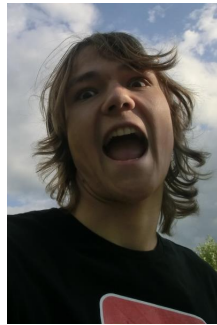


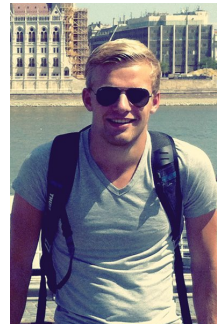
CONTEXTPROJECT PROGRAMMING LIFE
GROUP 2 - GEVATT
FINAL REPORT
TU DELFT



Ruben Bes
rbes
4227492



Mathijs Hoogland
mhhoogland
4237676



Jasper Denkers
jdenkers
4212584



Robbert van Staveren
rhvanstaveren
1527118



Willem Jan Glerum
wglерum
4141040

June 25, 2014

Abstract

This is the final report for the Programming Life Contextproject, a second year course from Computer Science at TU Delft. The Contextproject course is about applying all learned skills in a particular context at developing a piece of software. In this case the context was bioinformatics and we worked 10 weeks with a team of five people.

This document contains the main information about development, implementation and validation of the product. Main features of the product will be presented and it will be discussed how they satisfy the user needs. Furthermore, this document will contain an HCI module about the interaction of users with the product. Finally, an outlook will be given to show what possible improvements could be implemented if this project will continue in the future.

Besides this document several other documents are made covering other parts of the project. This is the final document consisting of the most information about the project.

Contents

1	Introduction	3
2	Overview of the product	3
2.1	Authentication	3
2.2	Context specific	4
2.2.1	Data analysisation	4
2.2.2	Data retrieval	4
2.2.3	Visualisation	4
3	Description of the functionality	5
3.1	Webapplication	5
3.2	Authentication	5
3.3	Context specific	5
3.3.1	Data analysisation	6
3.3.2	Data retrieval	6
3.3.3	Visualisation	6
4	Human Computer Interaction Design	7
4.1	Our action plan	7
4.2	The realisation of our experiment	7
4.3	Debriefing	8
4.4	Conclusion	8
5	Evaluation of the functional models	9
5.1	Evaluation of functional models	9
5.2	Failure analysis	9
6	Outlook	9
7	Appendix A: Screenshots	10

1 Introduction

The application we developed is called GEVATT, which stands for GEnetic Variations Analyzer Through Triodata, which is an application to be used by doctors to browse genetic information of patients. Such applications are also called genome browsers. It is a secured web application based on the Play Framework. Doctors need to login from any device with a browser in order to use the application.

The application makes it possible for a doctor to upload a VCF file and let the application analyse it. A VCF (Variant Call Format) file is a file used in bioinformatics for storing gene sequence variations. This storing of genetic information of patients is based on variations between the patient and a reference genome. This is done because storing all information would be taking too much space.

After uploading the file, the user of the application (a doctor) waits until the file is processed. In the meanwhile the doctor could browse other patients he uploaded information about earlier. After uploading, the first important part of the application is executed: analysing the data. The outcome of the analysis consists of mutations found in the genome of the patient and the relations between these mutations.

Secondly, the main focus of the application is about visualising the found mutations. This is firstly done by giving a main overview of the whole patient. An overview of all chromosomes is given, and per chromosome is indicated if and how many mutations it contains. Besides a visual overview of the chromosomes there's a tabular overview with all mutations and some extra information per mutation. This makes it easier to estimate which mutations are more harmful than others.

Most information is shown on the overview pages per mutation. The pages with these visualisations have two visualisations which both give another insight in the mutation. The first part shows the position of the mutation relative to a gene. The second visualisation is a graph-based interactive visualisation showing proteins related to the mutation and the connections between these proteins. Here we also have a tabular overview per protein with about which diseases they could cause and to which other mutations of the patient they are related.

2 Overview of the product

The developed product, a secured web application, is built on the Play Framework. It contains several parts that work together to deliver a user friendly environment to explore genetic information.

2.1 Authentication

Starting with the basis, it is a secured web application where a user needs to log in. After logging in some secured pages become accessible. Some other pages are already accessible before logging in, e.g. the 'documentation' and 'about' pages. We distinguish the following parts:

- A login page

- Secure pages that are not publicly accessible
- Redirect if secured page is requested while unauthenticated
- Prevent doctors from accessing data of patients of other doctors

2.2 Context specific

The web framework and its built-in secure module is not developed by ourselves so we could focus on developing the context specific parts. This was mainly separable in three parts: analysing data, retrieval of relevant data from databases and visualising data.

