Lehrstuhl fuer Systeme und Informationsverwaltung

Prof.Dr.Ing. Klemens Boehm Hoang Vu Nguyen Marco Neumann



Praxis der Softwareentwicklung (PSE) WS 2012/2013

Visualizing and Statistically Analyzing Access Behavior to Scientific Databases

Functional Specification



November 19, 2012

Functional Specification
Design
Implementation
Alexander Noe
Jonathan Klawi
Anas Saber
Nikologo Alexa

Design Jonathan Klawitter
lementation Anas Saber
QA / Testing Nikolaos Alexandros Kurt Moraitakis
Final Lukas Ehnle

About this document

General Information

PSE is a mandatory course for students of Bachelor Informatik in the Karlsruhe Institute of Technology (KIT). Therefor groups of five to six students are formed to write programs of 'medium size' (about 5000 Lines of Code). The task given in this course consists of analyzing and visualizing the server log of a database.

This specification is the one of PSE Group #10 in the winter semester 2012/13.

What this specification does

The purpose of this document is a outline of the functional specifications and requirements for the WHAT application. Is tries to give a complete and exact model of the future system. In the course of this it wants to give the developers answers to all possible question concerning what should be implemented.

What this specification does not

This specification does not give any information on how the system should be implemented. Also it doesn't contain any project planning or deadlines.

A 'Living Document'

Software development is a dynamic process. So with proceeding development the content of this specification will change continuously.

Skyserver

The concept of the system described in this specification should work with any database storing it's query-log. The SkyServer description will function as example and testing reference for this system. SkyServer is one of the biggest databases for astronomical data.

About us

Because no one of us speaks English as his first language, we can't guarantee that our documentation doesn't contain simple or incorrectly used language. Our focus lies more on having easy-to-understand and correct documentation than on perfect English by a native speaker.

If you have any questions or comments regarding this document feel free to send us an E-Mail.

Contents

| 1 | Goa | |
|---|-----|------------------------------|
| | 1.1 | Web page |
| | | 1.1.1 Core criteria |
| | | 1.1.2 Optional criteria |
| | | 1.1.3 Exclusion criteria |
| | 1.2 | Parser |
| | | 1.2.1 Core criteria |
| | | 1.2.2 Exclusion criteria |
| | 1.3 | Analyzer |
| | | 1.3.1 Core criteria |
| | | 1.3.2 Optional criteria |
| 2 | Usa | ige 7 |
| | | Area of Application |
| | 2.2 | Target groups / Audience |
| | 2.3 | Operating conditions |
| | | |
| 3 | | erating environment |
| | | Software |
| | | Hardware |
| | | Orgware |
| | 3.4 | Product interfaces |
| 4 | Fur | ctional requirements |
| | | Main functions |
| | | Extending functions |
| | | |
| 5 | Dat | |
| | | Static data |
| | 5.2 | Data-Warehouse data |
| 6 | Noi | nfunctional requirements |
| _ | | First Whatever |
| | | Next Whatever |
| | | Optional |
| | | |
| 7 | | t Cases 14 |
| | 7.1 | Global test cases |
| | | 7.1.1 General |
| | | 7.1.2 Administrator specific |
| | | 7.1.3 Parser specific |
| | | 7.1.4 Analyzer specific |

| 8 | Models | 16 | | | |
|----|-------------------------|----|--|--|--|
| | 8.1 Overview | | | | |
| | 8.1.1 Parser | | | | |
| | 8.1.2 Analyzer | 17 | | | |
| | 8.1.3 Web page | | | | |
| | 8.2 Dynamic models | 18 | | | |
| | 8.3 Web interfaces | 19 | | | |
| 9 | Development Environment | 20 | | | |
| 10 | 10Glossarv 2 | | | | |

1 Goals

With this program the user should be put in the position to visualize the prepared data of queries, made against his data base.

To structure the criteria the system has to fulfill, it is divided into three parts for this and other sections.

- Web page
- Parser
- Analyzer

See 8.1 for a overview of this parts and their relationships. In the following their specific goals and criteria are described.

