

Moving Towards Continuous Delivery and Continuous Experimentation: A Case Study in the B2B Domain

Olli Rissanen

Master's thesis
University of Helsinki
Department of Computer Science

Helsinki, October 12, 2014

| | | | |
|---|--|-----------------------------------|---|
| Tiedekunta — Fakultet — Faculty | | Laitos — Institution — Department | |
| Faculty of Science | | Department of Computer Science | |
| Tekijä — Författare — Author | | | |
| Olli Rissanen | | | |
| Työn nimi — Arbetets titel — Title | | | |
| Moving Towards Continuous Delivery and Continuous Experimentation: A Case Study in the B2B Domain | | | |
| Oppiaine — Läroämne — Subject | | | |
| Computer Science | | | |
| Työn laji — Arbetets art — Level | | Aika — Datum — Month and year | Sivumäärä — Sidoantal — Number of pages |
| Master's thesis | | October 12, 2014 | 65 |
| Tiivistelmä — Referat — Abstract | | | |
| <p>Delivering more value to the customer is the goal of every software company. To increase the value of a software product, a deep understanding of the customers behavior has to be achieved. Producing, testing and delivering software in real-time allows for faster reaction and shorter feedback loops. Delivering value in real-time requires a company to utilize real-time delivery, data-driven decisions and empirical evaluation of new products and features. This thesis is an exploratory case study investigating practices known as continuous delivery and continuous experimentation to tackle real-time delivery and data-driven decisions. Continuous delivery is a design practice where the software functionality is deployed continuously to customer environment. This process includes both automated builds and automated testing, but also automated deployment to a staging environment. Continuous experimentation is a design practice where the entire R&D process is guided by controlled experiments and feedback. Such experiments can include deploying multiple versions of software with different properties, measuring their success based on some criteria and then selecting the best candidate for further development. This paper conducts a deductive case study in a medium-sized software company to examine the development process of two different software products, and the transition towards continuous delivery and continuous experimentation. The scope of this study is in the B2B domain, where the customer is another company instead of an individual. As a result, challenges and best practices regarding the transition towards real-time value delivery in the B2B domain are identified.</p> <p>ACM Computing Classification System (CCS): General and reference → Experimentation TODO: add rest</p> | | | |
| Avainsanat — Nyckelord — Keywords | | | |
| Continuous delivery, Continuous experimentation, Development pipeline | | | |
| Säilytyspaikka — Förvaringsställe — Where deposited | | | |
| Muita tietoja — Övriga uppgifter — Additional information | | | |

Contents

| | | |
|----------|---|-----------|
| 1 | Introduction | 1 |
| 1.1 | Overview | 1 |
| 1.2 | Research questions | 3 |
| 1.3 | Scope and limitations | 4 |
| 2 | Literature: Continuous Delivery | 5 |
| 2.1 | Continuous delivery | 6 |
| 2.1.1 | Deployment pipeline | 8 |
| 2.1.2 | Impact on development model | 10 |
| 2.2 | Challenges regarding continuous delivery | 10 |
| 3 | Literature: Continuous Experimentation | 11 |
| 3.1 | Continuous experimentation | 11 |
| 3.1.1 | Components in continuous experimentation | 13 |
| 3.1.2 | Experimentation stages and scopes | 13 |
| 3.1.3 | Experimentation planning | 14 |
| 3.1.4 | Data collection | 15 |
| 3.1.5 | Data analysis | 15 |
| 3.1.6 | Roles in continuous experimentation | 16 |
| 3.2 | Continuous delivery and continuous experimentation collab- oration | 16 |
| 3.3 | Continuous experimentation in embedded software | 16 |
| 3.4 | Challenges regarding continuous experimentation | 18 |
| 3.5 | Motivation | 19 |
| 3.6 | Existing case studies | 20 |
| 4 | Case study as a research method | 21 |
| 4.1 | Case study | 21 |
| 4.2 | Qualitative research | 23 |
| 4.2.1 | Template analysis | 24 |
| 5 | Research design | 26 |
| 5.1 | Objective | 27 |
| 5.2 | Case descriptions | 28 |
| 5.2.1 | Dialog | 30 |
| 5.2.2 | CDM | 32 |
| 5.3 | Context | 33 |
| 5.3.1 | Development process | 34 |
| 5.4 | Methods | 37 |
| 5.4.1 | Data collection | 38 |
| 5.4.2 | Data analysis | 39 |

| | |
|---|-----------|
| 6 Findings | 40 |
| 6.1 Continuous delivery: B2B challenges | 40 |
| 6.1.1 Technical challenges | 40 |
| 6.1.2 Procedural challenges | 42 |
| 6.1.3 Customer challenges | 43 |
| 6.2 Continuous experimentation: B2B challenges | 44 |
| 6.2.1 Technical challenges | 44 |
| 6.2.2 Customer challenges | 46 |
| 6.2.3 A/B testing | 48 |
| 6.3 Continuous delivery's benefit to case company | 50 |
| 6.4 Continuous experimentations benefit to case company | 51 |
| 6.5 Implementing continuous experimentation as a development process | 55 |
| 6.5.1 Applying the general concept | 55 |
| 6.5.2 Applying the model | 55 |
| 7 Discussion | 58 |
| 7.1 Contribution | 59 |
| 7.2 Further research | 59 |
| A Interview questions | 64 |

It's hard to argue that Tiger Woods is pretty darn good at what he does. But even he is not perfect. Imagine if he were allowed to hit four balls each time and then choose the shot that worked the best. Scary good. – Michael Egan, Sr. Director, Content Solutions, Yahoo (Egan, 2007)

1 Introduction

1.1 Overview

To deliver value fast and to cope with the increasingly active business environment, companies have to find solutions that improve efficiency and speed. Agile practices [8] have increased the ability of software companies to cope with changing customer requirements and changing market needs [13]. To even further increase the efficiency, shortening the feedback cycle enables faster customer feedback. Additionally, providing software functionality to collect usage data can be used to guide the product development and to increase the value generated for customers.

The logical starting point of this thesis was an open issue by the case company: how do we validate whether a feature would be successful or not?

The requirement of developing value fast has led to the evolution of a new software development model [7]. The model is based on three principles: evolving the software by frequently deploying new versions, using customers and customer data throughout the development process and testing new ideas with the customers to drive the development process and increase customer satisfaction. Currently more and more software companies are transitioning towards this model, called innovation experiment systems.

A way to shorten the delivery cycles and to automate the delivery process is continuous delivery. It is an extension to continuous integration, where the delivery process is often entirely automated, and software functionality is deployed frequently to customer environment. While continuous integration defines a process where the work is automatically built, tested and frequently integrated to mainline [17], often multiple times a day, continuous delivery adds automated acceptance testing and deployment. Continuous delivery therefore attempts to deliver an idea to users as fast as possible by automating the deployment process.

Lean software development, which was inspired by the Toyota Production System [], is a development approach attempting to eliminate unnecessary work and to create value for the customer. The approach consists of seven principles: eliminate waste, build quality in, deliver fast, optimise the whole, create knowledge, defer commitment and respect people. Lean startup [38] is a method for developing products based on the lean software development principles. The key component of lean startup is the Build-Measure-Learn cycle, which guides the development process by defining a loop process where ideas are turned into products, which are then measured and the idea pivoted or persevered forwards.

Continuous experimentation attempts to validate that an idea is, in fact, a good idea. In continuous experimentation the organisation runs controlled experiments to guide the R&D process. The development cycle in continuous experimentation resembles the Build-Measure-Learn cycle of lean

startup [38]. The process in continuous experimentation is to first form a hypothesis based on a business goals and customer "pains" [7]. After the hypothesis has been formed, quantitative metrics to measure the hypothesis must be decided. After this a minimum viable product can be developed and deployed, while collecting the required data. Finally, the data is analyzed to validate the hypothesis.

Fagerholm et al. define that a suitable experimentation system must be able to release minimum viable products or features with suitable instrumentation for data collection, design and manage experiment plans, link experiment results with a product roadmap and manage a flexible business strategy [16]. The Build-Measure-Learn cycle starts with forming a hypothesis and building a minimum viable product (MVP) with tools for data collection. Once the MVP has been created, the data is analyzed and measured in order to validate the hypothesis. To persevere with a chosen strategy means that the experiment proved the hypothesis correct, and the full product or feature can be implemented. However, if the experiment proved the hypothesis wrong, the strategy is changed based on the implications of a false hypothesis.

In the context of this thesis, continuous experimentation refers to a development model following these principles:

- Data driven decisions
- Link product development and business aspect
- React to customers present needs
- Turn hypotheses into facts
- Steer development process

Bosch et al., introduce an innovation experiment system, in which the development process consists of frequently deploying new versions, using customers and customer usaged data in the development process and finally by focusing on innovation and testing ideas with customers to drive customer satisfaction and revenue growth [7]. "First, it frequently deploys new versions focusing on continuously evolving the embedded software in short cycles of 2-4 weeks" [15].

The study is an exploratory deductive case study, which aims to explore how continuous delivery and continuous experimentation can be integrated to the development process of the case company. The study specifically aims to identify the main requirements, problems and key success factors with regards to these approaches. Integrating these approaches to the development process requires a deep analysis of the current development process, seeking the current problems and strengths. Adopting both continuous delivery and continuous experimentaton also requires understanding the requirements of

continuous delivery and continuous experimentation, and restrictions caused by the developed software product.

In this research, the units under the study are two teams responsible for developing two different software products. By focusing on different products, a broader view on the application of these development approaches can be gained. The other product, a marketing automation called Dialog, is used through an extensive user interface. The other product under inspection is CDM, a Master Data Management [30] solution running as an integrated background application. The objective of this thesis can therefore be summarized as follows:

Research objective To analyze the software development process transition towards shorter feedback cycle and data-driven decisions from the perspective of the case company, with respect to continuous delivery and continuous experimentation practices.

1.2 Research questions

The broader research objective is explored here. The research questions are structured to progress from a higher-level perspective to a more detailed view. The first two research questions are needed to study the requirements and benefits of these practices to rationalize the decisions to embrace these practices. Then, the third research question corresponds to the operational adjustments required to adopt these practices.

RQ1: What are the B2B specific challenges of continuous delivery and continuous experimentation?

Software development practices and product characteristics vary based on the domain and delivery model. Typical B2C applications are hosted as SaaS applications, and accessed by users via a web browser. In the B2B domain, applications installed to customer environments are more common. Running experiments on customer environments also requires a new deployment of the software product, possibly causing downtime and changing some functionality of the software. The purpose of this research question is to identify the challenges in these practices in the B2B environment. The research question is answered by conducting an interview to discover B2B specific challenges, and using the available literature on continuous delivery and continuous experimentation to draw conclusions on the environmental differences.

RQ2: How do continuous delivery and continuous experimentation benefit the case company?

To rationalize the decision to adopt continuous delivery and experimentation in a company, the actual benefits to the business have to be identified.

Sections of the interview aim to identify the current perceived problems of the case company related to deployment, product development, collecting feedback and guiding the development process. This question is answered by comparing these problems to the benefits of these approaches found from the literature review, with respect to the challenges identified in RQ1.

RQ3: How can continuous experimentation be systematically organised in the B2B domain?

When the benefits and requirements are identified and the process is decided to be approved by a company, it has to be adopted and applied in practice. This research question investigates how to apply these practices in the case company, and aims to answer to questions such as what, how, who and when. This question is answered by applying Fagerholm et al. model [16] to two different teams and software products within the case company, and addressing the challenges found in the first research question. Then the deviations in the model are identified, and possible solutions to these deviations provided.

1.3 Scope and limitations

This study focuses on the development process of the case company and its two software products in the B2B domain. Case studies only allow analytic generalisations, and the findings cannot be directly generalised to other cases. This limitation applies especially to applying continuous delivery and continuous experimentation in practice, since software products and customers inevitably differ.

The background theory is drawn from literature focusing on continuous delivery, continuous experimentation, lean startup and innovation experiment systems.

The empirical research of this thesis has been carried out in the context of a medium sized information technology company, Steeri Oy. The company specializes in customer data on multiple fronts: CRM, BI, Master Data Management and data integrations, among others. The empirical research has applied case study methodology to offer deep understanding of the phenomena in their natural context. The tradeoff of a case study research is that it does not provide much room for direct generalizations outside the empirical context.

This thesis is organized as follows. The first chapter is the introduction at hand. The second and third chapter summarize the relevant literature and theories to position the research and to educate the reader on the body of knowledge and where the contributions are intended. The fourth chapter briefly explains the usage of case study and qualitative research as research methods, to provide background for the decision to utilize them in this

particular study. The fifth chapter presents the research design: objectives, context and applied methods. The findings are then presented in the sixth chapter, organized according to the research questions. Finally, the seventh chapter summarizes the thesis by explaining the contributions to current research and introducing further research avenues.

2 Literature: Continuous Delivery

Agile software development emerged from the iterative software development methods, with early implementations being development in the 1990's [12]. The general principles of agile software development are breaking the development to short iterations, maintaining a short feedback loop with the customer by having a customer representative, valuing working software over documentation and responding to change. Many agile models, such as Scrum [41], Extreme Programming [6] and Feature-Driven Development [36] further define their own strategies based on the agile principles. Scrum, for example, defines a framework for managing product development, without much focus on technical tools. Extreme Programming, on the other hand, includes many technical features, such as always programming in pairs, doing extensive code review and unit testing.

Lean manufacturing, derived from the Toyota Production System [35], is a streamlined production philosophy that aims to identify the expenditure of resources not used in creating value for the end customer, and eliminate these activities. Lean software development [37] has been developed based on the same idea. The purpose of lean software development is to eliminate waste, or anything that doesn't bring value to the customer, and to pursue perfection through smoothing the development flow. The principles in lean software development are eliminate waste, amplify learning, decide as late as possible, deliver as fast as possible, engage everyone, create quality and see the whole.

The software development models have moved towards even shorter and shorter feedback loops in the recent years. Lean startup [38] is a business and product development method with focus on business-hypothesis-driven experimentation and iterative product releases. A primary purpose of lean startup is to assess the specific demands of consumers and to fulfill those demands with the least possible resources. This is achieved by using a Build-Measure-Learn-loop, which is used to learn the customers actual needs and consecutively to steer business and product strategies to the correct direction. The tight customer feedback loop ensures that the product development doesn't include features and services not wanted by the customer.

Continuous delivery is a design practice where the software delivery process is automated, and thus the software can be deployed on short notice.

Olsson et al. have investigated the evolution path of software companies, and have identified continuous delivery and "R&D as an Experiment System" as final stages[34]. Continuous delivery is also seen as a prerequisite for continuous experimentation, as it provides a way to quickly deploy new features.

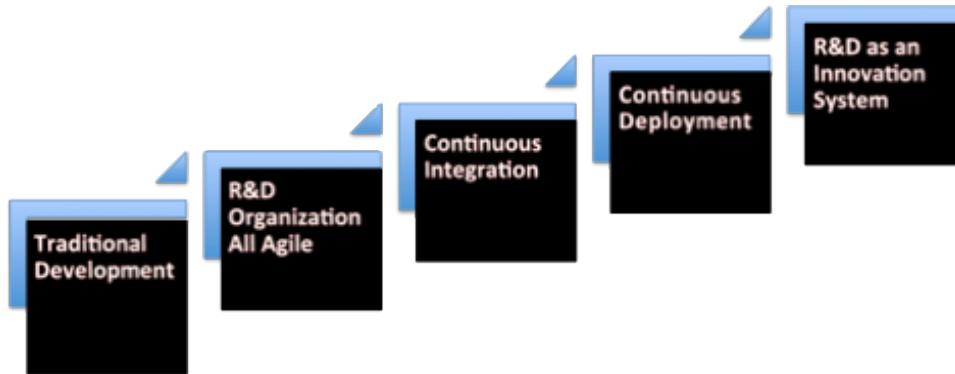


Figure 1: Organization evolution path [34].

