

Backup Management and Orchestration System

Final Year Project

Rob Shelly
Applied Computing
20068406

December 5, 2017

Contents

1 Introduction 1

1.1 Problem Statement 1

1.2 Aims & Objectives 1

2 Technologies 2

2.1 Docker 2

2.2 Amazon Web Services 2

2.3 Jenkins 2

2.4 Node 3

2.5 React 3

3 Design 4

3.1 System Architecture Overview 4

3.2 Formal Modelling 5

3.2.1 Sequence Diagrams 5

3.2.2 User Stories 7

3.3 Front End Design 9

3.3.1 Wireframes 9

4 Methodology 11

4.1 Agile 11

4.1.1 Scrum 11

4.2 CI/CD with Jenkins 11

4.3 JIRA 11

4.4 Tool n 11

4.5 Testing Approach 11

5 Implementation of Prototype 12

5.1 Sprint 1 12

5.2 Sprint 2 12

5.3 Sprint 3 12

6 Bibliography 13

1 Introduction

1.1 Problem Statement

In January 2017 GitLab suffered a data loss incident which was widely reported in media. It began with spammers targeting GitLab.com and culminated in an engineer erroneously deleting 300GB of PostgreSQL data in a production environment. The lost data included merger requests, users and comments ([GitLab, 2017a](#)). The bigger story was to come later however when it was realised that GitLabs backup process had failed silently. The backups did not exist, resulting in a total loss of the data. As it transpired, conflicting major versions of *pg_dump* (a utility for backing up PostgreSQL databases) in use for the backup procedure and the PostgreSQL database resulted in an error, and the procedure failing ([GitLab, 2017b](#)).

The incident was widely reported in the tech industry with the story being picked up by a number of outlets including TechCrunch ([2017](#)) and The Register ([2017](#)). For many, the focal point of the story was the failed backups. The incident highlighted the need for regular verification of backups. A simple way of performing this verification is to regularly restore data. The method of verification is to perform a restore of the data, which can be a mundane and time consuming task. The aim of this project is to create a solution to the issue. A system which can notify administrators when backups have failed may have prevented the data loss in the GitLab ordeal.

1.2 Aims & Objectives

The overall objective is to create a system to test that uncorrupted backups exist and contain valid, readable data. A system will be created that allows sysadmins to test backups and to schedule the regular testing of backups. This will be achieved by performing restoration on the backups. The main objectives of the system are as follows:

- AO1** Eliminate the mundane and time consuming task of backup testing by automating regular backup restorations and recording results;
- AO2** Catch silent failures of the backup procedure by notifying sysadmins of failed backups;
- AO3** Reduce the cost of backup restoration testing by automating the process of creating the necessary infrastructure (such as virtual machines on AWS), performing the restoration and destroying the infrastructure once results are obtained, thus minimising the uptime of infrastructure;
- AO4** Performs the restoration check in a secure manner by managing encryption keys and the movement and decryption of data only when necessary in safe environment.

The system will focus on backups of databases. For scope, design will focus on testing MongoDB data and MySQL data, thereby providing a sample of both relational and non-relational database management systems (DBMS). However, the system should be designed such that it can easily modified to test data from others forms of database management systems. As part of the system, the following should be implemented:

- **Web app:** This will act as a front end for the sysadmins to run and schedule tests and view results.
- **Automation Server:** This will be the backend of the system. It will take care of retrieving the backup data before performing some sorts of tests.
- **Container Platform:** This will be utilised by the backend to test the server. For example, when testing the data from a MongoDB database, the backend will spin up a container with MongoDB installed in order to verify the data.

2 Technologies

2.1 Docker

Docker is a container platform for building and managing applications. This project is interested not in Docker's platform but rather in the Docker images that run on the platform. A container image is a modular piece of software. It encapsulates all the code and tools needed to run the software packaged in the image. The image can then be run in a container on any environment using a container platform or service. Thus, it runs independent of the hardware or operating system. The container also isolates the software from other images and software running within the environment ([Docker, 2017](#)).

The modularity of software makes Docker images appealing for this project. It will allow testing various data base types (e.g. MongoDB, MySQL) through its software agnostic feature, by deploying an image with the corresponding DBMS software.

2.2 Amazon Web Services

The project will make extensive use of Amazon Web Services (AWS) with most or possibly all of the systems infrastructure deployed on AWS. More specifically the project will make use of two specific services; Elastic Compute Cloud (EC2) and EC2 Containers Service (ECS).