2.2.1 Data analysis

Analysing the data is done by processing the VCF file and detecting mutations. The outcome of the processing is used in the visualisation. At both of these parts there is information used given by some databases. The information is used to get more information about mutations and find relations between several mutations. The application does the following:

- Read the VCF file and find mutations of two types: de novo's and recessive homozygous
- Save metadata about the uploaded file
- Find relations between found mutations

2.2.2 Data retrieval

Information from multiple databases is used for both visualisation and finding relations between mutations. This includes:

- Manage connections to multiple databases
- Querying databases to get relevant information about mutations and proteins

2.2.3 Visualisation

The most important part of the application: visualisation. Multiple views of the same found mutations are made to give as much insight as possible. This includes the following visualisations:

- An overview of all mutations in a patients VCF file with a view per chromosome
- A distinct page per mutation, showing:
 - The position of the mutation relative to nearby genes
 - A graph with a protein related to the mutation, and proteins related to that protein

3 Description of the functionality

This chapter covers a detailed description of the developed functionality in the product. This description is based on the overview of all functionality given in the previous chapter and therefore has the same structure.

3.1 Webapplication

The application is a web-application based on the Play Framework. This means the application is accessible via devices with an internet connection and a browser. The layout is optimised for devices with a screen with a minimum width of 1024 pixels. The page structure is as follows:

- A login page (see authentication below)
- A dashboard page with an introduction to the application and links to the following context specific parts:
 - The patient overview with a sortable list of all patients with some meta data. If the user clicks on a patient in the overview, he gets redirect to a page specific to the selected patient
 - On the patient page an overview of the found mutations is given. It contains a visual overview of the chromosomes and the mutations found per chromosome
 - A separate page with visualisations and information (see visualisation below) per mutation
 - There is a separate page for adding new patients to the database

3.2 Authentication

The authentication part takes care of securing the pages and data of the application to prevent it from being accessible to everyone. It does the following:

- If a user is not authenticated and tries to open a secured page, he is redirected to the login page
- At the login page, a user needs to supply his username and password to sign in. When these credentials are not recognised, the user receives an error message
- When logged in, each page contains some user specific information, like his name. Furthermore, links to application specific pages become visible and a logout button appears

3.3 Context specific

We could separate the context specific functionalities in data analysis, data retrieval and visualisation. These will be discussed individually below.

3.3.1 Data analysis

The analysing of data is done after a VCF file is uploaded.

- When a user has added a patient by filling in the associated information and uploading a VCF file, the user gets redirected to the patient overview. In the background the file gets processed, so the user can continue browsing. After processing, the patient overview gets updated automatically and the patient view becomes available.
- The mutations that are found are stored in a databases dedicated to the application
- The application searches for relations between the found mutations and stores them in the database

3.3.2 Data retrieval

We use the following databases to retrieve relevant information:

- CADD[1]

This database is used for retrieving scores representing the dangerousness of mutations. It contains all 8.6 billion mutations compared to the reference genome. These scores range from 0 to 99, and estimate how dangerous a mutation is. A score between 0 and 10 means the mutation is among the 90% least dangerous ones, greater than 10 means among the top 10%, greater than 20 means top 1% and so on.

- dbSNP[2]

This database is used for retrieving all kind of information related to SNPs. This is by far the most used database in this project, as all we start with are rsIDs, which we have to look up in dbSNP. Here we retrieve information concerning the location of the SNP, the gene it is on and the frequency of its occurrence.

- STRING[3]

This database is used for retrieving everything related to proteins. We mostly use the genes connected to SNPs found with dbSNP to find connections with proteins. These proteins are then checked for links to other proteins, and how much interaction is between them. This is shown with a number between 0 and 999, as shown in Figure 5. We can also use the name of proteins to check if there is a disease associated with them, thanks to a file from the ClinVar database[4].

3.3.3 Visualisation

Several visualisations are used to show retrieved information in a clear way to the user.

- The first visualisation is about giving an overview of all found mutations, showing them per chromosome. This overview is given on the patient overview page. It contains a graphical representation of all chromosome

pairs and each pair is colored red or black if it respectively contains a mutation or not. While hovering over a chromosome pair, a list appears with links to the mutations found in that pair.

- On the pages for individual mutations there are two extra visualisations
 - The first mutation-specific visualisation is the top one found on the mutation page. It shows the part of the chromosome that the mutation is on, and displays the position of the base pair and the genes relatively close to this pair.
 - There's also an overview in the form of a graph containing proteins related to the mutation. The protein directly related to the mutation is marked and proteins related to the marked one are shown in the graph. This graph is interactive so proteins can be dragged around to get a clearer view. Connections between protein are shown darker and thicker if the connectivity between two proteins is high relative to the other connections in the graph.

