, with the same cohort but on a longer time-span with several observations. Second, we will replicate the same design presented here with professional developers—*e.g.*, in the form of workshops about Agile software development. Third, we will devise qualitative longitudinal studies to triangulate statistical results.

REFERENCES

- [1] Dave Astels. 2003. *Test driven development: A practical guide*. Prentice Hall Professional Technical Reference.
- [2] Kent Beck. 2003. *Test-driven development: by example*. Addison-Wesley Professional.
- [3] Jeffrey Carver, Letizia Jaccheri, Sandro Morasca, and Forrest Shull. 2003. Issues in Using Students in Empirical Studies in Software Engineering Education. In *Proceedings of International Symposium on Software Metrics (METRICS '03)*. IEEE, 239–249.
- [4] Adnan Causevic, Daniel Sundmark, and Sasikumar Punnekkat. 2011. Factors limiting industrial adoption of test driven development: A systematic review. In *Software Testing, Verification and Validation (ICST), 2011 IEEE Fourth International Conference on*. IEEE, 337–346.
- [5] Thomas D Cook, Donald Thomas Campbell, and William Shadish. 2002. *Experimental and quasi-experimental designs for generalized causal inference*. Houghton Mifflin Boston.
- [6] Hakan Erdogmus, Maurizio Morisio, and Marco Torchiano. 2005. On the effectiveness of the test-first approach to programming. *IEEE Transactions on Software Engineering* 31, 3 (2005), 226–237.
- [7] Davide Fucci, Hakan Erdogmus, Burak Turhan, Markku Oivo, and Natalia Juristo. 2017. A Dissection of the Test-Driven Development Process: Does It Really Matter to Test-First or to Test-Last? *IEEE Transactions on Software Engineering* 43, 7 (2017), 597–614.
- [8] Davide Fucci, Giuseppe Scanniello, Simone Romano, Martin Shepperd, Boyce Sigweni, Fernando Uyaguari, Burak Turhan, Natalia Juristo, and Markku Oivo. 2016. An External Replication on the Effects of Test-driven Development Using a Multi-site Blind Analysis Approach. In *Proceedings of International Symposium on Empirical Software Engineering and Measurement (ESEM '16)*. ACM, Article 3, 10 pages.
- [9] Davide Fucci and Burak Turhan. 2013. A replicated experiment on the effectiveness of test-first development. In *Empirical Software Engineering and Measurement, 2013 ACM/IEEE International Symposium on*. IEEE, 103–112.
- [10] Davide Fucci, Burak Turhan, Natalia Juristo, Oscar Dieste, Ayse Tosun-Misirli, and Markku Oivo. 2015. Towards an operationalization of test-driven development skills: An industrial empirical study. *Information and Software Technology* 68 (2015), 82–97.
- [11] Donald E Harter, Chris F Kemerer, and Sandra A Slaughter. 2012. Does software process improvement reduce the severity of defects? A longitudinal field study. *IEEE Transactions on Software Engineering* 38, 4 (2012), 810–827.
- [12] Martin Höst, Björn Regnell, and Claes Wohlin. 2000. Using Students as Subjects—A Comparative Study of Students and Professionals in Lead-Time Impact Assessment. *Empirical Software Engineering* 5, 3 (2000), 201–214.
- [13] Andreas Jedlitschka, Marcus Ciolkowski, and Dietmar Pfahl. 2008. Reporting Experiments in Software Engineering. In *In Guide to Advanced Empirical Software Engineering*. Springer, 201–228.
- [14] Natalia Juristo and Ana M. Moreno. 2001. *Basics of Software Engineering Experimentation*. Kluwer Academic Publishers.
- [15] Roberto Latorre. 2014. Effects of developer experience on learning and applying unit test-driven development. *IEEE Transactions on Software Engineering* 40, 4 (2014), 381–395.
- [16] Jingyue Li, Nils B Moe, and Tore Dybå. 2010. Transition from a plan-driven process to Scrum: a longitudinal case study on software quality. In *Proceedings of the 2010 ACM-IEEE international symposium on empirical software engineering and measurement*. ACM, 13.