In this chapter we first introduce the general principles of continuous delivery and continuous deployment. Then we introduce the software deployment pipeline and the effects of continuous delivery on the development model. Finally, we list challenges faced when applying continuous delivery.

2.1 Continuous delivery

To be able to react to customers changing requests and manage multiple customer environments, the software product has to be deployed frequently and with different configurations. As the case company is working in an agile manner, requests change and adjust often. To improve the software product, the case company has a need for a deployment approach that has the following benefits:

- Fast feedback on changes
- Automation of repetitive actions
- Validating that deployed code has passed automated tests
- Traceable. large history records and changelogs
- Configurable per customer environment

Fast feedback is used to validate the functionality of the software and to ensure that quality requirements are met. Automation of repetitive actions ensures that the actions are performed exactly as instructed, without

room for manual error. Reducing the amount of manual work also improves efficiency. With traceable history records, troubleshooting process can be shortened. With changelogs, the customer can be kept informed of the changes in new versions. A design practice that fills the aforementioned benefits is continuous delivery.

Continuous delivery is an extension to continuous integration, where the software functionality is deployed frequently at customer environment. While continuous integration defines a process where the work is automatically built, tested and frequently integrated to mainline [17], often multiple times a day, continuous delivery adds automated acceptance testing and deployment to a staging environment. The purpose of continuous delivery is that as the deployment process is automated, it reduces human error, documents required for the build and increases confidence that the build works [20].

In an agile process software release is done in periodic intervals [8]. Compared to waterfall model it introduces multiple releases throughout the development. Continuous delivery, on the other hand, attempts to keep the software ready for release at all times during development process [20]. Instead of stopping the development process and creating a build as in an agile process, the software is continuously deployed to customer environment. This doesn't mean that the development cycles in continuous delivery are shorter, but that the development is done in a way that makes the software always ready for release. Continuous delivery differs from continuous

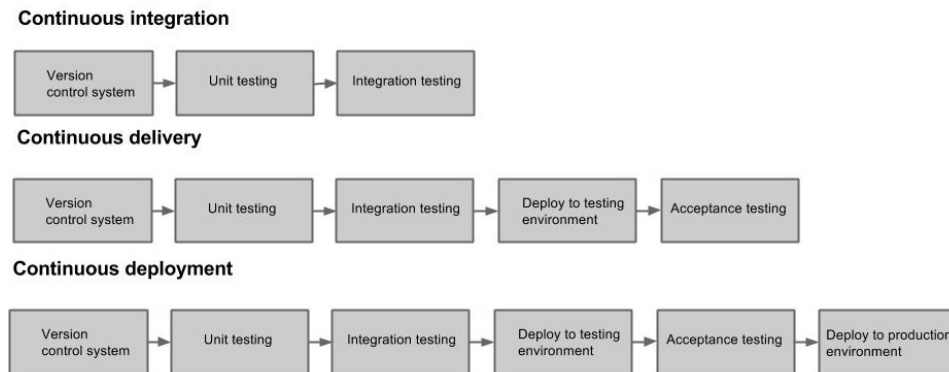


Figure 2: Continuous integration, delivery and deployment.

deployment. Refer to Fig. 2 for a visual representation of differences in continuous integration, delivery and deployment. Both continuous delivery and continuous deployment include automated deployment to a staging environment. Continuous deployment includes deployment to a production environment, while in continuous delivery the deployment to a production environment is done manually. The purpose of continuous delivery is to

prove that every build is deployable [20]. While it necessarily doesn't mean that teams release often, keeping the software in a state where a release can be made instantly is often seen beneficial.

Holmström Olsson et al. have researched the transition phase from continuous integration to continuous delivery, and have identified barriers that companies need to overcome to achieve the transition [34]. One such barrier is the custom configuration at customer sites. Maintaining customized solutions and local configurations alongside the standard configurations creates issues. The second barrier is the internal verification loop, that has to be shortened not only to develop features faster but also to deploy fast. Finally, the lack of transparency and getting an overview of the status of development projects is seen as a barrier.

Olsson et al. define continuous deployment as a step in the typical evolution path for companies [34]. The way to deploy versions continuously allows for faster customer feedback, ability to learn from customer usage data and to focus on features that produce real value for the customer.

Olsson et al. state that in order to shorten the internal verification loop, different parts in the organization should also be involved, especially the product management as they are the interface to the customer. Olsson et al. also note that finding a pro-active lead customer is also essential to explore and form a new engagement model.

2.1.1 Deployment pipeline

An essential part of continuous deployment is the deployment pipeline, which is an automated implementation of an application's build, deploy, test and release process [20]. A deployment pipeline can be loosely defined as a consecutively executed set of validations that a software has to pass before it can be released. Common components of the deployment pipeline are a version control system and an automated test suite.

Humble and Farley define the deployment pipeline as a set of stages, which cover the path from a committed change to a build [20]. Refer to Fig. 3 for a graphical representation of a basic deployment pipeline. The commit stage compiles the build and runs code analysis, while acceptance stage runs an automated test suite that asserts the build works at both functional and nonfunctional level. From there on, builds to different environments can be deployed either automatically or by a push of a button.

Humble et al. define four principles that should be followed when attempting to automate the deployment process [21]. The first principle states that "Each build stage should deliver working software". As software often consists of different modules with dependencies to other modules, a change to a module could trigger builds of the related modules as well. Humble et al. argue that it is better to keep builds separate so that each discrete module could be built individually. The reason is that triggering other builds can

be inefficient, and information can be lost in the process. The information loss is due to the fact that connection between the initial build and later builds is lost, or at least causes a lot of unnecessary work spent in tracing the triggering build.

The second principle states that "Deploy the same artifacts in every environment". This creates a constraint that the configuration files must be kept separate, as different environments often have different configurations. Humble et al. state that a common anti-pattern is to aim for 'ultimate configurability', and instead the simplest configuration system that handles the cases should be implemented. Another principle, which is the main element

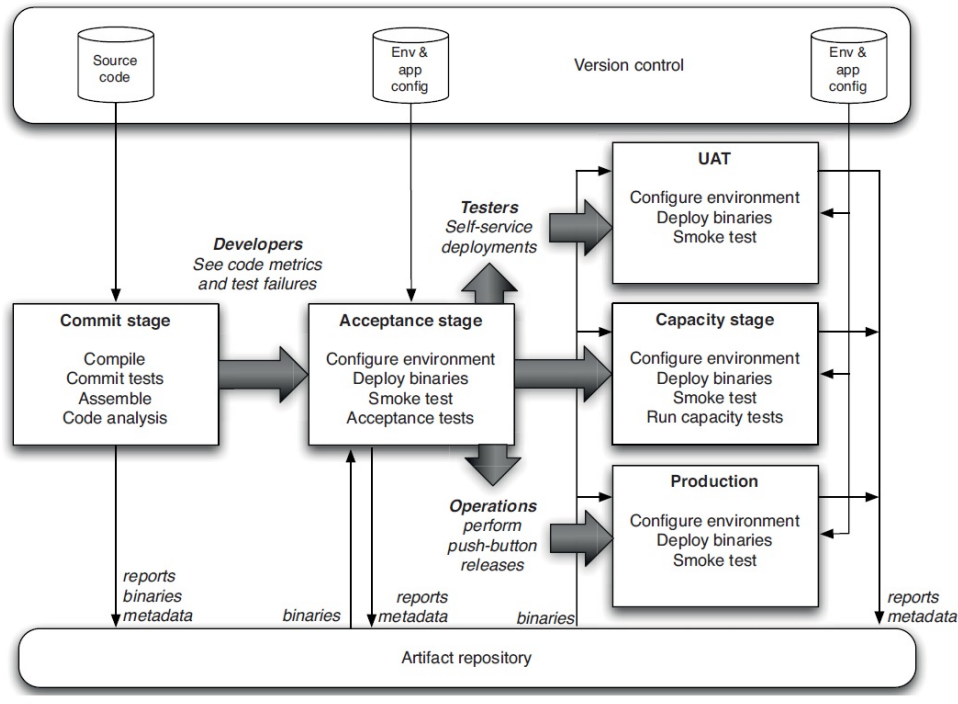


Figure 3: A basic deployment pipeline [20].

of continuous deployment, is to "Automate testing and deployment". Humble et al. argue that the application testing should be separated out, such that stages are formed out of different types of tests. This means that the process can be aborted if a single stage fails. They also state that all states of deployment should be automated, including deploying binaries, configuring message queues, loading databases and related deployment tasks. Humble et al. mention that it might be necessary to split the application environment into *slices*, where each slice contains a single instance of the application with predetermined set of resources, such as ports and directories. *Slices* make it possible to replicate an application multiple times in an environment, to

keep distinct version running simultaneously. Finally, the environment can be smoke tested to test the environments capabilities and status.

The last principle states "Evolve your production line along with the application it assembles". Humble et al. state that attempting to build a full production line before writing any code doesn't deliver any value, so the production line should be built and modified as the application evolves.

2.1.2 Impact on development model

A picture of the typical development process in continuous delivery is shown in Fig. 4. After the team pushes a change to the version control system, the project is automatically built and tests are triggered stage by stage. If a test stage fails, feedback is given and the deployment process effectively cancelled. In a continuous delivery process, the last stages are approved and activated manually, but in a continuous deployment process the last stages are triggered automatically as well.

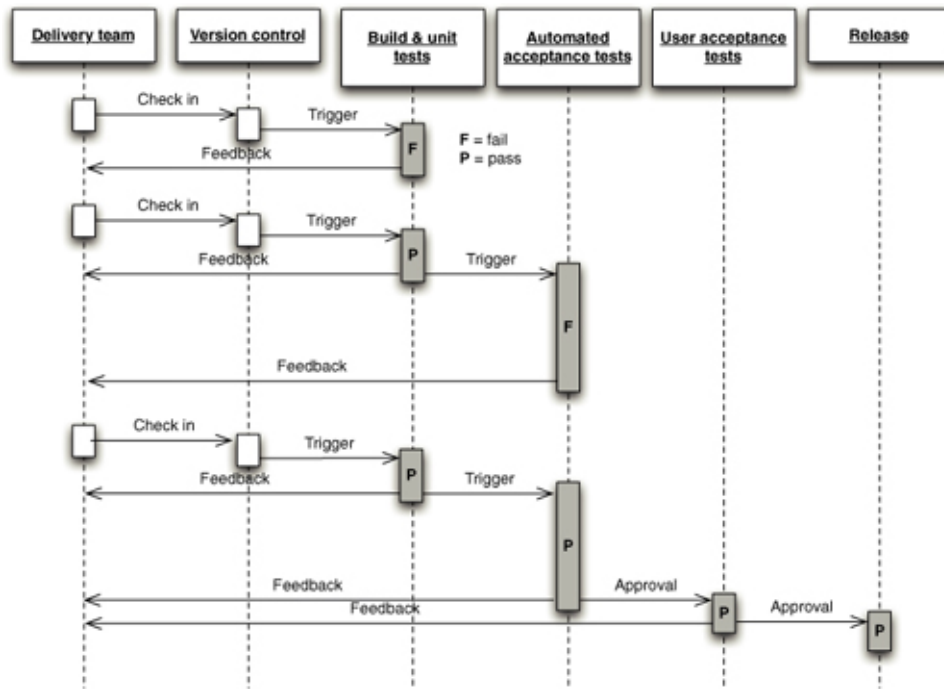


Figure 4: Components of the development process [20].

2.2 Challenges regarding continuous delivery

Neely and Stolt found out that with a continuous flow, Sales and Marketing departments lost the track of when features are released [33]. Continuous delivery requires the company as a whole to understand the continuous flow

process. Neely and Stolt solved this challenge by providing mechanisms by which the parties in need of knowledge regarding release times could track the work in real time [33].

-Technical implementation: Automate testing -Education on the subject: Developer training costs -Active participation of employees

3 Literature: Continuous Experimentation

3.1 Continuous experimentation

Continuous experimentation is a term created by Dane McKinley, a Principal Engineer at Etsy [2]. McKinley defines the key aspects of continuous experimentation as making small, measurable changes, staying honest and preventing the developers from breaking things. With staying honest, McKinley implies that the product releases are tracked and measured, so that it's possible to tell whether things went worse or better. McKinley also states that design and product process must change to accommodate experimentation, and that experimenting should be done with minimal version of the idea. When experimentations are implemented and measured, counterintuitive results can be found, and planning to be wrong should be considered [2].

"Most organizations have many ideas, but the return-on-investment (ROI) for many may be unclear and the evaluation itself may be expensive." [28]

While continuous experimentation is only loosely defined by McKinley, it poses resemblance to multiple development models. One of such models is the Innovation Experiment System [7]. The key principles in these development models are data-driven decisions, linking product development and business aspect, reacting to customers present needs, turning hypotheses into facts, failing fast and steering the development process. Fagerholm et al. define that a suitable experimentation system must be able to release minimum viable products or features with suitable instrumentation for data collection, design and manage experiment plans, link experiment results with a product roadmap and manage a flexible business strategy [16]. The Build-Measure-Loop of Lean startup is similar in the sense that controlled experiments drive the development.

In Lean startup methodology [38] experiments consist of Build-Measure-Learn cycles, and are tightly connected to visions and the business strategy. The purpose of a Build-Measure-Learn cycle is to turn ideas into products, measure how customers respond to the product and then to either pivot or persevere the chosen strategy. The cycle starts with forming a hypothesis and building a minimum viable product (MVP) with tools for data collection. Once the MVP has been created, the data is analyzed and measured in order to validate the hypothesis. To persevere with a chosen strategy means that the experiment proved the hypothesis correct, and the full prod-

uct or feature can be implemented. However, if the experiment proved the hypothesis wrong, the strategy is changed based on the implications of a false hypothesis.

Holmström Olsson et al. have researched the typical evolution path of companies [34]. The final stage of the evolution phase is when development is driven by the real-time customer usage of software. Defined as "R&D as an 'experiment system'", the entire R&D system acts based on real-time customer feedback and development is seen as addressing to customers present needs. At the same time deployment of software is seen as a "starting point for 'tuning' of functionality rather than delivery of the final product". While the evolution from continuous delivery to innovation system wasn't explored, Olsson et al. anticipate that the key actions required are the automatic data collection components and capability to effectively use the collected data.

An experiment is essentially a procedure to confirm the validity of a hypothesis. In software engineering context, experiments attempt to answer questions such as which features are necessary for a product to succeed, what should be done next and which customer opinions should be taken into account. According to Jan Bosch, "The faster the organization learns about the customer and the real world operation of the system, the more value it will provide" [7]. Most organizations have many ideas, but the return-on-investment for many may be unclear and the evaluation itself may be expensive [28].

Continuous delivery attempts to deliver an idea to users as fast as possible. Continuous experimentation instead attempts to validate that it is, in fact, a good idea. In continuous experimentation the organisation runs controlled experiments to guide the R&D process. The development cycle in continuous experimentation resembles the build-measure-learn cycle of lean startup. The process in continuous experimentation is to first form a hypothesis based on a business goals and customer "pains" [7]. After the hypothesis has been formed, quantitative metrics to measure the hypothesis must be decided. After this a minimum viable product can be developed and deployed, while collecting the required data. Finally, the data is analyzed to validate the hypothesis.