4 Human Computer Interaction Design

As described before, we developed an application for doctors to browse through genetic information. A big challenge in this assignment was to make our product intuitive for the end users: people that have a completely different background than we, the developers, have. We couldn't just deliver a product, and hope it was good, we had to have some feedback from a person that was not involved in the project (but still had biological background).

4.1 Our action plan

After meeting with Willem-Paul Brinkman, we decided that we were going to organise an empirical evaluation. This evaluation should be of the type experiment, the experiment would make use of the **think aloud protocol**. We had a fellow student, Max, test our application. He was the most suitable person we could find, as he studies nanobiology at the TU Delft. We were going to give him a computer with our application on it, and give the following assignment:

You want to browse through the genome of patient John Doe, he is sick. He has the symptoms of Marfan's syndrome, can you look if you can find any traces of this sickness?

We deliberately added the patient to the application, because adding a new patient and uploading his SNPs aren't the main concern in our project.

4.2 The realisation of our experiment

As expected, our test person didn't have any trouble to find his way through our dashboard (See Figure 1), this dashboard has a clear "Getting started" instruction that easily helps the user to navigate to the Patients page. Also when at the Patients page (See Figure 2), Max was able to directly navigate to the patient John Doe.

When clicking on the patient John Doe, the patient overview page appeared (See Figure 3). Here are a couple of quotes that our testing person made:

1. *“I see browse mutations, and a lot of vertical red lines. I suppose these lines represent the mutations that are found in the patient”.*
2. After hovering over these mutations, as Max called them (they represent chromosomes), a new block with the title Mutations in chromosome X appeared, it became clear that Max was wrong. *“Ah, I see, these vertical lines represent chromosomes. Not that obvious since the header says: Browse mutations”.*

Initially, our test person was a bit confused by all the information that appeared. What is a rsID, what is a frequency and what is the CADD score? After some explanation (we basically told that the CADD score is an indication about how harmful a mutation could be), Max immediately sorted the column with mutations on the CADD score, and clicked on the mutation with the highest value. Our tester saw the mutation view (See Figure 4 and Figure 5). Here are a couple of useful quotes that the tester made:

3. *“I see a big vertical red line with a id on it, and a couple of horizontal lines next to it. I can imagine that the red line symbolizes the mutation, and that the black lines symbolize the proteins around it”.*
4. *“So I guess that this central green bulb symbolizes the protein that has a mutation, but I’m not really sure how I can get more information to confirm that this is the mutation that causes the Marfan syndrome”.*

Our tester was clearly not going to come further, we explained that the green bulb was clickable and useful information appeared. We stopped the experiment and started with the debriefing.

4.3 Debriefing

The first thing we asked Max was: “So what did you think of our application?” He was quite positive, and found most parts were quite intuitive. The main problem in our experiment was that Max wasn’t exactly the end user. A couple of terms (like CADD score) were terms that he had never heard of. We asked if there were any observations that Max had made, and if he wanted to share them with us. Max answered that he didn’t have a lot, he did think that the navigation in a mutation view could be more intuitive.

4.4 Conclusion

The numbered quotes that Max made and are included in this report are all quotes that we used in making our product better. Quote number one for example made us realise that the heading “Browse Mutations” is inappropriate when it is placed above 23 pairs of chromosomes, this was confusing and thus we renamed it.

Quote number three was a confirmation that our mutation overview is quite intuitive, although quote number four still gave us some room for improvement: we should make it clear that the round Proteins are clickable. We improved this by implementing a tooltip that pops up when hovering over the protein.

5 Evaluation of the functional models

This section describes the evaluation of the functional modules and how they are tested to guarantee that each module performs as required. Furthermore, the failure analysis is discussed to show where the product does not perform as needed.

5.1 Evaluation of functional models

All of our functionality that is implemented, is tested using the test suite of the play framework. Each model is simply tested with JUnit, where every bit of functionality is checked for functioning correctly.

The controllers and views are tested by running a fake play application, which starts an application in the background so we can call the controllers and view. We check if each request to the controllers results in the expected action, such as a "Ok", "Page not found", "Redirect" or "Bad request". Next we check the content of the rendered page, this should match the information we sent to this page. For the views only the content is checked.

5.2 Failure analysis

We rely a lot on data coming from other systems. First of all we have to read VCF-files to load the data of a patient. This is done by a doctor which could load a corrupted file to our server. This causes the system to stop processing the current file and the error will be written to the log. At the moment there is no notification to warn the user of this error, however this could be simply added in a future release.