EC2 is Amazon's compute service. It allows easy deployment and management of virtual compute resources within the cloud. The flexibility of operating systems, virtual machines (or instances as they are known in AWS) and size of volume of storage make it ideal for this project ([Amazon, 2017a](#)). It will allow the system to create instances with only the necessary resources required (i.e. memory and storage) to test the restoration of a given backup. This keeps the cost of testing to a minimum in keeping with Aim [AO3](#).

ECS is Amazon's container management service. It allows Docker images to be easily deployed to and run on EC2 instances without the need to install Docker on the instances. ECS takes care of much of the container management issues that would arise when deploying a services if implemented through Docker alone. This includes managing port mappings between container ports and host ports, ensuring all containers are accessible if necessary. There is no added cost for using ECS. i.e. the customer only pays for the EC2 instances ([Amazon, 2017b](#)).

ECS is an appealing platform for running Docker images for the following reasons.

- Images can be deployed on EC2 instances, meaning there is no need to install Docker on the instance, without any extra charge.
- Containers are created within the customers own EC2 instances, meaning they are not exposed to other AWS customers. They are secured by the same infrastructure created by the customer for their instances, for example Virtual Private Clouds (VPCs) and Security Groups.

AWS also features a command line interface for building, modifying and destroying infrastructure across all of its service, including EC2 and ECS. This provides an programmatic method of creating the resources necessary to perform test restorations. The ability to do so allows for the automation of the infrastructure creation and destruction (after testing). This is an objective of the project set out in Aim [AO1](#).

2.3 Jenkins

Jenkins is an automated build server. It can be used to implement continuous integration (CI) and continuous delivery (CD). Configuration and management of the server can be achieved using both a web interface and an API. Jenkins service is also extensible through a library of plugins ([Jenkins, 2017](#)). One plugin which will be of particular interest to

this project is the Pipeline plugin. It allows the creation of pipeline as code. This means that for this project, pipelines can be used to create EC2 instances and deploy the necessary Docker images using ECS in order to test backup.

Jenkins was chosen as a backend service for this project as it, along with its library of plugins, presents many useful features which will be beneficial to the implementation of the system. These include the following:

- A built in email notification system which can be used to notify users of silently failed backups. This provided the functionality to implement a satisfactory solution to Aim [AO2](#);
- A Credentials plugin which provides a means of storing various credentials in various forms (e.g. username/password pairs, SSH keys) along with a standard API for Jenkins and other plugins to access and use these credentials ([Connolly, 2017](#)). This provides a secure manner for using SSH keys for backups servers as well as encryption keys for sensitive backups. This is a key objective of the project outlined in Aim [AO4](#).
- The ability to schedule jobs to run at regular intervals will provide the functionality described in Aim [AO1](#). This eliminates the need for the development of a scheduling system in order to fulfil the systems requirements.
- The REST API can be used to by a user-friendly web based frontend, allowing users who unfamiliar with Jenkins or AWS to perform test restorations.

2.4 Node

The frontend of the system will be designed using Node (also known as Node.js). Node is a JavaScript runtime environment for building network applications. It is light-weight and efficient framework through its event driven, non blocking I/O implementation.

The default package manager of Node is *npm* (for Node Package manager). It is the worlds largest software registry ([NPM, 2017](#)). The vast registry of free and open source packages available through make Node an attractive choice for this project. Of particular interest are the multiple Node clients for Jenkins. These are Node wrappers for the Jenkins REST API enabling easy integration of the frontend with the Jenkins backend.

Although any of a number of frameworks could have been used, for example Django, Node was chosen for this project due its light-weight design and extensibility through *npm*, including the aforementioned Jenkins API wrappers.

2.5 React

The UI element of the frontend will be built using React, a JavaScript library available through *npm* for building user interfaces. React is developed to work independently of other technologies, meaning it can be integrated easily with Node and other *npm* packages without the need for refactoring. React builds UI's as a set of components, each managing and their own state and implementing their own render function. This allows fast and efficient rendering as data changes as only components that are updated will be re-rendered.