Jan Bosch has widely studied continuous experimentation, or innovation experiment systems, as a basis for development. The primary issue he found is that "experimentation in online software is often limited to optimizing narrow aspects of the front-end of the website through A/B testing and inconnected, software-intensive systems experimentation, if applied at all, is ad-hoc and not systematically applied" [7]. The author realized that for different development stages, different techniques to implement experiments and collect customer feedback exist. Bosch also introduces a case study in which a company, Intuit, adopted continuous experimentation and has increased both the performance of the product and customer satisfaction.

3.1.1 Components in continuous experimentation

Kohavi et al. investigate the practical implementations of controlled experiments on the web [28], and state that the implementation of an experiment involves two components. The first component is a randomization algorithm, which is used to map users to different variants of the product in question. The second component is an assignment method which, based on the output of the randomization algorithm, determines the contents that each user are shown. The observations then need to be collected, aggregated and analyzed to validate a hypothesis. Kohavi et al. also state that most existing data collection systems are not designed for the statistical analyses that are required to correctly analyze the results of a controlled experiment.

The components introduced by Kohavi et al. are aimed primarily for A/B testing on websites. Three ways to implement the assignment methods are shown. The first one is traffic splitting, which directs users to different fleet of servers. An alternative methods is server-side selection, in which API calls invoke the randomization algorithm and branch the logic based on the return value. Last alternative is a client-side selection, in which the front-end system dynamically modifies the page to change the contents. Kohavi et al. state that the client-side selection is easier to implement, but it severely limits the features that may be subject to experimentation. Experiments on back-end are nearly impossible to implement in such manner.

A different model for general experiments, not limited to A/B testing, has also been recently researched. While Fagerholm et al. state that "A detailed framework for conducting systematic, experiment-based software development has not been elaborated" [16], have they created an initial model consisting of required building blocks for a continuous experimentation system, depicted in Fig. 4.

The primary components of the introduced infrastructure are experiment database, analytics tools and instrumentation. The database stores raw data collected from the instrumentation and information related to the experiments, such as the logic for conducting the experiments and the experiment results. The analytics tools are used to perform data analysis on raw data from the experiment database to produce results regarding the experiment. The instrumentation is a part of the delivered software used to gather data regarding pre-defined metrics.

3.1.2 Experimentation stages and scopes

Fig. 7 introduces different stages and scopes for experimentation. For each stage and scope combination, an example technique to collect product performance data is shown. As startups often start new products and older companies instead develop new features, experiments must be applied in the correct context. Bosch states that for a new product deployment, putting

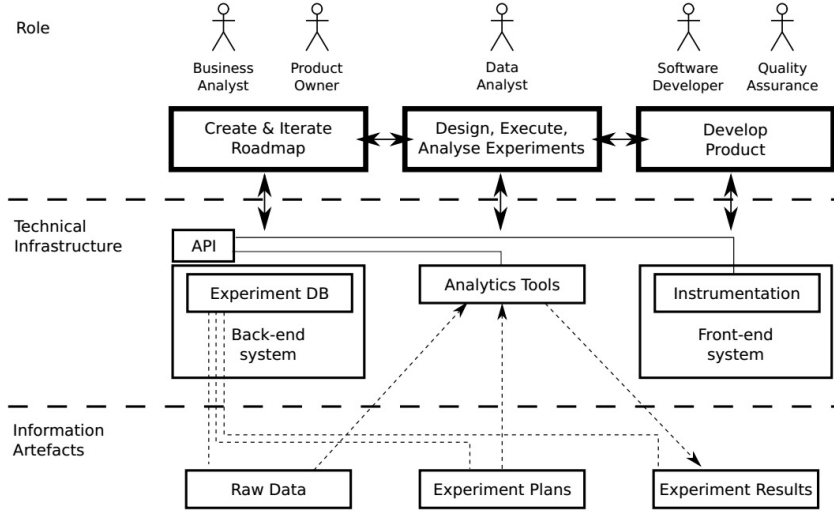


Figure 5: Continuous experimentation infrastructure [16].

a minimal viable product as rapidly as possible in the hands of customers is essential [7]. After the customers can use the product, it is often not yet monetizable but is still of value to the customer. Finally, the product is commercially deployed and collecting feedback is required to direct R&D investments to most valuable features.

| | Feature Optimization | New Feature Development | New Product Development |
|---------------------------|-----------------------|-------------------------|--------------------------|
| Pre-Development | Upsell advertising | Participatory design | Collaborative innovation |
| Non-Commercial Deployment | A/B testing | Feature alpha | Usage behavior tracking |
| Commercial Deployment | Multi-variate testing | Usage metrics | Cross-sell actions |

Figure 6: Scopes for experimentation [7].

3.1.3 Experimentation planning

In continuous experimentation literature, the common way to plan experiments is to utilize the product strategy or roadmap [38, ?]. In the Lean Startup Build-Measure-Learn loop, the Learn phase is essentially a decision to whether to persevere or pivot the product strategy. However, this requires that the Measure phase is extensively planned in advance, so that a clear hypothesis exists. Essentially, the Build-Measure-Learn loop is very

similar to the Plan-Do-Study-Adjust cycle [?]. The single exception is that the Build-Measure-Learn loop is based on an assumption that the customer needs cannot be known even with planning.

Fagerholm et al. emphasize the importance of analysing the data to test business hypotheses and assumptions, which have been created from the roadmap. These proven hypotheses can therefore be used to react to customers needs.

3.1.4 Data collection

TODO: finish this "Active customer feedback is concerned with surveys and other mechanisms where the customer is aware that he or she is providing feedback. Passive feedback and usage data is collected while the customer is using the system. Examples include the amount of time a user spends using a feature, the relative frequency of feature selections, the path that the user takes through the product functionality, etc. The low cost and ease of data collection leads to the next major difference between IES-based and traditional software." [7]

To collect, aggregate and analyze the observations, raw data has to be recorded. According to Kohavi et al., some raw data could be for example page views, clicks, revenue, render time and customer-feedback selections [28]. The data should also be annotated to an identifier, so that conclusions can be made from it. Kohavi et al. present three different ways for collecting raw data. The first solution is to simply use an existing data collection tool, such as Webmetrics. However, most data collection systems aren't designed for statistical analyses, and the data might have to be manually extracted to an analysis environment. A different approach is local data collection, in which a website records data in a local database or log files. The problem with local data collection is that each additional source of data, such as the back-end, increases the complexity of the data recording infrastructure. The last model is a service-based collection, in which service calls to a logging service are placed in multiple places. This centralizes all observation data, and makes it easy to combine both back-end and front-end logging.

An important aspect of collecting data is to pick the correct metrics. One of the common pitfalls for experiments is to choose the wrong overall estimation criteria [9]. TODO: follow this up

3.1.5 Data analysis

To analyze the raw data, it must first be converted into metrics which can then be compared between the variants in a given experiment. An arbitrary amount of statistical tests can then be run on the data with analytics tools in order to determine statistical significance.

"organisation needs a workfore to analyze qualitative + quantitative data"

3.1.6 Roles in continuous experimentation

Fagerholm et al. define five different roles in continuous experimentation: Business Analyst, Product Owner, Data Analyst, Software Developer and Quality Assurance [16]. The business analyst along with the product owner are responsible for creating and updating the strategic roadmap, which is the basis of the hypotheses. The basis for decisions are the existing experimental plans and results stored in a database. A data analyst is used to analyze the existing experiments and results and to create assumptions from the roadmap. The analyst is also responsible for the design and execution of experiments. The data analyst is in tight collaboration with the software developer and quality assurance, who are responsible for the development of MVPs and MVFs. The software designers create the necessary instrumentations used to collect the data required by the analyst.

Another approach for the roles is introduced by Adams et al., who conducted a case study on the implementation of Adobe's Pipeline. Adobe's pipeline is process that is based on the continuous experimentation approach [5, 3]. The innovation pipeline introduced in the study includes the following roles: designer, developer, researcher, product manager [3].

3.2 Continuous delivery and continuous experimentation collaboration

As the experiments are run in a regular fashion, integrating experiments to the deployment pipeline should be considered. This requires changing the development process in such fashion that functionality is developed based on some actual data. The components required to support continuous experimentation include tools to assign users to treatment and control groups, tools for data logging and storing, and analytics tool for conducting statistical analyses.

In the organization evolution path [34] continuous deployment is seen as a prerequisite for the experimental system.

3.3 Continuous experimentation in embedded software

The purpose of analysing experimentation in embedded software in the context of this thesis is that B2B applications and embedded applications share similarities in the release cycle. The innovation ideas for embedded software are prioritized as a part of the yearly release cycle. This resembles the release model found in the case company, where some customers prefer to have new versions shipped every three months. By analyzing how experiments

are implemented and the data collected in embedded software, information on how long and short release cycles differ can be gained.

Eklund and Bosch have applied the innovation experiment system to embedded software [15]. In the study an embedded architecture for realising an innovation experiment system is presented. The study concentrates on what to evaluate in experiments in embedded software. The infrastructure for enabling the innovation experiment system is shown in Fig. IV.

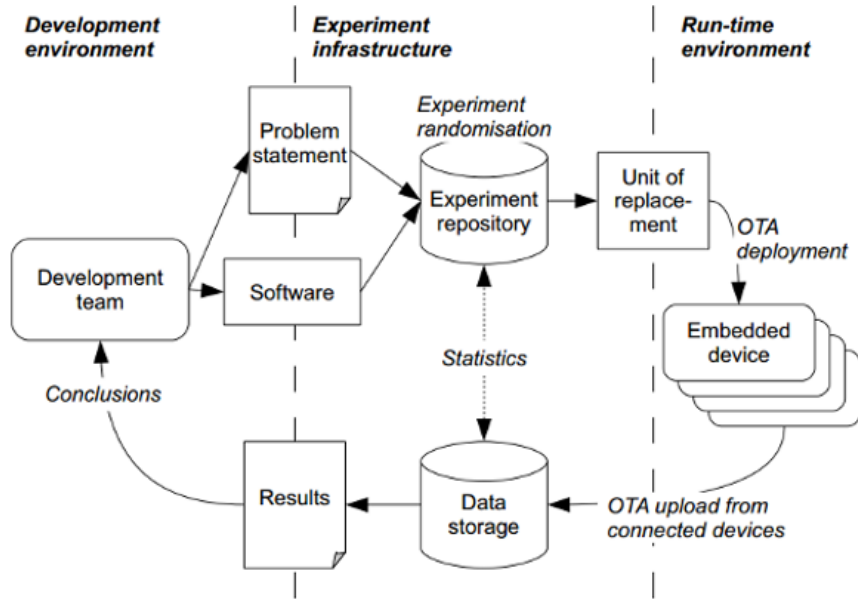


Figure 7: Continuous experimentation architecture for embedded software [15].

TODO: loppuun

"Improvement occurs in individual software parts, but the underlying design concept remain mostly unchanged [3]" "The experiment infrastructure allows developers to deploy new software and collect data how it behaves in areal-world settings being used by actual users. The infrastructure support deployment of software experiments and collection of data over-the-air on a scale comparable to the entire customer base". "The infrastructure supports with automated randomisation and factorial designs [6] sufficient to draw statistical conclusions from the experimental scenarios." "The experiment manager architecture, seen in Figure 3, supports the deployment of multiple experimental software parts to the same device and autonomously controls when to run which experiment, even allowing for local A/B-testing. Measurements and analysis is done on-board in real-time. The experiment

scenario to be answered is implemented on the embedded device (i.e. how long does it take to . . .)"

The study also lists possible experiment scenarios supported by the architecture [15]. Possible experiment cases with the architecture are as follows:

- How long does it take to . .
- Which of . . . is most often used/accessed/. . .
- Identify behaviour that is not intended, e.g. menu selection followed by "back" indicates that the user made a mistake.
- Are there any features that are not used?
- Be able to evaluate competing designs based on the answers above, i.e. A/B testing (AKA. split testing).

3.4 Challenges regarding continuous experimentation

It is commonly found in the studies that the practiced experimentation is limited to optimizing narrow aspects of the front-end through A/B testing [7]. Additionally, experimentation is not often systematically applied [7]. Typically the imbalance between number of ideas that exist in the organization and the number of concepts that are tested with customers is enormous [7].

Some companies use traditional metrics such as the Net Promoter Score [?] to gather feedback. These metrics however fail to guide the development process by providing timely feedback, and instead are backward looking and focusing on the whole product [7].

One of the biggest challenges is also the architectural requirements set by the product to allow experimentation. According to Jan Bosch: "Collecting customer or performance data early in the innovation process requires changes to the R&D processes, as customers need to be involved much earlier and deeper in the process. This also requires architectural changes of the product and platform to determine customer preference and interest." [7].

Crook et al. investigate the common pitfalls encountered when running controlled experiments on the web [9], which are listed in Table II. The Overall Evaluation Criteria (OEC) used in the first pitfall is a quantitative measure of the experiment's objective. In experimental design, Control is a term used for the existing feature or a process, while a Treatment is used to describe a modified feature or a process. As a motivator for the first pitfall, Crook et al. introduce an experiment where the OEC was set to the time spent on a page which contains articles. The OEC increased in the experiment implementation, satisfying the objective. However, it was soon realized that the longer time spent on a page might have been caused

by confusion of the users, as a newly introduced widget was used less often than a previous version of it. Aside from picking a correct OEC, common pitfalls deal with the correct use of statistical analysis, robot users such as search engine crawlers, and the importance of audits and control.

| | |
|-----------|--|
| Pitfall 1 | Picking an OEC for which it is easy to beat the control by doing something clearly “wrong” from a business perspective |
| Pitfall 2 | Incorrectly computing confidence intervals for percent change and for OECs that involve a nonlinear combination of metrics |
| Pitfall 3 | Using standard statistical formulas for computations of variance and power |
| Pitfall 4 | Combining metrics over periods where the proportions assigned to Control and Treatment vary, or over subpopulations sampled at different rates |
| Pitfall 5 | Neglecting to filter robots |
| Pitfall 6 | Failing to validate each step of the analysis pipeline and the OEC components |
| Pitfall 7 | Forgetting to control for all differences, and assuming that humans can keep the variants in sync |

Table 1: Pitfalls to avoid when running controlled experiments on the web [9].

-Roles -Environment -Deployment

3.5 Motivation

The software industry is growing due to increasing digitalisation, and it is especially important for new and old companies alike to identify the problems relevant to customers and to stay flexible. Identifying the customers problems allows the companies purely focus on creating value, while staying flexible means quickly diverting development efforts towards currently desired features. The importance of continuous innovation is increasingly recognized in the literature [44]. Companies are also increasingly interested in utilizing real-time customer usage of the software to shorten feedback loops [34].

The process of guiding product development with continuous experiments diminishes guesswork and instead provides a systematic approach. Avinash Kaushik states in his article Experimentation and Testing primer [1] that "80% of the time you/we are wrong about what a customer wants". Mike Moran also found similar results in his book Do It Wrong Quickly,

stating that Netflix considers 90% of what they try to be wrong [32]. A way to tackle this issue is to adopt a process of continuous experimentation. Ron Kohavi at Microsoft wrote "Features are built because teams believe they are useful, yet in many domains most ideas fail to improve key metrics" [27]. Additionally, Regis Hadianis from Quicken Loans wrote that "in the five years I've been running tests, I'm only about as correct in guessing the results as a major league baseball player is in hitting the ball. That's right - I've been doing this for 5 years, and I can only "guess" the outcome of a test about 33% of the time!" [31].