Second we need access to external databases to retrieve relevant data. As these databases are too large to host them on the same webserver as the application, we cannot expect to always have a connection to the remote server. However we also store basic information in a small database on the webserver which stores the basic information such as authentication, user data and patient data. This helps us to always have a basic running system. If for some reason the connection to the database is dropped the application will continue to function, but no new information is loaded. In the future, useful notifications could be shown to the user to indicate such problems.

Thirdly, a user could direct his web browser to an invalid URL, in this case a "Not found" page will show, to prevent retrieving private information from the web-application. Also when something goes wrong in the application, the user will be sent to an "Error" page, which will only show a basic error and not the full error trace, which could expose our application logic.

6 Outlook

GEVATT is developed in quite a short time. The course Contextproject lasts about ten weeks and in this ten weeks the team had to become skilled in the context, make accompanying documents and present the final product. So only a limited part of the time is actually used to develop the application (week 4 to 9). This means the team focusses on the main parts of the application and there

is more functionality that could be added later if development were to continue. This chapter is about what functionality could be added.

One thing is exporting; this is currently not possible in the application. The found mutations and information from the VCF files are only visible in the application itself and require logging in. For a doctor, it might be handy to export some data and visualisations to make it easier to present them. E.g. it's easier for a doctor to show a print of visualisations to his patients than showing a screenshot of the application. It also might be handy to have the data exported to a spreadsheet format.

The application is web-based and therefor principally platform independent: each device with a modern browser and an internet connection could open the application. Something that is limited to the use of the application now is the minimum screen width. When a device with a resolution of lower than 1024px is used, the user need to scroll to see all information. The application could be modified so that the layout adopts to smaller screens and become thereby on lower resolution tablets and phones. A nice feature, but it is questionable whether a doctor would use the application on such devices.

Currently the application is mainly focused on individual mutations found in the uploaded data, and less on the relations between these mutations. If development of this application were to continue, it will probably make sense to focus on implementing/extending visualisations to show these relations.

Testing is only done on a low scale. If GEVATT would be used in production, there should be some testing be done based on heavy use. Furthermore, if the application would be used in production, some managing features should be added. For example, creating accounts en revoking access to specific doctors. Due to time considerations this is left out in this project.

7 Appendix A: Screenshots

Here are a couple of screenshots of our application.