3 Design

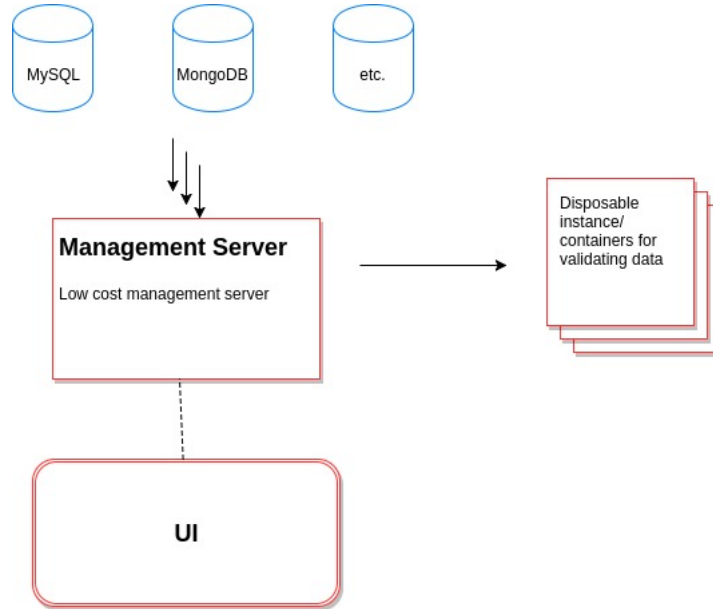
3.1 System Architecture Overview

The system will comprise of three main components:

- Management Server
- User Interface
- Disposable instances/containers

The system will also use existing infrastructure. This is where the backups are stored. Depending on the user of the system there may be multiple backup server in different location (such as AWS regions) or for different data types (relational and non-relational databases). Backup data may be stored in a variety of ways such as on EC2 instance or S3 buckets.

Figure 1: Diagram of System Architecture



Management Server: This will be a small low cost AWS instance on which the Jenkins automation server will be installed. The majority of the systems functionality will be carried out and/or orchestrated by this server. Jenkins jobs will copy the backups from their location to a disposable instance and implement the necessary steps to validate them such as importing and and reading.

User Interface: This will provide a simple user interface (UI) for the system, implemented as a simple web app, hosted on AWS. It will allow users with little knowledge of Jenkins and AWS to perform backup restoration checks by adding a layer of abstraction. Users will be able to run restorations by providing the parameters such as the backup file and it's location. The UI will utilise the Jenkins API to run execute the restoration with the parameters provided.

Disposable Instances or Containers: Disposable infrastructure will be used to perform the restoration. EC2 instances or containers can be used to quickly and easily deploy the necessary software to perform the restoration (i.e. the correct DB management system). They can also be destroyed afterwards, destroying the data and therefore maintaining confidentiality.

3.2 Formal Modelling

3.2.1 Sequence Diagrams

The main function of the systems have been demonstrated below in sequence diagrams. Figure 2 shows the process of running a single backup restore. This involves a user manually triggering a restoration using the web interface. The trigger a Jenkins job automates the remaining steps. The backup is copied from the backup server to the test restoration server where it is imported to a Database Management System such as MongoDB. A read of the data is then performed to verify that the data is uncorrupted and readable. Finally it is destroyed from the restoration server.