3.6 Existing case studies

Existing case studies regarding continuous experimentation have been conducted with varying aspects under concentration. Olsson et al. conducted a multiple-case study to explore companies moving towards continuous deployment and continuous experimentation [34]. In the transition phase from continuous integration stage to continuous deployment, the study investigates a telecommunications company with a highly distributed organization. The study considers transition from continuous delivery to continuous experimentation as a future research interest.

Fagerholm et al. analytically derive a continuous experimentation model and apply it in practice in a startup company [16]. The study introduces an initial continuous experimentation model, examines the preconditions for setting up an experimental system and describes the building blocks required for such a system. The case company in this study is a relatively new startup company, and the company operates in the B2C domain. The study states that it is yet not clear how to determine the suitability of an experimental approach in specific situations. This is in line with the goal of this thesis: analyze the continuous experimentation approach in the B2B domain with two different software products.

Steiber et al. [44] have studied continuous innovation in Google. The study defines continuous innovation as "an experimental iterative process that operates successively to solve problems in markets characterized by turbulence, uncertainty and complex interactions". The study explores the organizational characteristics for continuous innovation in rapidly changing industries. The study finds that the entire organisation of Google is involved in continuous innovation. The study also states that there is a need for a more comprehensive analytical framework for continuous innovation, and also a need to study how to organize continuous innovation and continuous improvements.

Adams et al. conducted a case study on the implementation of Adobe's Pipeline, a process that is based on the continuous experimentation approach [5, 3]. The innovation pipeline introduced in the study consists of the following steps:

- Collect user problems
- Identify problems worth solving
- Prototype possible solutions
- Live user testing + public releases
- Refine, pivot or kill the idea
- Build products that understand the user

The pipeline team consists of the following roles: designer, developer, researcher and product manager. The general principles found in the study are as follows:

- The idea will change
- We're only right when the market tells us so
- Create fast. improve constantly
- Find the simplest thing that could possibly work
- Define "success" before we build
- Fail forward

The study is carried out in the company with over 10,000 employees. The study states that "there is a wealth of understanding about how small companies can innovate in a low-cost way in order to make small products viable".

Kohavi et al. have widely studied online experimentation at microsoft [29]. The experiments have been mostly performed in Microsoft's web sites with A/B testing, and the efforts have been proven succesful. This success further motivates the goal of this thesis of applying experiments in a different domain and different software products. However, as the study only consists of A/B testing, this thesis also inspects other kind of experiments performed in a continuous manner.

4 Case study as a research method

4.1 Case study

Case study is a way to collect data through observation and to test theories in an unmodified setting [46]. The case in a case study is the situation, individual, group or organization that the researchers are interested in [39]. Kitchenham et al. found case studies important for industrial evaluation of software engineering methods and tools [25]. Runeson and Höst define it

as a suitable research method for software engineering, since it studies contemporary phenomena in its natural context [40]. It might not be possible to generate causal relationships from case studies as compared to controlled experiments. However, they provide a deeper understanding of the unit under study [40]. An opposite type of research would be a formal experiment, with a narrow focus and control on variables.

Case studies can serve different purposes. Robson defines four purposes as Exploratory, Descriptive, Explanatory and Improving [39]. In an exploratory case study, the purpose is to seek new insights, find out what is happening and generate ideas and hypotheses for new research. Descriptive case study attempts to portray a situation or a phenomenon. Explanatory case study attempts to seek an explanation to a problem or situation, mostly in a causal relationship. An improving case study tries to improve an aspect of the studied phenomenon.

Different approaches to the order of data collection, analysis and generalization can be categorized into inductive, deductive and abductive approaches [11]. In an inductive case study the research is conducted by first collecting the data, then looking for patterns and forming theories that explain those patterns. Therefore an inductive approach moves from data to theory. A deductive approach, on the other hand, is in reversed order; the approach moves from general level to a more specific level. A deductive approach begins with studying what others have done, and then testing hypotheses that emerge from those theories. An abductive approach starts with the consideration of facts, or particular observations. These observations are then used to form hypotheses that relate them to a fact or a rule that accounts for them. The facts are therefore correlated into a more general description, that relates them to a wider context.

Runeson and Höst define the case study structure in five steps [40]. First the case study is designed: objectives are defined and the study is planned. Then procedures and protocols for data collection are defined. Then the evidence is collected by executing the data collection on the studied case. The collected data is then analysed, and finally the results are reported.

Triangulation is a way to increase the reliability and validity of the findings. Triangulation means using different data collection methods or angles, and providing a wider picture of the case. Triangulation is especially important for qualitative data, but can also be used for quantitative data to compensate for measurement or modeling errors [40]. Stake defines four ways to apply triangulation: Data triangulation, Observer triangulation, Methodological triangulation and Theory triangulation [43]. In data triangulation, multiple data sources are used, or the data is collected at different occasions. In observer triangulation more than one observer is used in the study. In methodological triangulation different types of data collection methods are combined, such as qualitative and quantitative methods. In theory triangulation alternative theories or viewpoints are used.

Data collection in case studies tend to lean towards qualitative data that provides a richer and deeper view as compared to quantitative data [40]. As a data collection method, semi-structured interviews are common in case studies [40]. However, generalizing the results of a case study is often a subject of internal validity [26]. It is especially important in a case study to address reliability and validity of the findings. Even with good faith and intention, biased and selective accounts can emerge [39].

4.2 Qualitative research

Qualitative research attempts to answer to questions "Why?" and "How?" instead of "What?", "Where?" and "When?" as compared to quantitative research. The most common way to collect qualitative data is via an interview. Seaman defines interviews as a way to collect historical data, opinions and impressions about something [42]. The interview can be either structured, unstructured or semi-structured. In a structured interview the interviewer asks all of the questions, and the objectives are very specified. The focus of a structured interview is to find relations between constructs, and the objective is descriptive and explanatory [40]. In an unstructured interview the topic is broadly defined, and questions are asked also by the interviewee. In an unstructured interview open-ended questions are typical, and unforeseen types of information can be gained. The focus of an unstructured interview is on how the interviewee qualitatively experience the phenomenon, and the objective is exploratory [40]. A semi-structure interview is a mixture of both open-ended and specific questions. Semi-structured interviews focus on how individuals qualitatively and quantitatively experience a phenomenon, and the objective is both descriptive and explanatory [40].

Seaman has explored qualitative research in software engineering context, and he defines different ways to analyse qualitative data [42]. A typical way of analyzing is coding, which means extracting values for quantitative variables from qualitative data. Coding can also be used to organize the data based on themes and views. Seaman states that qualitative data is a lot richer than quantitative data, and using qualitative methods increases amount of information, diversity of information and confidence in results.

The basic objective of data analysis is to form conclusions and theories from the data based on a chain of evidence. Qualitative data analysis can be divided into two different parts, hypothesis generating techniques and hypothesis confirmation techniques [40]. Both of these techniques can be exploratory and explanatory case studies [40]. In hypothesis generating techniques, a statement or proposition is extracted from the field notes that is supported in multiple ways by the data, and then used as a hypotheses [42]. The data can also be divided into cases using cross-case analysis, where two groups based on some attribute are examined for similarities and differences. In the confirmation of theory, weight of evidence is built using

triangulation of different data sources.

Robson introduces common features of qualitative data analysis in a sequential list [39]. The first step is to code the initial set of data. Then, comments and reflections (referred to as 'memos') are added. After this the data is processed, similar phrases, patterns, themes, relationships and differences between sub-groups are identified. These identified patterns are then used to focus the next data collection phase. Gradually a set of generalizations are formed, that cover the consistencies discerned in the data. These generalizations are then linked to a formalized body of knowledge in the form of constructs or theories.

In the sequential list, hypotheses are identified after the data has been coded. The hypotheses are then used to guide the following data collection process in an iterative approach. During the iterations, some generalizations can be formed, which eventually form a formalized body of knowledge as a final result.

Hypothesis generating techniques are for example "constant comparison method" and "cross-case analysis" [42]. In the constant-comparison method the field notes are coded, text pieces are grouped into patterns and propositions that are strongly supported by the data are made. The propositions are then validated against any new data. In the cross-case analysis method data is divided into cases, such as two groups based on the same attribute. Pairs of cases are then compared to determine validations and similarities.

The data analysis can be conducted with multiple different approaches, with each having a varying degree of formalism. Robson lists four such approaches, from least formal to most formal: Immersion approaches, Editing approaches, Template approaches and Quasi-statistical approaches [39]. Immersion approaches are the least formal, relying mostly on the interpretive skills of the researcher and having a low level of structure. Editing approaches include the usage of some codes which the researcher decides based on findings during the analysis. Template approaches, also known as template analysis, use priori codes to develop a coding template, organising the data based on themes. Quasi-statistical approach is the most formal approach, including quantitative calculations, such as word frequencies. Runeson and Höst state that editing approaches and template approaches are most suitable for software engineering case studies, as a clear chain of evidence is hard to obtain in immersion approaches, and word frequencies are hard to interpret [40].

4.2.1 Template analysis

Template analysis doesn't describe a clearly defined method, but rather refers to a related group of techniques used to thematically organise and analyse textual data [24]. Template analysis is a relatively structured process to analyse textual data, and has mostly been used to analyse data from

individual interviews [23]. The key component in template analysis is the generation of a coding template, which is initially formed on the basis of a subset of the data and then applied to further data.

King defines code as a label attached to a section of text to index it as relating to a theme or issue in the data [24]. As an example, codes can be defined to identify points in the text where an interviewee mentions particular categories of presenting problems. These kind of codes can be descriptive, and require no further analysis by the researcher. However, also codes requiring more interpretation can be defined. In template analysis, a key feature is to form a hierarchical organization of codes [24]. This way groups of similar codes are clustered together to form higher-order codes. This way the researcher can analyse the data at different levels of specificity, with higher-order codes giving a broad view and lower-order codes a more specific view. Parallel coding, where the same text segment is categorised with two different codes, is also allowed.

The process of template analysis is to first create an initial template based on some pre-defined codes [24]. Then the template is revised by working through the data systematically while identifying and coding sections of interest in the text. Based on these findings, the initial template is modified, and finally develops to its final form. Modifications can be either insertion, deletion, changing scope or changing higher-order classification [24]. The final template is such that no sections of text related to research questions remains uncoded.

The template is then used in interpretation to provide an account. King lists various guidelines that can be used in the interpretation of coded data [24]. The first guideline is to compile a list of all codes with some indication of frequency. The distribution of codes can then help to draw attention to the aspects of data with either multiple codes or missing codes. For example, if one interview in a set of interviews is missing a certain theme, analysis for possible reasons behind that missing code can be made. However, King suggests that while patterns can be used to draw attention to certain parts of the textual data, frequencies alone cannot be to gain any meaningful information.

Another guideline is selectivity. King states that every code cannot be interpreted to equal degree of depth, and themes most relevant to understanding the phenomena under study must be identified. Prior assumptions of the researched shouldn't limit the analysis. Yet another guideline is openness, which King describes as taking themes judged as marginal into account. Themes that aren't of direct relevance to the initial research questions shouldn't be disregarded, as they can play a useful role in adding to the background of the study. Also themes lying outside the scope of the study should be included in the analysis, if they're considered to cast light on the interpretation of central themes in the study. Finally, King explains that purely linear relationships between codes, such as hierarchical relation-

ships with subsidiary codes next to parent codes, "may not reflect the kinds of relationships a researcher may want to depict in his or her analysis". Instead, maps, matrices and other diagrams can be used to explore and display findings, and analysis doesn't have to stop when a full linear template is produced.

Finally, an account of the interpreted data has to be presented. As with interpretation, King provides three common approaches to presentation [24]. First of the approaches is a set of individual case studies, followed by discussion of differences and similarities between cases. This approach gives a lot of perspective of individuals, but can be confusing if the amount of participants is large. The second approach is an account structured around the main themes identified, drawing illustrative examples from each transcript as required. King states that this is the approach that most readily provides a clear thematic discussion, but the dangers is in drifting towards generalizations and losing sights of individual experiences. The third approach is a thematic presentation of the findings, using a different individual case-study to illustrate each of the main themes. The main problem here lies in selecting the cases which represent the themes in the data as a whole. King also adds that no matter which approach is used, direct quotes from participants is essential.

5 Research design

Research design depicts the strategic decisions that guide how research is carried out [10]. This includes the studied phenomena, method selection, data-gathering and analysis. Different perspectives towards the researched subjects are usually categorized into two paradigms: positivist and constructivist [18]. Positivism contents that an objective reality, which can be studied, captured and understood, exists. Postpositivism argues that reality can never be fully apprehended, but only approximated. Postpositivism therefore relies on multiple methods to capture as much of reality as possible. Constructivism assumes that multiple realities exist. A critical paradigm lies somewhere between positivism and constructivism, assuming that our ability to know of this reality is imperfect. To understand as much of reality as possible, it must be subject to wide critical examination.

This thesis is based on the critical realist perspective. Critical realism studies people's perceptions to gain understanding of the reality that exists beyond these perceptions [19]. Critical realism is especially suitable for case study research if the process involves thoughtful in depth research with the object of understanding why things are as they are [14].

During this research, the author was an employee at Steeri. Consequently, this case study is based on the participant-observer research method [45]. Data triangulation is ensured by combining the observer results and

interview results.

The logical starting point of this thesis was the need by the case company to find a way to validate whether a feature is succesful or not.

5.1 Objective

The purpose of this study is to seek ways to improve the product development process of two products within the case company, using practices known as continuous delivery and continuous experimentation. Existing documented implementations of continuous experimentation are primarily executed in the B2C domain, often with a SaaS product. Examples are the Microsoft EXP platform [4] and Etsy []. The focus of this study is in the B2B domain, with two applications that are not used as SaaS.

Research objective To analyze the software development process transition towards shorter feedback cycle and data-driven decisions from the perspective of the case company, with respect to continuous delivery and continuous experimentation practices.

The study is an exploratory deductive case study, which aims to explore how continuous delivery and continuous experimentation can be implemented in the case company. The study specifically aims to identify the main requirements, problems and key success factors with regards to these approaches in the context of the case company. Integrating these approaches to the development process requires a deep analysis of the current development process, seeking the current problems and strengths. Adopting both continuous delivery and continuous experimentaton also requires understanding the requirements of continuous delivery and continuous experimentation. To narrow the scope of the thesis, the focus of this thesis is in the development process of two different software products. The two different products were chosen so that cross-case analysis and therefore more generalizations can be made.

The research process starts by reviewing literature on both continuous delivery and continuous experimentation. In the literature review, main goals are to identify existing requirements and success factors in these approaches. Then interview is used to analyze and identify the pain points the company currently has in its development process. Then continuous delivery and continuous experimentation processes are viewed from the viewpoint of the case company, comparing the previously found benefits to the pain points in the case company's development process. As a result of viewing the practices from case company's viewpoint, necessary restrictions and requirements encountered in the B2B domain are obtained. Finally, the continuous experimentation infrastructure [16] is applied to the case company, and deviations listed.

The research questions and research methods are summarize in Table 2.

| Knowledge gap | Research question | The focus of analysis | Data |
|--|---|--|---|
| B2B specific challenges in the practices | RQ 1: What are the B2B specific challenges of continuous delivery and continuous experimentation? | Empirical case supported by conceptual analysis | 13 interviews, documentations of the development process, literature |
| Benefits of the practices after reviewing the challenges | RQ 2: How do continuous delivery and continuous experimentation benefit the case company? | Empirical case supported by conceptual analysis | 13 interviews, literature |
| What deviations are there to existing models | RQ 3: How can continuous experimentation be systematically organised in the B2B domain | Conceptual analysis of continuous experimentation infrastructure by Fagerholm et al. [16], and application of the model to case company. | 13 interviews, documentations of the development process, author observations |

Table 2: Research questions and research methods.