- [17] Artem Marchenko, Pekka Abrahamsson, and Tuomas Ihme. 2009. Long-term effects of test-driven development a case study. In *International Conference on Agile Processes and Extreme Programming in Software Engineering*. Springer, 13–22.
- [18] Laurie McLeod, Stephen G MacDonell, and Bill Doolin. 2011. Qualitative research on software development: a longitudinal case study methodology. *Empirical software engineering* 16, 4 (2011), 430–459.
- [19] Matthias M Müller and Andreas Höfer. 2007. The effect of experience on the test-driven development process. *Empirical Software Engineering* 12, 6 (2007), 593–615.
- [20] Hussan Munir, Misagh Moayyed, and Kai Petersen. 2014. Considering rigor and relevance when evaluating test driven development: A systematic review. *Information and Software Technology* 56, 4 (2014), 375–394.
- [21] Yahya Rafique and Vojislav B Misić. 2013. The effects of test-driven development on external quality and productivity: A meta-analysis. *IEEE Transactions on Software Engineering* 39, 6 (2013), 835–856.
- [22] Outi Salo and Pekka Abrahamsson. 2005. Integrating agile software development and software process improvement: a longitudinal case study. In *Empirical Software Engineering, 2005. 2005 International Symposium on*. IEEE, 10.
- [23] Julio Cesar Sanchez, Laurie Williams, and E Michael Maximilien. 2007. On the sustained use of a test-driven development practice at ibm. In *Agile Conference (AGILE), 2007*. IEEE, 5–14.
- [24] S. Shapiro and M. Wilk. 1965. An analysis of variance test for normality. *Biometrika* 52, 3-4 (1965), 591–611.
- [25] Helen Sharp, Cleidson deSouza, and Yvonne Dittrich. 2010. Using ethnographic methods in software engineering research. In *Proceedings of the 32nd ACM/IEEE International Conference on Software Engineering-Volume 2*. ACM, 491–492.
- [26] Forrest Shull, Grigori Melnik, Burak Turhan, Lucas Layman, Madeline Diep, and Hakan Erdogmus. 2010. What do we know about test-driven development? *IEEE software* 27, 6 (2010), 16–19.
- [27] Ayse Tosun, Oscar Dieste, Davide Fucci, Sira Vegas, Burak Turhan, Hakan Erdogmus, Adrian Santos, Markku Oivo, Kimmo Toro, Janne Jarvinen, and Natalia Juristo. [n. d.]. An industry experiment on the effects of test-driven development on external quality and productivity. *Empirical Software Engineering* 22, 6 ([n. d.]), 2763–2805.
- [28] Layman Lucas Diep Marie Erdogmus Hakan Shull Forrest Turhan, Burak. [n. d.]. How Effective is Test-Driven Development? In *Making Software: What Really Works, and Why We Believe It, year = 2006, pages =*, O'Reilly Media (Ed.).
- [29] Jari Vanhanen, Casper Lassenius, and Mika V Mantyla. 2007. Issues and tactics when adopting pair programming: A longitudinal case study. In *Software Engineering Advances, 2007. ICSEA 2007. International Conference on*. IEEE, 70–70.
- [30] Sira Vegas, Cecilia Apa, and Natalia Juristo. 2016. Crossover Designs in Software Engineering Experiments: Benefits and Perils. *IEEE Transactions on Software Engineering* 42, 2 (2016), 120–135.
- [31] Geert Verbeke, Geert Molenberghs, and Dimitris Rizopoulos. 2010. *Random Effects Models for Longitudinal Data*. Springer Berlin Heidelberg, 37–96.
- [32] Claes Wohlin, Per Runeson, Martin Hst, Magnus C. Ohlsson, Bjrn Regnell, and Anders Wessln. 2012. *Experimentation in Software Engineering*. Springer.
- [33] Robert K. Yin. 2009. Case study research: Design and methods (applied social research methods). *London and Singapore: Sage* (2009).