Figure 2: Run Restore

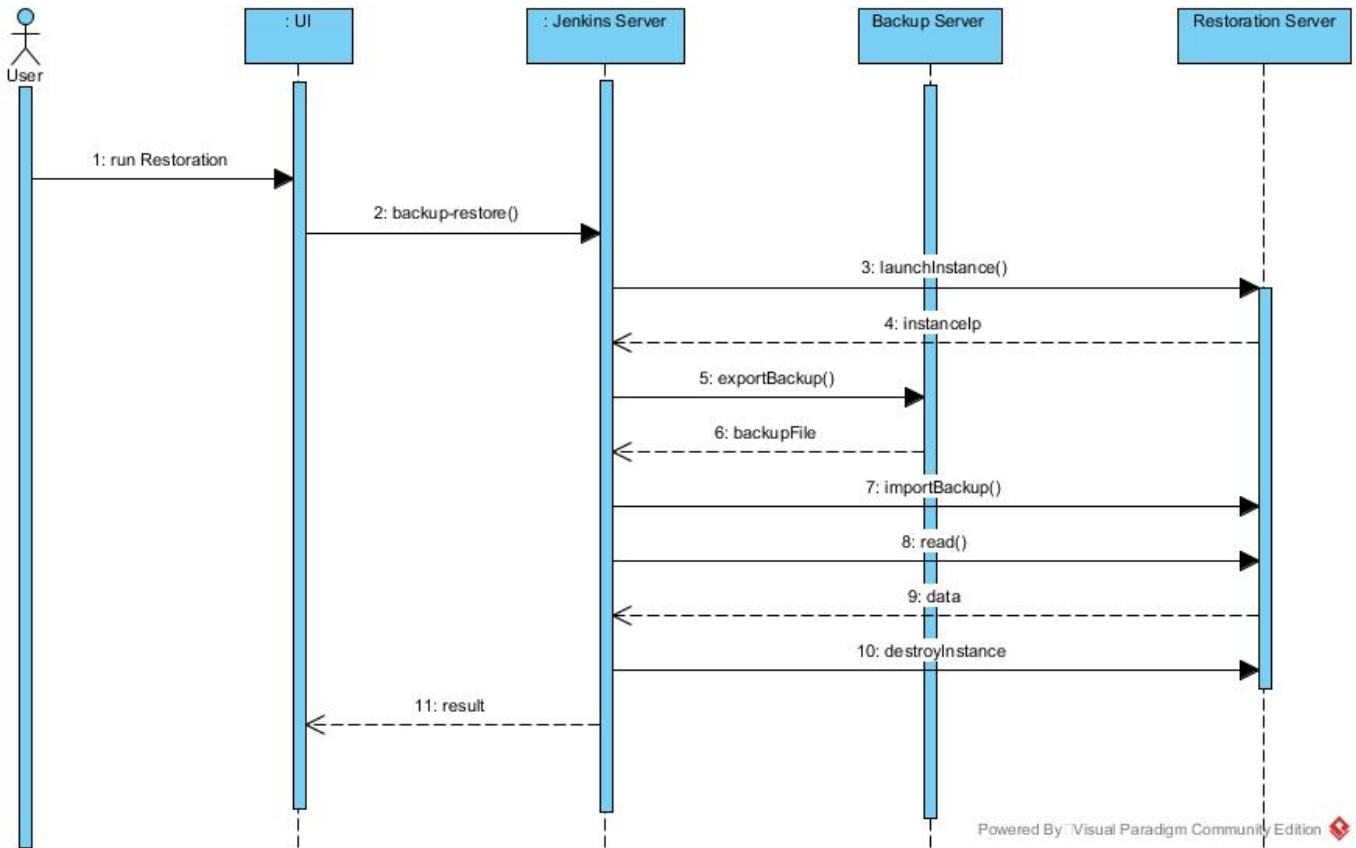


Figure 3 shows the process of a scheduling regular backup restoration tests. Again, this is triggered by a user from the web interface. The web interface will pass the JSON or XML configuration for a job to the Jenkins server. The server will then create and save the job. A status indicating whether the job was created successfully is returned to the user.