In this thesis along with confirming existing theories we do aim to generate new concepts and theories. The theory and case analysis are continuously matched to gain insights coherent with both theory and empirical observations.

5.2 Case descriptions

Experimental context needs the three elements: background information, discussion of research hypotheses, and information about related research [26]. The two former will be discussed here, and the latter is introduced in the second and third chapters.

The company in question is Steeri Oy, which is a medium-sized company specializing in managing, analyzing and improving the usage of customer data. In this research, the units under the study are two teams responsible for developing two different software products. The first team is CDM, which is focusing on the respectively named software application. The other team

is Dialog, also working for a respectively named software application. Both of the teams are managed by a similar structure, consisting a team leader and a product owner, quality assurance and commercialization experts.

The organisation of Steeri Oy is of a divisional type, with each business area forming independent teams based on the products and projects. The business areas of the whole company range from IT consulting to CRM implementations and development of own software products. To scope down the thesis, only teams doing development of own software products are in the scope of this thesis.

The unit of analysis is the development process of the two teams. The whole development process consists of the development framework used, but also of the interaction with customers, tools used in the current development process and the roles of individuals in the teams. The unit of analysis is studied by focusing on interviewing individual members at different positions in the organisation. The purpose of the interview is to identify the pain points in the current development process, which are then viewed against the benefits of continuous delivery and continuous experimentation found from existing literature. Analysis of implementing continuous delivery and continuous experimentation is then done based on the interview results.

The development process of the team is elaborated in detail in the chapter "State of the practice". In the following sections the general characteristics of company, the two products and applied development processes are identified. The general characteristics of the two teams under scope are then compared. The general characteristics are as follows.

- | | |
|---|---------------------------------|
| • Product description | • Programming languages |
| • Product environment | • Development tools |
| • Product architecture | • Version releasing |
| • Users of the software component | • Unit testing |
| • Team composition | • Acceptance testing |
| • Development practice | • Feedback collection |
| • Business model TODO: validate if this is confidential | • Usage data collection |
| | • Feedback and usage data usage |

The first three characteristics give us information on how the teams develop their software products. Testing is used to.. TODO: explain the general characteristics

5.2.1 Dialog

The Dialog subteam focuses on developing a multi-channel online marketing automation, iSteer Dialog. It is a software product designed to allow effective marketing on multiple channels, such as e-mail and websites, and to automate repetitive tasks. The product can be integrated to existing CRM solutions, and the data stored in CRM can then be effectively used for marketing purposes.

Product description. Dialog is a multi-channel online marketing automation tool. The purpose of the product is to allow effective marketing on channels such as e-mail and websites, and to automate repetitive tasks. The product can be integrated to existing CRM solutions, and the data stored in the CRM can then be effectively used for marketing purposes. The product is primarily used through a comprehensive user interface by marketing professionals of customer companies.

Product environment. The product can be run on either servers provided by the case company or on the customers internal servers. Most of the customers use the product on the case companys own servers. Each customer using a provided server has an individual server and configurations, due to different integrations to the customers own systems. Customers running the product on their own servers form a minority.

Product architecture. The product consists of a single maven-packaged project. Along with the actual product, several additional third-party or custom components are required. Such components are for example Typo3, which is an open source web content management framework, and a custom component used for SMS integration. Most of the components reside in the version control system, and the main product additionally resides in the CI system. Custom integrations to customers data sources, often based on a third party ESB tool, are individual components as well.

Users of the software component. Dialog is used by marketing professionals to automatize repetitive tasks. The user amount varies, but it is generally less than five human users per customer company.

Team composition. The team consists of 3 software developers. The software developers don't have designated roles, and each developer is involved in every aspect of the product development process.

Development practice. Development is based around a prioritized kanban board. Features are split into tasks by the product owner or team leader, and inserted into the kanban board. The developers then pick a single task from the sprint backlog column, complete it, and move it to

the review column. From there the task undergoes quality assurance and delivery to the customer environment. The detailed development practice is introduced later in this chapter. The team is in the state of continuous integration, but the quality of automated tests isn't yet at the required level to lay confidence in builds working.

Business model. The product is still in the development phase, and is not sold as a ready package. The productization process is however ongoing, and the ideal situation is that the product can be sold as a ready package with only minor configurations per customer. The product is currently sold as a custom project, and requires installment to the customer environment and integration to the customers systems. The customer has also be trained to use the product, and after the product goes live the cooperation is continued with a support contract.

Programming languages. The product is written mainly in Java, with the front-end built using JavaScript.

Development tools. The team uses GitHub as the version control system, and Jenkins as the integration server. Trello is used as the kanban board to manage projects, and Bugzilla is used as the bugtracker.

Version releasing. A new version is released when a customer requires a certain feature, and that feature has been completed. Updates for other customers are then released every two to three months on average.

Unit testing. Unit testing is done by the developers whenever a feature is written.

Acceptance testing. Acceptance testing is performed by the Concepting and Management team whenever a version is deployed to the customer environment.

Feedback collection. Feedback is collected from customers in retrospective meetings after the projects. Feedback is also received directly from customers via e-mail and meetings, but isn't collectively gathered in other ways. Feedback isn't automatically collected in any way.

Usage data collection. Usage data isn't collected at all.

Feedback and usage data usage. The gained feedback is discussed within the team to improve the development process or certain parts of the product. The feedback isn't systematically used,.

-java, javascript front? -frontend -binaries are built -jenkins deployed
-which env decides -testing (written by coder, tested by uat?) -management

The software is configured and integrated to the customer environment as project work. This deployment doesn't require additional code as per customer, but only different configuration files.

5.2.2 CDM

The CDM subteam focuses on building a Master Data Management [30] solution, which integrates multiple customers data sources such as CRM, ERP and billing system to create a single point of reference.

Product description. CDM is a Master Data Management [30] solution, which is used to integrate multiple customer data sources to create a simple point of reference. These data sources can include for example CRM, ERP and a billing system. The product also manages the data by removing duplicate records, matching same records and cleaning and validating the data with the help of external data such as the resident registration database. The product has an user interface that can be used for testing, but the product is primarily a background application without users.

Product environment. The product is installed to a customer environment, and certain custom specific configurations and rules have to be implemented for each customer. Such configurations are for example the integrations to customers systems and external identification services used.

Product architecture. The product consists of several maven-packaged projects, which in the build stage are added as a dependency under a single main project. The whole product can thus be packaged into a single web archive file.

Users of the software component. CDM is an integrated application, which is primarily only used by other applications. Therefore humans are only using CDM for debugging purposes.

Team composition. The team consists of 4 software developers and a team leader. The software developers don't have designated roles, and each developer is involved in every aspect of the product development process.

Development practice. Development is based around a prioritized kanban board. Features are split into tasks by the product owner or team leader, and inserted into the kanban board. The developers then pick a single task from the sprint backlog column, complete it, and move it to the review column. From there the task undergoes quality assurance and delivery to the customer environment. The detailed development

practice is introduced later in this chapter. The team is in the state of continuous integration, but the quality of automated tests isn't yet at the required level to lay confidence in builds working.

Business model. The nature of the product requires customer integrations to customers systems, which are currently mostly done as manual work. In the ideal situation the product contains automated integration, and could be sold as a ready product with minor configuration work. Currently the product is sold as a custom project, and requires installment to the customer environment and integration to the customers systems. After the product goes live the cooperation is continued with a support contract.

Programming languages. The product is written in Java.

Development tools. The team uses GitHub as the version control system, and Jenkins as the integration server. Trello is used as the kanban board to manage projects.

Version releasing. A new version is released when a customer requires a certain feature, and that feature has been completed. Updates for other customers are then released every two to three months on average.

Unit testing. Unit tests are required for new features, and are written by the developers during of after the feature implementation.

Acceptance testing. Acceptance testing is performed by the Concepting and Management team and the team leader whenever a version is deployed to the customer environment.

Feedback collection. Feedback is collected from customers in retrospective meetings after the projects. Feedback is also received directly from customers via e-mail and meetings, but isn't collectively gathered in other ways. Feedback isn't automatically collected in any way.

Usage data collection. Usage data isn't collected at all.

Feedback and usage data usage. The gained feedback is discussed within the team to improve the development process or certain parts of the product. The feedback isn't systematically used,.

5.3 Context

The previous chapter introduced the two teams and products under the scope of this thesis. The more general context of the whole company under study is defined by the following properties.

Company size. Steeri has approximately 80 employees at the writing of this thesis. The employees are located in two cities in Finland, and many of the employees are assigned to projects at other organizations.

Business area. The company in question focuses in managing, analysing and improving the usage of customer data. This includes various CRM-systems, such as Oracle Siebel, but the company also has its own marketing automation and master data management products. The company also does data analytics within the Business Intelligence team.

Teams under study. As Steeri covers a wide business area, the scope for the thesis is narrowed down to two teams. The other team focuses on product development of a marketing automation product and a master data management product, while the other team consists of managers and Q&A dealing with the same products.

Team setup. The development team consists of a team leader and developers, while the management team consists of two product owners, commercialization manager and a requirement engineer.

Development practices. Currently the team in question uses continuous integration, with feature branches merged into mainline as soon as the feature is finished. This could differentiate from a couple of hours up to multiple days.

Development model. The development model is essentially a Lean model, the core being a prioritized Kanban board.

Tools used. The teams under study use Git as a version control system and GitHub code review tool, Trello as the Kanban board for task management, Jenkins as the continuous integration server and Bugzilla for issue tracking. Various other tools, such as the hour logging system, CRM system and various other ticket logging tools exist as well.

5.3.1 Development process

The development process of the case company has evolved from traditional, waterfall-style development through an agile phase to the current Lean approach based on a kanban board. In the very first phases, when the company was young, tasks were simply put into excel. Eventually, Agilefant backlog product management tool was introduced and tasks were moved into it. Scrum and agile practices were also adapted at the same time. Agilefant allowed sprint management, but caused too much overhead according to the developers. However, during the time Agilefant was used as the product

and sprint backlogs, a part of Dialog development was offshored. Agilefant allowed easy tracking and managing the offshore development efforts.

After Agilefant was determined to cause too much overhead, Scrum whiteboard was used for a while. The physical whiteboard however caused issues with offshore participants. After the Development and Integration team was split further into two sub-teams focusing on each project, Scrum meetings were also seen as overhead. This was caused by only half of the team being interested and attached to tasks considering a single product. After it was decided that the team is split into two sub-teams, practicing Scrum was again seen as an overhead due to the low amount of participants. Eventually the whiteboard was replaced with an online kanban board, Trello. By the same time the Scrum meetings were abandoned for a more streamlined process, that is able to respond to changing requirements faster.

Work items are first added to customer backlogs in iSteer Contact. In this phase they might be just initial skeletons and do not need to contain all needed information. The work items are added by product or business owners or project managers. At this phase the created backlog items are not yet visible in team backlogs in iSteer Contact. The idea is to list them in the customer backlog as early as possible and start refining the requirements collaboratively both offline and online with the help of tools such as Chatter. Chatter can be used to discuss a single story or the complete backlog.

When the backlog item is ready for the development team to start working on it, it should contain at least the following information:

- A descriptive name
- Specifications that explain both the technical side and the business side
- Agilefant link four hour reporting

Team backlog is a view to all stories assigned to one selected team. It is a tool to collect and organize stories from all customer/project backlogs without cloning details into multiple places. It is a view of all “in team backlog” stories that are not yet completed and are assigned to the selected team. Team backlog items are prioritized and estimated once a week, and new stories will be moved to Trello when accepted by the development team. Stories can be moved to Trello by project managers or product owners, but they must be checked by team leaders who then move them forward into the sprint backlog or prioritized list columns.

Overview. Trello is a project management application that uses Kanban to control the production chain from development to release. Kanban attempts to limit the work currently in progress by establishing an upper limit of tasks in the backlog, thus avoiding overloading of the

team. Trello consists of multiple boards, each representing a project or a development team. A board consists of a list of columns, and each column consists of cards. Columns each contain a list of tasks, and cards progress from one column to the next when each task has been completed. A card is essentially a task, which is added by the Backlog Owner, and can be checked out by a developer. In Steeri, the columns used are Sprint Backlog, In Progress, Review, Ready, Verified and Done.

Sprint backlog. The Backlog Owner moves the cards that have the highest priority into the sprint backlog. The sprint backlog should always contain enough cards so that whenever a developer completes a task something new is available in the backlog.

In progress. The actual development work is done in this stage. The actor here is the developer. The checklist to move a card into review stage is:

- All tasks are complete
- Changes are pushed to a feature-specific branch
- Unit tests are written and pass in the CI server
- Feature is well-written and does not need refactoring
- Pull request is created and a link is added to the comment field
- Feature has been documented as needed

If the feature needs refactoring a task list must be created and the card moved back to In Progress column.

Review. Here other developers review the new code and deploy it to a development environment. Checklist for moving the card into ready stage is:

- Pull request is reviewed by at least two (2) persons
- Pull request is merged to the development branch
- Feature is deployed to a development environment
- Source code quality has to be good enough!

After a developer has reviewed the pull request, he leaves himself as an assignee. The second person who reviews the pull request is responsible for cleaning all the assignees and moving the feature to Ready column. If there is a major problem in the pull request the feature should be moved back to Sprint Backlog and the yellow “Boomerang” tag added. Person who created the pull request is responsible for implementing the necessary remarks. If only small fixes are needed, they should be

implemented within the Github pull request workflow. The second person who has reviewed and accepted the pull request is responsible for deploying the feature to a development environment.

Ready. Here the product owner verifies the new functionality in the development environment. The card can be moved to Verified if the feature has been verified by the Product Owner in the development environment. Product owner is responsible for moving the feature to the Verified column. If verification for the feature fails, the Product Owner should move the feature back to Sprint Backlog column with the highest priority. In addition the Product Owner should add a yellow “Boomerang” label with a comment describing the results in the feature.

Verified. Here the backlog owner collects the timestamps and trello flow data. The timestamps depicts the duration it took from a card to process through the whole chain. The data is then used to analyze which columns the card spent the longest time in, and to identify the pain spots.

Done. This column simply states that the task has been completed, and should eventually be archived. There’s currently no general validation required from the customer, as the customer projects each have a different schedule and process for builds.

Prioritized lists. The backlog owner adds tasks to prioritized lists from team backlog as soon as the tasks meet the required criterias, contain the required information and are inspected by both the stakeholders and the backlog owner.

5.4 Methods

The primary source of information in this research are semi-structured interviews [40] performed within the CDM and Dialog teams. The interview consists of pre-defined themes as follows: (1) current development process, (2) current deployment process, (3) current interaction with customers, (4) problems and strengths in the current development process, (5) software product, (6) future ways with continuous delivery and continuous experimentation. Data triangulation is implemented by interviewing multiple individuals. Methodological triangulation is implemented by collecting documentary data regarding the development process from internal documents.

The interview was chosen as a data collection method because of the nature of the research questions. Since the study focuses on applying two development methods to the development process of a team, individual perceptions of members are required. As the case study doesn’t include a

technical implementation, quantitative measurements to properly measure the effects before and after implementation isn't an option. The research questions cannot be properly supported with quantitative data. There is also a uncertainty about how much information the interviewees are able to provide, and thus the questions are mostly open-ended. The nature of interviewees opinions on the research questions are also not known in advance, and quantifying it isn't simple.

In order to answer the research questions, information regarding the development process, customer interaction and feedback, deployment process and technical product details are required. The interview is divided into themes to address these aspects. The interview questions address, among other things, the specific situations and action sequences of the interview rather than general opinions.