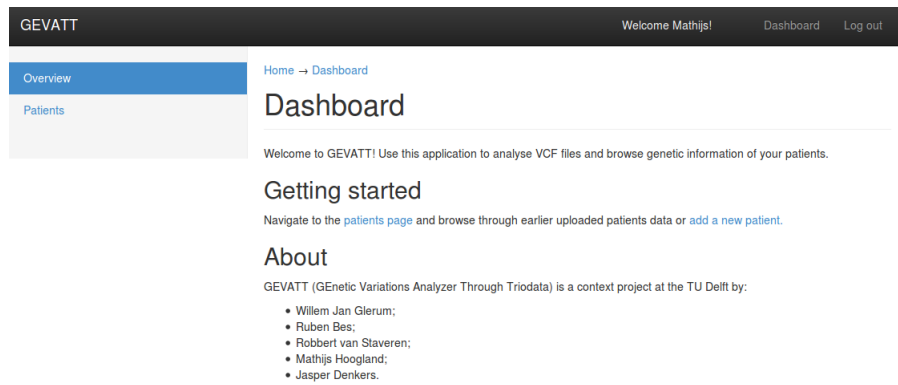


Figure 1: A screenshot of the dashboard

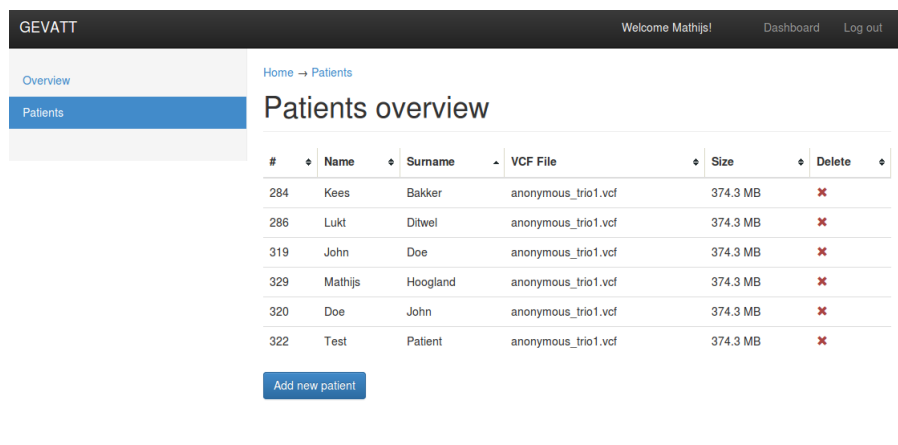


Figure 2: A screenshot of the patients view

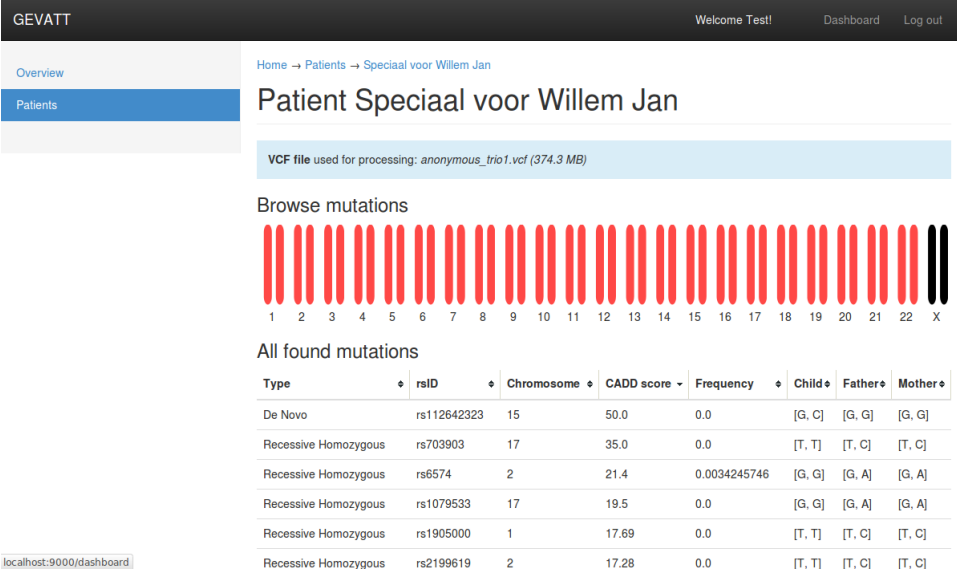


Figure 3: A screenshot of the patient view



Figure 4: A screenshot of the mutation view

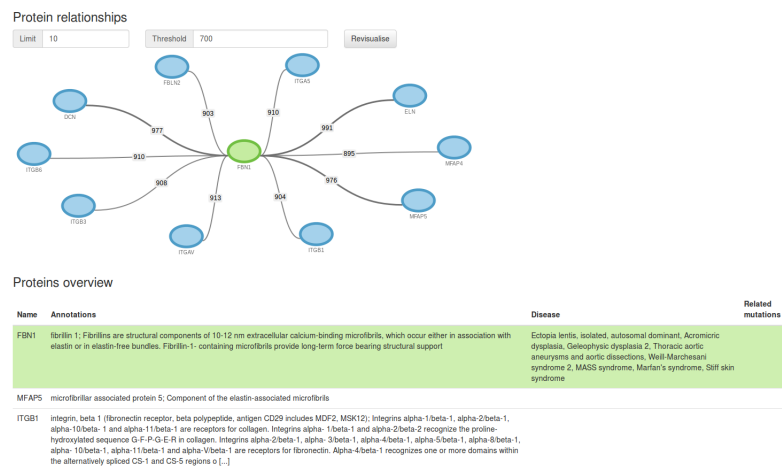


Figure 5: A second screenshot of the mutation view

References

- [1] Martin Kircher, Daniela M Witten, Preti Jain, Brian J O’Roak, Gregory M Cooper, and Jay Shendure. A general framework for estimating the relative pathogenicity of human genetic variants. *Nature genetics*, 46(3):310–315, 2014.
- [2] Stephen T Sherry, M-H Ward, M Kholodov, J Baker, Lon Phan, Elizabeth M Smigielski, and Karl Sirotkin. dbsnp: the ncbi database of genetic variation. *Nucleic acids research*, 29(1):308–311, 2001.
- [3] Andrea Franceschini, Damian Szklarczyk, Sune Frankild, Michael Kuhn, Milan Simonovic, Alexander Roth, Jianyi Lin, Pablo Minguéz, Peer Bork, Christian von Mering, et al. String v9. 1: protein-protein interaction networks, with increased coverage and integration. *Nucleic acids research*, 41(D1):D808–D815, 2013.
- [4] Melissa J Landrum, Jennifer M Lee, George R Riley, Wonhee Jang, Wendy S Rubinstein, Deanna M Church, and Donna R Maglott. Clinvar: public archive of relationships among sequence variation and human phenotype. *Nucleic acids research*, page gkt1113, 2013.