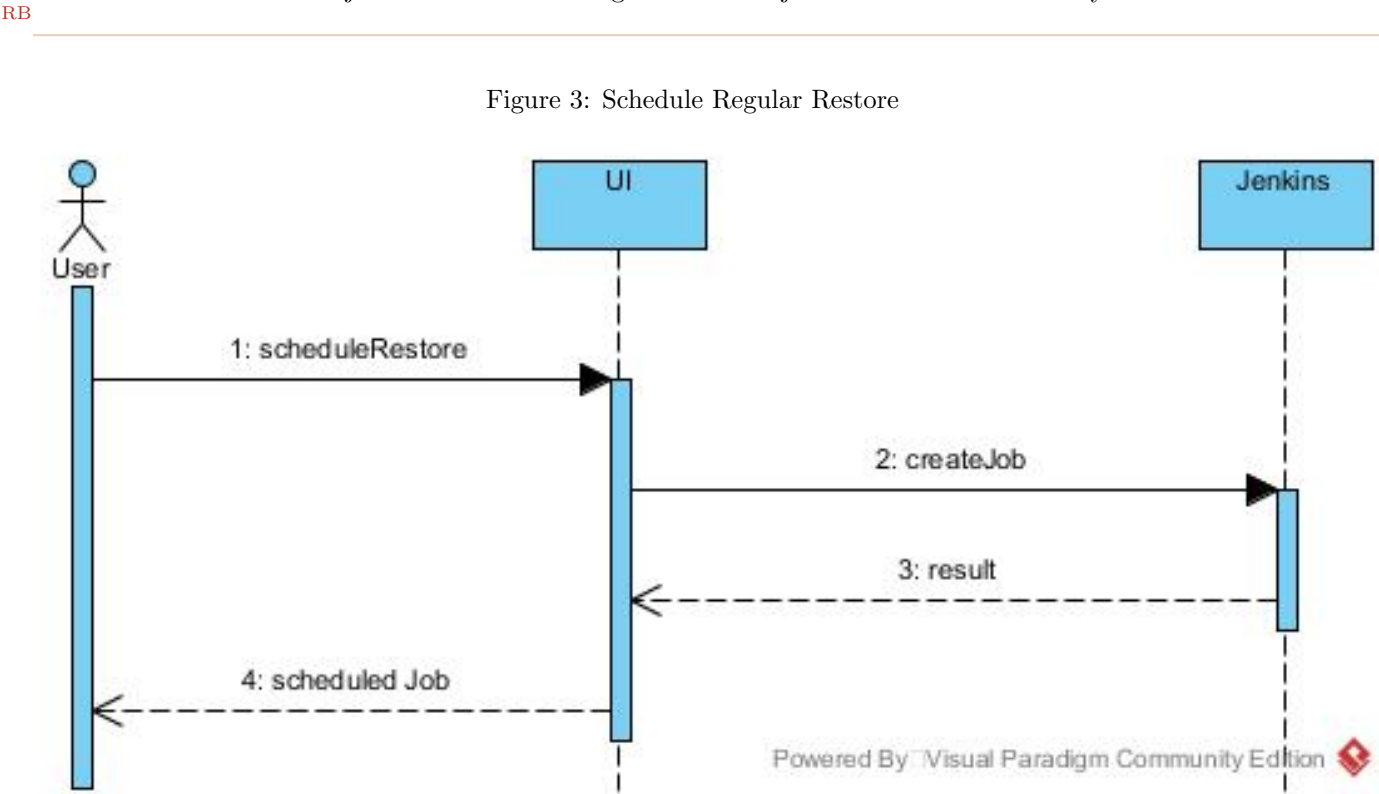
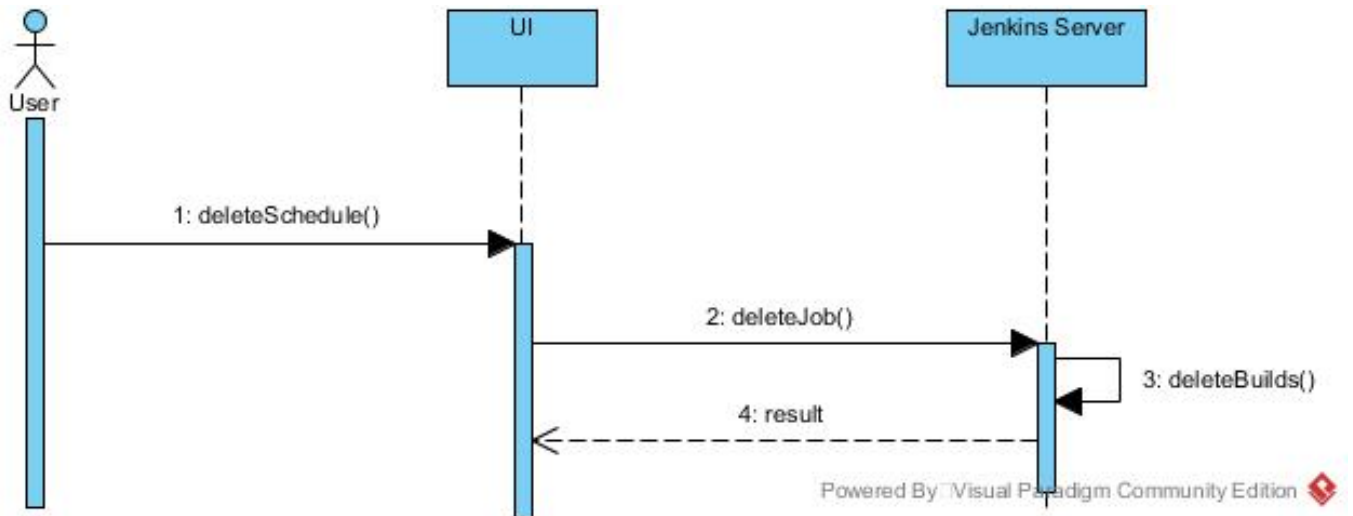


Figure 4 Show the process of deleting an existing scheduled job. This is required if a user not longer want to run scheduled restoration of a particular backup (for example if that backup is non longer needed and deleted). The user must delete the job on the Jenkins server in order to prevent further attempted restorations running. The user triggers this process from the web interface. This sends a delete commands to the Jenkins server via the API to remove the schedule job. The status of the command, indicating a successful or failed restore, is returned to the user.