5.4.1 Data collection

The interview is a semi-structured interview with a standardized set of open-ended questions, which allows deep exploration of studied objects [40]. The interviews are performed once with every subject. Leading questions are avoided on purpose, and different probing techniques such as "What?"-questions are used. The interviews are performed in the native language of the interviewee if possible, otherwise in English, and are recorded in audio format. The audio files are then transcribed into text.

The interview begins with a set of background questions, used in coding the interview data per subject. After the introductory questions, the main interview questions make up most of the interview. The interview session is structured based on areas under study rather than based on a specific model. The interview roughly consists of four different areas: development process, deployment process, feedback and usage data, and experimentation.

The development process questions aim to identify the current ways of working, including variances between different roles. The questions also identify the perceived problems and strengths in the development process, which are then used in analysing the potential process for continuous experimentation.

The deployment process is examined from multiple point of views. This includes the customer point and view, developer point of view and management point of view. Emphasis is based on B2B specific issues and practices, which are often caused by differences in the software environment as compared to B2C.

Feedback and usage data collection section analyse the procedure of collecting and utilising the data. The section aims to identify the manners in which feedback is received, collected and systematically used. The section also analyses how different roles interact with the customer.

The experimentation section analysis the possibilities of conducting ex-

periments to guide the development process in the case company. This includes technical details, such as instrumenting the software components with data collection mechanisms and external usage data collection components, and storing the data in customer environment. The section also includes more specific questions, such as possible challenges faced in A/B testing and issues in quickly deploying experiments to the customer environment.

Interview data is primarily sought from the developers of the development teams and the managers of the team. As the focus of the thesis is on the development process of two different teams, all participants of the team and its management team are interviewed. Process data is sought from the process documents made by the team. Quantitative data, such as the Trello ticket flow time, is sought from Trello. Data regarding the development process is also collected from the internal documents.

During the interviews, it was soon detected that the interviewees are able to give accurate answers to certain contexts based on their positions. This eventually contributed to the template analysis, where the data was better organised based on roles. Developers were most aware of the product and project context, including technical details and project details. Team leaders were mostly aware of the process and project context. They were the persons who were most able to explain the development process in detail. The team leaders were also most proficient in feedback collection. The management was most aware of the company context, including the status of the organization and available resources in the company.

5.4.2 Data analysis

In this case study the data analysis is based on template analysis, which is a way of thematically analysing qualitative data [22]. The initial template was first formed by exploring the qualitative data for three themes: development process, deployment of software and experimentation. Through multiple iterations of the data, multiple subthemes were then added to the three existing themes by further coding the data. Attention was paid to different roles of the interviewees.

Prior to the template analysis the data was inserted into a spreadsheet, where the rows represent codes of interest and columns represent interview subjects. This eased the analysis of data. The final template was formed after three iterations, which each increased the depth of the template by one or two levels, and the amount of codes by several.

Template analysis was chosen as an analysis approach because it makes it easy to organise and analyse the data according to themes. It is also particularly effective when comparing the perspectives of different participant groups [24]. Since the thesis focuses on multiple themes and the interview on multiple groups, template analysis was seen as a suitable approach.

The results are formed from the final template, and structured according

to themes of the research questions to answer the research questions as comprehensive as possible. The results are then validated from the original interview data, so that none of the sections left out of the template are in conflict with the findings.

6 Findings

This section summarises key findings of the case study. It is structured according to the research questions, allowing the progress from a higher-level perspective to more detailed view. First, the identified challenges in continuous delivery and continuous experimentation are listed and explained. Then, the benefits of applying continuous delivery and continuous experimentation, taking the challenges into consideration, are analyzed. Finally a continuous experimentation infrastructure by Fagerholm et al. [16] is mapped to the case company context, and deviations identified.

6.1 Continuous delivery: B2B challenges

Software development practices and product characteristics vary based on the domain and delivery model. Typical B2C applications are hosted as SaaS applications, and accessed by users via a web browser. In the B2B domain, applications installed to customer environments are very common.

The challenges regarding continuous delivery will be analyzed in three areas: technical, procedural and customer.

Technical aspect includes the environmental challenges, configural challenges and other challenges related to the software product and its usage. Procedural aspect includes the challenges regarding the development process of software. Customer aspect consists of the customer interaction,

6.1.1 Technical challenges

The technical challenges for continuous delivery are derived from the interview. A part of the interview focused on the current deployment process and customer interaction of the case company, and this process is then compared to the continuous delivery practice from literature review.

| Problem area | Specific problem |
|---------------------|--------------------------------------|
| General | Downtime |
| General | Automate testing |
| Configuration | Third party application integration |
| Configuration | External components |
| Configuration | Different configuration per customer |
| Connection | Firewall configurations |
| Connection | Required certifications |
| Connection | Missing user rights |
| Connection | VPN user limit |

One of the main differences between B2B and B2C is the production environment. While a company focusing on B2C often owns the product environment, in B2B domain it is common to install the software on customers environment. This introduces multiple challenges: accessing the environment, integrating the software to third party components and downtime, to name a few.

The largest technical task is to set up automated testing and test environments. While this task can, depending on the project and current state of tests, take up most of the resources in transition to continuous delivery, in B2B context some of the testing costs can be pushed partially to customer projects.

Both of the case company's software products are integrated to various 3rd party applications and APIs. As CDM is a component that integrates various data source together, changes to the APIs communicating with these applications must be planned and discussed in advance. Based on the interview results, automatically updating the integrations requires an unduly amount of work seeing the results.

It is common for B2B applications to have external components that have to be configured when the software is installed or the APIs to these components changed. Therefore the configurations for these external components either have to be manually updated or automated as well. In the case company, both Dialog and CDM are integrated to multiple external components. Currently upon version change these integrations are manually updated when necessary. In continuous delivery, changes that modify the API and therefore break these integrations might have to be manually updated.

Both of the case company's products are used in multiple different customer environments. This introduces a problem of managing different configurations per customer environment and software instance. Multiple possible solutions exist for the configurations, however, such as creating custom profiles per customer with different configuration for each profile.

Transferring the software product to a customer-owned environment instead of own environment introduces multiple possible challenges. These challenges include firewall configurations, required certifications and missing user rights. The customers also have to be informed on actions performed in their servers. For certain environments, a VPN connections is required. A VPN connection often has a user limit, which can prevent automatic version deployment.

6.1.2 Procedural challenges

| Problem area | Specific problem |
|---------------|--------------------------------|
| Environmental | UAT |
| Procedural | Feature stubs |
| Procedural | Multiple customer environments |
| Procedural | Versioning |
| Roles | Developer responsibility |

The basic deployment pipeline in the case company first includes a deploy to UAT server, which is then tested manually by either the team or the customer. Only after the UAT version has been validated to work properly, can the production version be released. Based on the interview results, continuous deployment to production is seen very risky due to the applications playing a major role in running the customers business, and continuous delivery with deployments to a staging environment is seen as very beneficial.

CDM is a software component that is compiled from multiple components. Therefore the event triggering a deployment shouldn't be the merge of a pull request, since pull requests are project specific. Some features require changes to multiple projects, and all of the pull requests have to be merged for the feature to work as intended.

Both of the case company's products are used in multiple different customer environments. This introduces a procedural challenges: whenever a feature is ready to be delivered, it can either be deployed to a single customer or to every customer. Attention has to be paid when configuring the workflow to allow deploying to an arbitrary customer. If only certain customers accept updates on a continuous basis, it has to be possible to prohibit deploying to every customer at once.

Multiple customer environments affect versioning of the software product. In the case company, each customer has a unique configuration of the product, with possibly different versions of components. However, according to Jan Bosch, in an IES environment only a single version exists: the currently deployed one. Other versions are retired and play no more [7].

Another procedural challenge is that continuous delivery drifts response towards the developer, and the developers decide what is ready to be re-

leased. If the developer can single-handedly deploy a feature, the management can quickly loose track on the features available to customers. This also requires the developers to deeply understand the details of the version control system and automated testing.

6.1.3 Customer challenges

| Problem area | Specific problem |
|--------------|---|
| Training | Customers have been trained to use a certain version of the product |
| New release | Customer reluctancy towards new versions |
| New release | Deployment time and down-time |
| New release | A customer might be doing a critical task |

Customer reluctance towards new features is common, since new releases occasionally contain new bugs. In the case company, customers using Dialog have been trained to perform certain tasks with a certain user interface. The customer might perform these tasks daily, once every two weeks or even less frequently. Thus if the UI changes often, the customers initially take more time to perform the tasks.

In the case company many of the customers have a reluctancy towards new versions, caused by the versions occasionally containing new bugs and the customers having just learned to perform certain tasks with the last version. Customers need to understand what continuous delivery is all about. While customers initially are against constant modifications, if modifications were made more often, it would lead to smaller modifications and it might reduce the customer reluctancy towards new versions.

Customers might be using the software when a new version is deployed, and the deployment process shouldn't interfere with ongoing usage. Therefore the most viable solution is to initially deploy the versions to the UAT environment. To prevent downtime and obstructing the usage of the software, the production deployment date should be negotiated with the customer. Production deployment might cause downtime, and cannot be made during times when critical tasks are performed. The test environments can be continuously updated, and once the customer is satisfied with the software and feels a need for a new version, deployment to production can be triggered.

In the case company, the production deployment date is often negotiated with and specified by the customer, while the product can be deployed to UAT at will. UAT is also used by the customer to validate the functionality

of the software, since the customer has ordered the product to perform certain tasks and fill certain criterias. However, the customers don't validate the versions at the UAT environment unless they are informed that new versions are available, since the UAT is not used for any other purpose than testing.

6.2 Continuous experimentation: B2B challenges

The case company's two software products differ on how they are used by the customers. In this section, continuous experimentation challenges are analyzed for general cases encountered in both products, and special cases encountered in a single product due to a specific product characteristic.

The challenges regarding continuous experimentation will be analyzed in three areas: technical, customer and A/B-testing. A/B testing was chosen as an examined subject of its own, since it is a commonly known experiment, and the difference between the case company's products mostly affects A/B testing.

Technical aspect includes the data collection, measurable metrics and other challenges related to the software product and its usage. The customer aspect includes the challenges regarding the experiments from customers point of view. A/B testing section analyses the viability of guiding the development process of both products with A/B tests.

6.2.1 Technical challenges

| Problem area | Specific problem |
|-----------------|---|
| CDM | No user interface |
| CDM | Only used via API calls |
| General | Measurable metrics |
| Data collection | Data collection in customer environment |
| Data collection | Required permissions |

The case company's products represent two totally different styles of software products. While Dialog is used by users via a UI, CDM is used by other systems via a REST API. As CDM is a software component which is primarily used by other systems via API's and certain calls, tracking user behavior doesn't provide much input. CDM also isn't used by an user interface, which limits the possibilities for A/B testing.

The difference between the products limits the possible scope of experiments. While it is possible to experiment with a software component that has a UI and human users by employing A/B testing, other experiments have to be conducted on software without direct user access. Such experiments include for example feature alpha and usage metrics.

The business model for the products limits the wanted improved metrics. Since the purpose of the software components isn't to generate revenue per se, generated revenue for example is not a valid measurable metric. The products also have custom solutions for different customers, which means a measurable property might only apply to a single customer. This makes it harder to pick the correct measurement variables. This is what a developer has to say:

"We don't aim for increase in clicks or in purchases. Instead, if we wanted to measure how Dialog could be improved, the most useful metric would be the customer opinion. A measure could also be how fast certain work tasks can be performed, but it's hard to measure it. The customer doesn't necessarily perform the task as fast as possible."

A possibility is to tune the product in cooperation with the customer, using their internal measures as metrics. According to a Dialog developer:

"If the customers had their own metrics that they have been tracking, and we could affect this trend, we could figure out whether this is caused by changes in the software. Many customers monitor their employers and measure them."

A CDM developer also has an idea of surveying the customers daily routine as a basis for building experimental features:

"We should always think what the customers status is now, what kind of features the customer is using and what is the system that CDM replaces. Against these things we can develop the experimental features. In the old system the customer might have done things that they felt was inconvenient. We could conduct the experiments customer wise by surveying the customers daily routine and build experimental features based on this."

While it is common for B2C applications, too, to fluctuate between different Overall Estimation Criterias (OEC), in the B2B the possibility to co-operate with the customer should be considered. Eventually, the OEC has to be well thought. For example pure profit is not a good OEC, since raising price can increase a short-term profit but have a negative impact in the long run [27].

The case company's products are mostly installed to customer environments, with customer-owned servers. The data then has to be effectively stored and acquired from an external feature, if the analysis is to be done at the company's own servers. However, the case company's applications deal with millions of records of data, and the developers don't see storing additional experimentation-related data as a challenge:

"The data would be on the customer companies servers, where we have our own access and right to install software. If the data is collected on the server where we fetch it from, it's not a challenge. The data shouldn't encumber the system."

"The customers are motivated that we are trying to improve their actions. In B2B customers probably let us track their actions, because they have already bought the product."

Data collection directly from a customer might also require permissions and negotiations with the customer.

"We need a permission from the customer to run experiments. The customers conduct mostly routine actions. The software has to be stable, so that the customers know where everything is."

6.2.2 Customer challenges

| Problem area | Specific problem |
|--------------------|---|
| Origin of features | Customers have actual needs |
| Downtime | Customers don't want down-time |
| Downtime | Users perform routine tasks and have been trained for them, therefore the UI shouldn't change often |
| General | Customers have to be informed for major changes |
| Data collection | Environment |
| Data collection | Feedback collection |
| Data collection | Proper agreements need to be made |

In the case company, features have been primarily created for actual customer needs. This somewhat limits the direction of product development. In the case company, almost 100% of the completed features have been determined by the customer. This is what a product owner has to say:

"The goal ratio would be 50-50. Half would be our vision of how to perform business, the rest requirements from customers. In the current situation the ratio is 100-0, and for over a year every feature has been the result of a customers needs regarding Dialog."

The product roadmap has advanced by being led by the customer projects, since the revenue comes from the customers. The actual customer requirements are a large part of the product development stage, and customer projects often consists of only features required by the customer. As the product development is in the ideal case guided by both the customer needs and own ideas, the proportion of data-driven decisions is smaller than in B2C.

Another issue in the customer projects is downtime. Since the customers are using the software during normal work hours, critical tasks cannot be interrupted by version changes. Currently the time slow for version deployment is negotiated with the customer, and the downtime is positioned during times that the software isn't used. A simple solution for this is to have the test environment continuous updated, while the production environment update is negotiated with the customer. However, the test environment isn't always actively used. This is what a developer has to say on the issue:

"Usually the customers test environment isn't in heavy use, and it might take for the customer to react to changes. In our own test server we get a lot faster feedback. The customer has also be informed if we perform major changes, since the customer needs to know how to use the software."

Therefore continuously updating the test environment either has to alarm customers that new features are to be tested, or production environments has to be fixed in such manner that continuous updates don't cause downtime.

The customers have been trained to perform certain tasks with the software, and if the UI constantly undergoes major changes the customers might get lost. There is already some reluctancy towards new versions, since new versions possibly introduce new bugs and certain UI elements might have their places changed. Therefore customers have to be informed for major changes, so that possible new training sessions can be arranged. For minor updates, release notes and required guides are a must.

The case company's products are mostly installed to customer environments, with customer-owned servers. The typical deployment flow includes a customer testing phase in the UAT environment before a production version can be deployed. The UAT environment is not actively used, and the customers perform UAT only when they are informed that a new version is available. The production environment is seen critical enough that no major bugs are allowed, and thus the UAT testing cycle is seen favorable by both the customer and the company.