Figure 4: Delete Scheduled Restore



3.2.2 User Stories

User stories are provided in [Table 1](#). There are two users of the system; managers and users.

Managers will have the ability to configure security aspects of the systems. This includes not only granting users access to the system but also adding credentials such as SSH keys and decryption keys. They will also be able to perform all the tasks that a regular user can perform without needing a set of regular user credentials.

Regular users of the systems will be able to perform backup restoration after logging in. They will have the ability to run a restore at will and view the results. However, they will also be able to schedule the regular restoration of a backup and view the results of each restore that has occurred.

These are the stories for general user such as running and scheduling restores. Also included are administration user stories. As the system is created and destroys infrastructure on AWS, it would be necessary to limit access to the system.

Table 1: User Stories

#	As a	I want to be able to	so that
US1	manager	implement a user system	I control who can run backup restores
US1.1	manager	add my team members to the system	they can run backup restores
US1.2	manager	remove users from the system	former team members no longer have access
US2	manager	add and control sensitive information within the system	I can implement a security policy
US2.2	manager	securely store credentials within the system	they don't need to be entered every time a restore is executed
US2.3	manager	add SSH keys for backup server	the system has secure access backup server
US2.4	manager	add decryption keys for backups	encrypted backups can be decrypted for testing
US2.5	manager	delete SSH keys	expired/outdated credentials are no longer stored
US2.6	manager	delete decryption keys	expired/outdated credentials are no longer stored
US3	manager	execute all same tasks as a regular user	I don't need a second set of credentials to run restores myself
US4	user	login	I can run restores
US5	user	logout	I avoid potential unauthorised access
US6	user	run a test restoration of a backup	I can verify that the backup exists, is a valid file, and is readable
US6.1	user	run a test by filling out a simple form with basic parameters (location, filename) of the backup to test	I can easily run a restore of a specific backup without needing to worry about the implementation
US6.2	user	view the current status a running restoration	I can review the progress of long running restores
US6.3	user	check if a backup failed or succeeded	I can immediately investigate any failed backups
US7	user	create a schedule of automated restores for a given backup	I don't have to manually execute them myself on a regular basis
US7.1	user	choose the frequency of automated restores within a schedule, from daily through weekly to monthly	I control how often different backups are tested
US7.2	user	check if an automated restoration has started	I can verify my schedule is working correctly
US7.3	user	check the results of an automated restore	I can immediately investigate any failed backups
US7.4	user	view the all past results of an automated restore schedule	view the consistency of my backups success
US7.5	user	modify a scheduled restore	I can change the frequency of a scheduled restore
US7.6	user	the parameters of a schedule	any changes to the backups, such as location, will be reflect in the restoration schedule
US7.7	user	delete regularly scheduled restores	old backups/deleted backups are no longer tested
US8	user	view feedback of a failed restore	I might gain an insight into the fault in the backup
US9	user	notified when a restoration fails	silent, unnoticed fails are avoided

3.3 Front End Design

3.3.1 Wireframes

Figure 5: Homepage

Backup Restoration Test System

Run a Test Restoration on Backup

Backup File

Location

Type

Run

Schedule Regular Testration Test

Backup File

Location

Type

Frequency

Time

Repeat

Schedule

Scheduled Restores

Name	File	Location	Last Run	Successful	
					Run Now

9

Backup Restoration Test System

[illegible]

Delete

Date	Time	Status	Link to Jenkins
██████████	██████████	██████████	██████████
██████████	██████████	██████████	██████████
██████████	██████████	██████████	██████████
██████████	██████████	██████████	██████████

4 Methodology

4.1 Agile

The design methodology chosen for this project is Agile. Agile takes an iterative approach to designing and delivery products. It is a goal driven methodology that aims to build and deliver software iteratively and incrementally from the beginning of the project. This is in contrast with more traditional approaches such as Waterfall which deliver in one final stage. A notable aspect of Agile is user stories. The project is broken down in small sections of functionality which can be independently developed and delivered upon completion ([Rasmusson, 2017](#)).