The customer feedback is currently continuously received but not actively collected. Product owners are in constant connection with the customers, and customers can send comments and requests via them. The

customers can also insert tickets into a bug tracker. Most of the received feedback is about some functionality not working or a bug. Weekly status meetings with the customer are also arranged. Feedback is received via e-mail as well. In the end of the project, a retrospective meet is held with the customer, where the whole project is discussed.

The feedback collection in B2B has multiple challenges. The feedback doesn't come straight from the end user. Feedback relies on what the other companies representatives say, and this often includes their own opinions. Companies want the software for their own use, and they have requirements for property rights, security and others. If there are customer projects with multiple suppliers, it's hard to get feedback since maintaining communication and keeping track on up-to-date status of the project with multiple suppliers is hard.

In case of Dialog, the user is not the payer. The user has to use the software even if it was bad. It is hard to get a neutral point of view.

According to the interview results, the product owners state that to collect data from customers, proper agreements might need to be made. Some of the customer companies have strict rules for which data can be collected.

6.2.3 A/B testing

| Problem area | Specific problem |
|----------------|--|
| User amount | Low user volume |
| CDM | No users, only integrated systems |
| CDM | No user interface to tune |
| Environment | UAT is the norm |
| General | Measurable metrics |
| User education | Users have been educated to use certain features |

Low user amount in both products, with Dialog ranging from 3-5 users per customer and CDM having zero human users, complicates statistical analysis. Kohavi et al. [29] state that a sufficient user amount is a requirement for controlled experiments on the web: "Very small sites or products with no customer base cannot use experimentation, but such sites typically have a key idea to implement and they need quick feedback after going live. Because new sites are aiming for big improvements, the number of users needed to detect the desired effects can be relatively small (e.g., thousands of users). Large sites, which are typically better optimized, benefit even from small improvements and therefore need many customers to increase their experiment's sensitivity level". The lower a sample size, harder it is to draw meaningful conclusions from it. For a company operating in busi-

ness with low user amounts, due to either new market or the type of the product, attention has to be paid on when statistical significance is reached. The company either has to work with a lower confidence interval, or majorly increase the conversion rate to decrease the amount of subjects needed for statistical importance.

However, a product owner views the user amount from a different point of view:

"The amount of statistical data isn't much (due to the user amount), but the users actions can be tracked, and they continuously perform repetitive tasks. This customer would then provide a lot of usage data. If the time spent by this user could be halved, could the customer make more money. This doesn't therefore necessarily require a large user amount."

The main difference between the case company's software products is the intended use of the software. CDM, being an integrated system, doesn't have a UI unlike Dialog. The types of A/B tests that are viable to CDM are mostly related to performance, which is seen as the only viable metric by the developers.

Along with low user amount, the issue with Dialog is that the customers have been taught to use the UI to perform certain tasks. The UI of Dialog is quite complicated, and some customers are very reluctant to changes to the UI. A possible solution for this issue is a detailed use guide with each new version, where the changes are easily readable.

As stated earlier, new versions undergo testing in the UAT environment before a production release is made. This makes it possible to either perform A/B testing in the UAT environment, with a possibly very low amount of users, or instead only validate that the A/B testing does not cause complications and can be moved to production. Another possibility is to skip the UAT, which is however seen as very risky by the product owners and developers.

In A/B testing, the measurable metrics greatly vary between the case company's products. In CDM, the team finds A/B testing the query machine beneath the API to be the most beneficial. These kind of A/B tests would divide the users to different queries, and the performance metric could be optimized based on the result. Another possibility found from the interview results is offering multiple query features, and then tracking which one is used more by the customer. However, as the customers clients are hard-coded, the users would also have to modify their systems accordingly.

In Dialog, the need for data-driven decisions is greater than that of CDM due to the visible parts of the software. Even though Dialog doesn't generate profit per-se, optimizing the time it takes for customers to perform certain tasks increases the benefit of Dialog for the customers. Therefore the A/B testing in Dialog could include metrics such as time spent on performing

certain tasks, amount of errors encountered in certain features or amount of back-selections indicating mistakes made by the user.

"If we had more customers, we could for example do A/B testing. Some customers are really enthusiastic of Dialog, especially that they could have an impact on the development. If we sold them a model where they could try things out, they could really be willing to do that. On the other hand we have customers that want no changes at all."

Some features in Dialog and CDM are specified to work exactly as is, and the customers have been educated to use them in training session. This partially limits the pool of possible A/B cases, since the B way might have to be taught as well, if the change is major.

6.3 Continuous deliverys benefit to case company

A major problem found in the interviews was that currently the reliability in new versions is low. The low reliability both increases customers reluctance towards version updates, and increases the amount of user access testing that is performed after version release. This makes the duration of building a new version relatively long and unstable, since the bug amount differs upon each version deployment. The version releases are also often performed in a hurry, which also increases the error rate.

The manual build process involves human error factor, which can cause multiple bugs in new versions. One example of such an error is using the wrong commands to compile and package the project, which causes wrong binaries getting under the new versions. An important factor is also partially missing documents. Some steps of the deployment pipeline are documented, but there are still parts that some developers memorize. According to a developer:

"There are many components that have to be configured manually. The documents for how to execute these manual actions exist, but many things still rely on memory."

Since some of the manual steps are only memorized by certain developers, the deployment requires assistance from these persons. With continuous delivery only a handful of developers might have knowledge of the whole configuration for build deployment, but the rest can still trigger the deployment process.

The customers are unaware of the problems related to deployment issues unless the deployment has errors. However, if the version contains many UI changes the customer might feel frustrated by being unable to perform certain actions in the old way:

"If something changes, communication with the customer is hard. The customers feel that even if the version included very good changes, if something visible changes it gets a lot of negative feedback".

"Customers have been doing this for their living - they're getting old, they have their own working histories and they're in the working mode. These persons don't approve alterations in the software, and they're not ready for the constant changes."

Continuous delivery process solves many of the stated problems. By adopting continuous delivery, the changes can be kept small and frequent so that the customer doesn't have to train the staff with sparse, big updates. Instead, the changes can be introduced more frequently with smaller updates. Continuous delivery allows the release of changes in time spans of hours to days. Therefore the most important features and fixes can be delivered faster. With faster problem fixing due to increases in bug discovery and fixing rate, the software ends up with fewer bugs.

Continuous delivery allows the improvement of speed and quality. Releases can be quickly released, or put on hold for later release by feature switches [20]. Feature switches allow developers to increase the quality of certain features until they're ready for release.

The amount of manual work can be diminished to couple of button clicks defining which customer profile to deploy.

6.4 Continuous experimentations benefit to case company

Applying continuous experimentation reduces the amount of guesswork and instead provides a systematic approach towards product development. In the case company, ideas for new features are often discussed in the steering group. However, Avinash Kaushik states in his article Experimentation and Testing primer that "80% of the time you/we are wrong about what a customer wants" [1]. The company Netflix has also found similar results, considering 90% of what they try to be wrong [32]. A way to tackle this issue is to adopt a process of continuous experimentation.

Instrumenting software products and then gathering and analysing the data to guide development is the main process of continuous experimentation. The general benefits of continuous experimentation lie in validating the features with customers. Generally speaking, the benefits can be categorized into five areas:

- Data driven decisions
- Link product development and business aspect
- React to customers present needs

- Turn hypotheses into facts
- Steer development process into right direction

Data driven decisions are a way to gain competitive advantage by utilizing usage data or feedback in product development. Linking product development and business aspects verifies that the product is moving towards a direction with a higher business value. The consequence of utilizing data for product development is better knowledge of customer behavior, which can be used to better react to customers needs. Conducting experiments also allows companies to quantify business hypotheses and analytically derive answers, effectively allowing the use of hypotheses to steer development process.

In the case company, selection of the development ideas is based almost exclusively on the opinions of leaders in the organization. Studies report similar results in other companies as well [28]. However, based on the earlier stated conclusions that large parts of what companies do is wrong, business value can be better increased by moving towards data-based decision making.

Dialog is seen as a suitable candidate for experiments due to the extensive user interface and human users. Currently the product development is primarily guided by direct customer needs. Features planned for Dialog that are not direct customer needs are discussed in the product steering group. Especially A/B testing is useful to determine the best web-page version [7]. The metric for success differs based on the business goal. In Dialog the important metrics stated by the product owner and management include shortest path to complete a task, and amount of paths available to perform certain features. A path consists of clicks required to perform a single task, such as to create a marketing campaign for a new customer group.

"Developers can make improvement suggestions, and occasionally product owners invent feature ideas. The ideas are based on a guess that the customer might find this valuable."

While the development is therefore mainly guided by customer needs that directly bring in revenue, the feature development is guided by the need to generalize the product for general customer needs. While the low user amount for A/B testing - 3-5 users per customer - does not suit experiments based on statistical importance, the user actions can be tracked.

"The amount of available statistical data is not high much due to the user amount, but the users actions can be tracked, and they continuously perform repetitive tasks. This customer would then provide a lot of usage data. If the time spent by this user could be halved, the customer would make more money. It doesn't therefore necessarily require a large user amount."

A possibility expressed by a member of the management is to aid marketing by giving a MPV implementation of the product to new customers. This makes it is possible to pivot the product with customers, who can then test it in practice and make a decision whether to pay for the product or not. This strategy is similar to the lab website concept, where companies offer early versions of products free of charge for the purpose of collecting feedback [7].

Another experimental feature that is seen suitable for Dialog is an in-product survey. These surveys are then placed at strategic locations in the product to obtain active customer feedback when the customer is performing a certain task. It can provide remarkably valuable information as compared to purely quantitative passive data.

The list of things to consider for Dialog regarding continuous experimentation:

- Amount of own tasks vs customer tasks
- Types of experiments
- Statistical significance

CDM is perceived as a more challenging target for experiments due to the hard-coded interfaces and the nature of the application. The fact that there are no human users or a UI removes the possibilities of front-end A/B testing, the most typical experiment. However, the developers acknowledge that there is a need for better understanding the customers needs.

"Most of the work is strived to be done in a customer-oriented way. Even with a very good product owner, hardly ever does anyone know what the customer wants if it hasn't been asked."

Regardless of excluding most of A/B testing, continuous experimentation can be practiced with different types of experiments that can guide the development of CDM. Building alternative implementations of existing features with attached instrumentation can be used to determine if the current implementation or alternative implementation is preferable. Instead of relying on opinions, the decision can be made based on the feedback received. Most of the desired improvements are related to performance, and thus the important metrics are related to amount of queries performed in a given timeframe and the duration of a single query.

"CDM is a background system, and the performance related to queries is the most important metric for us. Imagine a situation where we're at a point analyzing how a query should be implemented, and whether to use a relational database or a document database. We can conduct an experiment where we implement both, test it with a data set and choose the more efficient query."

-CDM Team leader

Another possibility is to track the typical errors and to then focus the development towards fixing the most common errors.

"The software is composed of APIs, through which different components send and receive data. We could track the types of messages that are received and if they successfully pass. This data can then be used to develop the API further, and to for example track typical errors."

-CDM Developer

The lab website concept is less applicable for connected and embedded systems [7]. It would require the installation and integration of whole CDM, which is a laborious process.

The experiment scenarios listed for embedded systems by Eklund and Bosch [15], shown in chapter 3.4, contain cases that can be used for CDM development via experiments. Tracking the access amount of different queries in the API can be used to focus development towards most used calls. Identifying features that are not used at all is also a way to steer development into more beneficial cases. However, due to different customers using the product in different ways, has cross-case analysis be made to prevent failing to distinguish differences between users.

The list of things to consider for CDM regarding continuous experimentation:

- "Currently we do basically no product development, and all development is customer-oriented"
- b
- c

Feature alpha [7]

Minimum viable product principle is essential

current problems in roadmapping, choosing features, product development

"Vaatimuksia asiakkailta, tulee kuitenkin myös oman tiimin sisältä. Is-tumme usein viikottain alas, mietimme ja voi tulla omia kehitysideoita joita tiedämme että asiakkaat tulevat hyödyntämään. Asiakkailla on kuitenkin paljon omia juttuja joita he tarvitsevat. Nämä ovat kuitenkin sellaisia, joita muutkin asiakkaat voivat hyödyntää. Dialogissa pyritään siihen, että on vain yksi Dialog-versio kaikilla asiakkailla."

"Tavoite olisi 50/50. Puolet olisi meidän näkemystä millä voidaan tehdä bisnestä, loput asiakkailta. Nyt kuitenkin 100/0, ja yli vuoden aikana kaikki on tullut asiakkaiden tarpeista Dialogiin. Heikon palvelukuvauksen takia on tullu kiire tehdä asioita. A"

" Meillä on kuitenkin roadmappi, joka ei ole edennyt meidän johtamana vaan asiakasprojektien johtamana. Jos ajatellaan experimenttejä ja muita, pitäisi palautteen kerääminen saada tulemaan asiakkailta. Dialogin osalta tarvitaan sellaisia käyttäjiä jotka haluavat auttaa meitä, ja heillä on motivaatiota auttaa meitä koska he hyötyvät siitä itse."

discussion

"Asiakasfirmojen suhtautuminen heidän datan keräämiseen voi olla kysymysmerkki."

6.5 Implementing continuous experimentation as a development process

As requirements and benefits of continuous experimentation have been investigated in the previous sections, a plan to implement the practice in the case company is presented here. The continuous experimentation process is implemented by partially applying the continuous experimentation infrastructure by Fagerholm et al. [16].

6.5.1 Applying the general concept

As a general concept, continuous experimentation requires a paradigm shift from the teams and other essential parts of the organizations towards real-time delivery of value. Seeking feedback through the entire development process is essential, and entirely possible utilizing the scopes of experiments introduced by Bosch et al. [7]. However, due to the changes of the product development, marketing personnel and customers must be informed of the alterations. Customers might need to authorize the collection of their usage data, while the marketing team must be provided a way to track feature progress when moving towards real-time delivery from iterative model.

6.5.2 Applying the model

Fagerholm et al. define that a suitable experimentation system must be able to release minimum viable products or features with suitable instrumentation for data collection. The experimentation system must also allow the design and management of experiment plans, linking experiment results with a product roadmap and managing a flexible business strategy [16]. The Build-Measure-Learn cycle starts with forming a hypothesis and building a minimum viable product (MVP) with tools for data collection. Once the MVP has been created, the data is analyzed and measured in order to validate the hypothesis. To persevere with a chosen strategy means that the experiment proved the hypothesis correct, and the full product or feature can be implemented. However, if the experiment proved the hypothesis wrong, the strategy is changed based on the implications of a false hypothesis.

Kohavi et al. have.. [27] Tenets: Tenet 1: The Organization wants to make data-driven decisions and has formalized the Overall Evaluation Criterion (OEC) "For example, profit is not a good OEC, as short-term theatrics (e.g. raising price) can increase short-term profit, but hurt it in the long run" "Sessions per user, or repeat visits, is a much better factor in the OEC, and one that we use at Bing" Tenet 2: Controlled experiments can be run and their results are trustworthy Tenet 3: We are poor at assessing the value of ideas "Features are built because teams believe they are useful, yet in many domains most ideas fail to improve key metrics" [27]

Due to the very general nature of the infrastructure introduced by Fagerholm et al., no major modifications have to be made to address specific B2B needs.

- 1) Limiting the amount of experiments per project

| Type of experiment | Dialog | CDM | Rationalization |
|--------------------|--------|-----|-----------------|
| A/B testing | x | - | - |

A/B testing with small sizes

- 2) "A special test environment may be needed particularly in business-to-business markets, where the implications of feature changes are broad and there may be reluctance towards having new features at all"

"The length of the test cycle may thus have to be longer in business-to-business markets. Direct deployment could be more suitable for consumer markets, but we note that the attitude towards continuous experimentation is likely to change as both business and consumer customers become accustomed to it"

- 3) The practice of roadmapping in the case company is operated in product development steering groups, with the purpose to foresee the product development. The first step towards continuous experimentation is to inform the decision-makers of the practice and its benefits.