4.1.1 Scrum

A particular Agile framework which will be used for this project is Scrum. Scrum organises development in to cycles of *sprints*. A sprint consists short time-limited periods of development each with it's own development goals based on work within the backlog. Each sprint will consists of regular update and a final review/retrospective before beginning the next sprint. The Scrum master keeps the sprint focused on its development goal ([ScrumAlliance, 2016](#)).

Scrum is an ideal model for developing this project. The project supervisor plays a role in line with the concept of a Scrum master. Also, the use of user stories means sprints can be aligned with the implementation of stories.

4.2 CI/CD with Jenkins

Continuous Integration/Continuous Deployment (or continuous Delivery) is a development concept that focuses on the frequent and automated testing building and releasing of code. It aims to remove the large workload required when it is time to release a version or update of a product by performing the same process in a automated manner on every code commit ([Pittet, 2017](#)).

Continuous integration refers to preparing the code for release and often as code commits are performed. For example, running tests and building Docker images on each commit. IT means that code is prepared for release at each stage of development, instead of when it come to release time. This may occur often as many times a day ([Ramos, 2016](#)). Continuous Deployment is a step beyond Continuous Integration. After code is prepared for release, the built code is deployed to a server. However, this may be a development server. Pushing the built code to production required a manual trigger. Continuous Delivery automates this final manual trigger, meaning the entire process of moving code through testing, building and deployment to production is entirely automated ([Ellingwood, 2017](#)).

For this project, CI/CD (continuous development in this case) will be implemented using a Jenkins automated build server. The fronted web app will be built as a Docker image and deployed to ECS by Jenkins on every code commit.

4.3 JIRA

// TODO

4.4 Tool n

4.5 Testing Approach

//TODO

5 Implementation of Prototype

5.1 Sprint 1

5.2 Sprint 2

5.3 Sprint 3

6 Bibliography

- Amazon. Amazon EC2, 2017a. URL <https://aws.amazon.com/ec2/>.
- Amazon. Amazon EC2 Conatiner Service, 2017b. URL <https://aws.amazon.com/ecs/>.
- Stephen Connolly. Credentials plugin, 2017. URL <https://wiki.jenkins.io/display/JENKINS/Credentials+Plugin>.
- Docker. What is a container, 2017. URL <https://www.docker.com/what-container>.
- Justin Ellingwood. An introduction to continuous integration, delivery, and deployment, May 2017. URL <https://www.digitalocean.com/community/tutorials/an-introduction-to-continuous-integration-delivery-and-deployment>.
- GitLab. Gitlab.com database incident, February 2017a. URL <https://about.gitlab.com/2017/02/01/gitlab-dot-com-database-incident/>.
- GitLab. Postmortem of database outage of january 31, February 2017b. URL <https://about.gitlab.com/2017/02/10/postmortem-of-database-outage-of-january-31/>.
- Jenkins. Jenkins, 2017. URL <https://jenkins.io/>.
- Natasha Lomas. Gitlab suffers major backup failure after data deletion incident, February 2017. URL <https://techcrunch.com/2017/02/01/gitlab-suffers-major-backup-failure-after-data-deletion-incident/>.
- NPM. Npm, 2017. URL <https://www.npmjs.com/>.
- Sten Pittet. Continuous integration vs. continuous delivery vs. continuous deployment, 2017. URL <https://www.atlassian.com/continuous-delivery/ci-vs-ci-vs-cd>.
- Marcia Ramos. Continuous integration, delivery, and deployment with gitlab, 08 2016. URL <https://about.gitlab.com/2016/08/05/continuous-integration-delivery-and-deployment-with-gitlab/>.
- Jonathan Rasmusson. Agile in a Nutshell:what is agile, 2017. URL http://www.agilenutshell.com/how_does_it_work.
- ScrumAlliance. Learn about Scrum, 2016. URL <https://www.scrumalliance.org/why-scrum>.
- Simon Sharwood. Gitlab.com melts down after wrong directory deleted, backups fail, February 2017. URL https://www.theregister.co.uk/2017/02/01/gitlab_data_loss/.