(back-end system consists of an experiment database..) -Customer environment?

Data collection: Kohavi et al. The first solution is to simply use an existing data collection tool, such as Webmetrics. However, most data collection systems aren't designed for statistical analyses, and the data might have to be manually extracted to an analysis environment. A different approach is local data collection, in which a website records data in a local database or log files. The problem with local data collection is that each additional source of data, such as the back-end, increases the complexity of the data recording infrastructure. The last model is a service-based collection, in which service calls to a logging service are placed in multiple places. This centralizes all observation data, and makes it easy to combine both back-end and front-end logging.

"Finally, experimentation may be conducted with several kinds of stakeholders. Apart from customers and end users, experiments could be directed

towards investors, suppliers, sales channels, or distributors. Companies whose product is itself a development platform may want to conduct experiments with developers in their platform ecosystem. These experiments may require other kinds of experimental artefacts than the MVP/MVF. Research on the types of experimental artefacts and associated experimental designs could lead to fruitful results for such application areas."

Necessary ingredients for agile development with controlled experiments [29]

1. A clear objective that can be evaluated. Controlled experiments require an OEC, and organizations need to agree on what to optimize. As Lewis Carroll said, "If you don't know where you are going, any road will take you there".
2. Easy to collect data about the user behavior. With client software, user behavior is hard to track and usually requires consent. As more of the user experience moves online, it becomes easier to track user behavior since server-side logging and client-side JavaScript logging are commonly used in industry and accepted as reasonable practices.
3. Easy to change and experiment with real users. Moving from hardware devices through standalone client software and software+services to SaaS applications, iterative experimentation gets easier. In standalone client software, experiments might be limited to beta cycles. For the software+services, experimentation capability has to be properly baked into the client so that server-side changes can be made to impact the client.
4. Sufficient users exist. Very small sites or products with no customer base cannot use experimentation, but such sites typically have a key idea to implement and they need quick feedback after going live. Because new sites are aiming for big improvements, the number of users needed to detect the desired effects can be relatively small (e.g., thousands of users). Large sites, which are typically better optimized, benefit even from small improvements and therefore need many

"The simplest experiment cycle would be" [7] TODO: fix this according to the findings!!

1. IDEA: Define OECs for both projects
2. The team(s) defines the experimental problem and produces new software which implements a feature thought to improve a user perceived quality. In this case there is just one factor with two levels, present and new implementation
3. The response variable(s) are selected and implemented

4. The software is deployed to a statistically relevant number of devices, with unmodified devices being the control
5. The software is run for a period of time
6. The data is uploaded and analysed through the infrastructure
7. The development team draws conclusions about the new software and decides whether to develop the new software design further, keep it as it is, or to drop it and revert to the unmodified software

-Apply in context

-List deviations

"Active customer feedback is concerned with surveys and other mechanisms where the customer is aware that he or she is providing feedback. Passive feedback and usage data is collected while the customer is using the system. Examples include the amount of time a user spends using a feature, the relative frequency of feature selections, the path that the user takes through the product functionality, etc. The low cost and ease of data collection leads to the next major difference between IES-based and traditional software." [7]

Trello

7 Discussion

This thesis was motivated by the general increased interest towards innovative experimental systems [34], and the lack of related studies in the B2B domain. While especially continuous delivery and partly continuous experimentation have been utilized and researched in the B2C domain, the benefits and challenges are still vague for B2B companies. Understanding the central aspects of continuous delivery and continuous experimentation will be a must for software companies willing to stay ahead of competitors in the current rapidly moving industry.

This thesis applied an exploratory view on the development model of the case company, attempting to deconstruct it into smaller parts.

The analysis of continuous delivery and continuous experimentation proceeded from higher-level considerations of B2B requirements to more detailed focus of applying the models in B2B context. The findings provide insights into the differences between B2C and B2B companies, the software products in these domains and the development models used. In addition to the difference in two domains, the findings also compare the development models from the viewpoint of two different software products.

7.1 Contribution

This thesis contributes to continuous experimentation [16] and organization evolution path [34] discussions by analysing the theories in B2B domain. These research areas are increasingly interesting in research, but none of the existing studies focuses purely on the B2B domain. Thus, the primary contribution of this thesis is to analyze the effects of this phenomena in a new domain.

The contributions to .. theories .. are ..

| Theory | Contribution |
|--------|--------------|
| x | y |
| x | y |

7.2 Further research

The scope was limited to the development model of two software products within the case company, and direct generalizations to other domains and products should be avoided. Since case studies only allow analytic generalisations, cannot the findings be directly generalised to other cases. This limitation applies especially to applying continuous experimentation in practice, since software products and customers inevitably differ.

However, the phenomena was deeply understood through gathering large amount of qualitative data and systematically analysing it. Therefore the core findings should be applicable to similar research problems outside of the empirical context of this thesis. This means that the B2B challenges and benefits of continuous delivery and continuous experimentation can be considered as a starting point for further studies in other contexts where these development models take place.

References

- [1] <http://www.kaushik.net/avinash/experimentation-and-testing-a-primer/>.
- [2] march 2014. <http://mcfunley.com/design-for-continuous-experimentation>.
- [3] august 2014. <http://tv.adobe.com/watch/max-2013/the-innovation-pipeline-how-adobe-defines-the-next-big-thing-in-web-design>.
- [4] *Experimentation platform*, march 2014. <http://www.exp-platform.com/>.
- [5] Adams, Rob J, Evans, Bradlee, and Brandt, Joel: *Creating small products at a big company: adobe's pipeline innovation process*. In *CHI'13 Extended Abstracts on Human Factors in Computing Systems*, pages 2331–2332. ACM, 2013.
- [6] Beck, Kent: *Extreme programming explained: embrace change*. Addison-Wesley Professional, 2000.
- [7] Bosch, Jan: *Building products as innovation experiment systems*. In *Software Business*, pages 27–39. Springer, 2012.
- [8] Cockburn, Alistair: *Agile software development*, volume 2006. Addison-Wesley Boston, 2002.
- [9] Crook, Thomas, Frasca, Brian, Kohavi, Ron, and Longbotham, Roger: *Seven pitfalls to avoid when running controlled experiments on the web*. In *Proceedings of the 15th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1105–1114. ACM, 2009.
- [10] Denzin, Norman K and Lincoln, Yvonna S: *The discipline and practice of qualitative research*. Handbook of qualitative research, 2:1–28, 2000.
- [11] Dubois, Anna and Gadde, Lars Erik: *Systematic combining: an abductive approach to case research*. Journal of business research, 55(7):553–560, 2002.
- [12] Dybå, Tore and Dingsøyr, Torgeir: *Empirical studies of agile software development: A systematic review*. Information and software technology, 50(9):833–859, 2008.
- [13] Dzamashvili Fogelström, Nina, Gorschek, Tony, Svahnberg, Mikael, and Olsson, Peo: *The impact of agile principles on market-driven software product development*. Journal of Software Maintenance and Evolution: Research and Practice, 22(1):53–80, 2010.

- [14] Easton, Geoff: *Critical realism in case study research*. Industrial Marketing Management, 39(1):118–128, 2010.
- [15] Eklund, Ulrik and Bosch, Jan: *Architecture for large-scale innovation experiment systems*. In *Software Architecture (WICSA) and European Conference on Software Architecture (ECSA), 2012 Joint Working IEEE/IFIP Conference on*, pages 244–248. IEEE, 2012.
- [16] Fagerholm, Fabian, Guinea, Alejandro Sanchez, Mäenpää, Hanna, and Münch, Jürgen: *Building blocks for continuous experimentation*. In *Proceedings of the 1st International Workshop on Rapid Continuous Software Engineering (RCoSE 2014), Hyderabad, India*, 2014.
- [17] Fowler, Martin and Foemmel, Matthew: *Continuous integration*. Thought-Works) [http://www.thoughtworks.com/Continuous Integration.pdf](http://www.thoughtworks.com/Continuous%20Integration.pdf), 2006.
- [18] Gephart, Robert P: *Qualitative research and the academy of management journal*. Academy of Management Journal, 47(4):454–462, 2004.
- [19] Healy, Marilyn and Perry, Chad: *Comprehensive criteria to judge validity and reliability of qualitative research within the realism paradigm*. Qualitative market research: An international journal, 3(3):118–126, 2000.
- [20] Humble, Jez and Farley, David: *Continuous Delivery: Reliable Software Releases Through Build, Test, and Deployment Automation*. Addison-Wesley Professional, 1st edition, 2010, ISBN 0321601912, 9780321601919.
- [21] Humble, Jez, Read, Chris, and North, Dan: *The deployment production line*. In *Agile Conference, 2006*, pages 6–pp. IEEE, 2006.
- [22] King, Nigel: *Template analysis*. 1998.
- [23] King, Nigel: *Doing template analysis*. Qualitative organizational research: Core methods and current challenges, pages 426–250, 2012.
- [24] King, Nigel, Cassell, C, and Symon, G: *Using templates in the thematic analysis of texts*. Essential guide to qualitative methods in organizational research, pages 256–270, 2004.
- [25] Kitchenham, Barbara, Pickard, Lesley, and Pfleeger, Shari Lawrence: *Case studies for method and tool evaluation*. IEEE software, 12(4):52–62, 1995.
- [26] Kitchenham, Barbara A, Pfleeger, Shari Lawrence, Pickard, Lesley M, Jones, Peter W, Hoaglin, David C., El Emam, Khaled, and Rosenberg,

- Jarrett: *Preliminary guidelines for empirical research in software engineering*. Software Engineering, IEEE Transactions on, 28(8):721–734, 2002.
- [27] Kohavi, Ron, Deng, Alex, Frasca, Brian, Walker, Toby, Xu, Ya, and Pohlmann, Nils: *Online controlled experiments at large scale*. In *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 1168–1176. ACM, 2013.
 - [28] Kohavi, Ron, Henne, Randal M, and Sommerfield, Dan: *Practical guide to controlled experiments on the web: listen to your customers not to the hippo*. In *Proceedings of the 13th ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 959–967. ACM, 2007.
 - [29] Kohavi, Ronny, Crook, Thomas, Longbotham, Roger, Frasca, Brian, Henne, Randy, Ferres, Juan Lavista, and Melamed, Tamir: *Online experimentation at microsoft*. Data Mining Case Studies, page 11, 2009.
 - [30] Loshin, David: *Master data management*. Morgan Kaufmann, 2010.
 - [31] Moran, Mike: *Multivariate testing in action: Quicken loan’s regis hadiaris on multivariate testing*. Biznology Blog by Mike Moran.
 - [32] Moran, Mike: *Do it wrong quickly: how the web changes the old marketing rules*. IBM Press, 2007.
 - [33] Neely, Steve and Stolt, Steve: *Continuous delivery? easy! just change everything (well, maybe it is not that easy)*. In *Agile Conference (AG-ILE), 2013*, pages 121–128. IEEE, 2013.
 - [34] Olsson, Helena Holmström, Alahyari, Hiva, and Bosch, Jan: *Climbing the "stairway to heaven"—a multiple-case study exploring barriers in the transition from agile development towards continuous deployment of software*. In *Software Engineering and Advanced Applications (SEAA), 2012 38th EUROMICRO Conference on*, pages 392–399. IEEE, 2012.
 - [35] Ōno, Taiichi: *Toyota production system: beyond large-scale production*. Productivity press, 1988.
 - [36] Palmer, Steve R and Felsing, Mac: *A practical guide to feature-driven development*. Pearson Education, 2001.
 - [37] Poppendieck, Mary and Poppendieck, Tom: *Lean software development: an agile toolkit*. Addison-Wesley Professional, 2003.
 - [38] Ries, Eric: *The lean startup: How today’s entrepreneurs use continuous innovation to create radically successful businesses*. Random House LLC, 2011.

- [39] Robson, Colin: *Real word research*. Oxford: Blackwell, 2002.
- [40] Runeson, Per and Höst, Martin: *Guidelines for conducting and reporting case study research in software engineering*. Empirical software engineering, 14(2):131–164, 2009.
- [41] Schwaber, Ken and Beedle, Mike: *Agile software development with scrum*. 2002.
- [42] Seaman, Carolyn B.: *Qualitative methods in empirical studies of software engineering*. Software Engineering, IEEE Transactions on, 25(4):557–572, 1999.
- [43] Stake, Robert E: *The art of case study research*. Sage, 1995.
- [44] Steiber, Annika and Alänge, Sverker: *A corporate system for continuous innovation: the case of google inc.* European Journal of Innovation Management, 16(2):243–264, 2013.
- [45] Strauss, Anselm and Corbin, Juliet M: *Basics of qualitative research: Grounded theory procedures and techniques*. Sage Publications, Inc, 1990.
- [46] Zelkowitz, Marvin V and Wallace, Dolores R.: *Experimental models for validating technology*. Computer, 31(5):23–31, 1998.

A Interview questions

| ID | Question |
|----|--|
| 1 | Name of the interviewee |
| 2 | Team of the interviewee |
| 3 | Position of the interviewee |
| 4 | Years in the company |
| 5 | Years of experience in the industry |
| 6 | What is the software product you're working with? |
| 7 | Describe your personal daily work routine. |
| 8 | Describe a normal week with your team. |
| 9 | What is the current development model in your team? |
| 10 | How has the development model evolved during your stay in the company? |
| 11 | Where do the development ideas come from? Are they mostly requirements from customers? |
| 12 | Are the current development ideas based on evidence or guesswork? |
| 13 | What is the current deployment process like? |
| 14 | How is it decided when to deploy? |
| 15 | How is it decided what to deploy? |
| 16 | How often is new version released? |
| 17 | How often are the new versions deployed to customer? |
| 18 | How are the deployment dates chosen? |
| 19 | Which parts of the deployment process are manual? |
| 20 | Which parts of the deployment process are hard to measure automatically? |
| 21 | Is the customer involved in the deployment process? |
| 22 | Does the customer have to do something when a version is deployed? |
| 23 | What are the strengths in the current deployment process? |
| 24 | What are the problems in the current deployment process? |

| ID | Question |
|----|--|
| 25 | Describe the customer interaction process. |
| 26 | From your point of view, what are the challenges with the customer interaction? |
| 27 | From your point of view, what are the strengths with the customer interaction? |
| 28 | From your point of view, how could the interaction with the customer be improved? |
| 29 | Is it common for customers requirements to change? |
| 30 | Is the development team aware of the customers present requirements? |
| 31 | How is feedback collected from the customer? |
| 32 | What are the B2B specific challenges in feedback collection? |
| 33 | From your point of view, how could the feedback collection be improved? |
| 34 | How is the customer feedback used? |
| 35 | From your point of view, how could the feedback usage be improved? |
| 36 | How many end-users does the product have? |
| 37 | How is usage data collected? |
| 38 | How could usage data collection be improved? |
| 39 | What challenges would real-time deployment have? |
| 40 | What challenges would plugged-in data collection instruments have? |
| 41 | If data were to be collected in customer environment, what challenges would be faced storing it? |
| 42 | What are the strengths in the current development process? |
| 43 | What are the problems in the current development process? |
| 44 | Could your product be instrumented with a data collection service? If not, why? |
| 45 | In your product, could experiments be quickly deployed to the customer environment? |
| 46 | In your product, could the development process be guided by A/B testing? |
| 47 | Does your team have a data analyst? |