1.1 Web page

1.1.1 Core criteria

- The web page is the graphical interface for users and administrators.
- Users can choose the charts they want to see and also choose and filter the dimensions and measures for those. It has to support at least scatter plots, histrograms and bubble charts.
- The web page provides help for how to use the diagrams and variables (dimensions and measures). Also the navigation on the page is easy to use and straight forward.
- A admin-login provides administrative rights, which are needed to use the parser.

1.1.2 Optional criteria

- The language of the website can be changed into different languages. (e.g. German)
- The administrator gets the ability to handle the data warehouse and load new log-files.
- A little history of the last requested charts is stored and viewable.

1.1.3 Exclusion criteria

- The web page does not allow normal users to load new data into the data warehouse.
- The web page does not provide statistic tabular.

1.2 Parser

1.2.1 Core criteria

- The Parser is able to perform the ETL-process on CSV-formatted log files (from Skyserver). This means he extracts the data he needs from the files, transforms them and loads them into the data-warehouse. Dimensions and measures are specified in 5.2.
- The Parser will recognize invalid logs and won't add them to the data-warehouse. Every
 log with a mistake won't be accepted, because an error-message is not as bad as a
 corrupted data-warehouse.
- The Parser will be fed with log files from the administrator via web page.

1.2.2 Exclusion criteria

- This parser is only able to operate on CSV-formatted logs from SkyServer. It can neither read logs in another formats nor logs from another source.
- There is no way to avoid using this Parser when adding data to the warehouse. This can stop corrupting the warehouse to guarantee correct data in the warehouse.
- The Parser doesn't correct mistakes in the log file.

1.3 Analyzer

1.3.1 Core criteria

- The analyzer is the gate to the data warehouse. It extracts the specific data, needed for the diagrams, passing them to the java-script front-end.
- The analyzer can take filtered information via web page from the user to use only certain data for the charts.
- The charts that are supported are at least:
 - scatter plots
 - histograms
 - bubble charts

1.3.2 Optional criteria

- The analyzer will support more different chart-types, especially the combination of a histogram and a scatter plot.
- The analyzer does a little bit of data mining, and presents some potentially interesting information. This function will be called Niko's data mining.

2 Usage

2.1 Area of Application

The application area of this product are scientific databases that need to be analysed, understood, and visualised. The program has been designed with the Skyserver in mind, so that is the primary application area. However, as the only major requirement is the existence of correct access logs of the database, the program may also be used with many other scientific databases, although minor modifications may be needed.

2.2 Target groups / Audience

The target group of this application is people that want to analyse and visualize queries made against a database. Specifically in our case, the Skyserver.

This includes

- the owners of the database who want to optimize the access to their data,
- people that are interested in what others use the database for,
- people new to the database, who want to understand it or discover interesting information.

To summarize, WHAT will be of interest for many people, whether they want to see how their database is used, analyse current trends, or just love statistics and diagrams.

This is a rather technical audience.

That said, this does not prevent the (web) user interface from being functional, usable and prettier than what you would expect a group of computer science students to design.

2.3 Operating conditions

The program is mainly used as a website, with the primary difference being that the server has to be started if the capacity for it to run all the time on a dedicated machine doesn't exist. The program needs a server to run.

If a dedicated server exists, the program can be used from anywhere with a decent network connection with the server.

If not, the program can still be run on the same computer as the server (on localhost), but the server will have to be started first.

3 Operating environment

Whereas the program and the data warehouse run on a server, the access to it will be via a web browser on the users computer. This implies that the webpage will need to be hosted on a webserver. However, a dedicated server machine is not required, as the clients computer can host both the client and the webserver, and then access the webpage locally.

3.1 Software

- The server hosting the data warehouse needs MySQL
- The server hosting the web server needs a recent version of the Java Runtime Environment / Java VM

Following software is required on the users computer:

- Latest version of a modern browser
 - Though Chrome and Firefox will be officially supported, other browsers might work as well
- JavaScript (enabled)

3.2 Hardware

The server has to be fast enough to support all clients. This depends on the expected number of clients. Most computations will be done on the server.

The client needs to be fast enough to visualize the data received from the server, as that's almost all it does. The required hardware thus depends greatly on the amount of data that needs to be visualized. That said, any recent computer, with for example, an Intel® Core $^{\text{TM}}$ 2 Duo CPU E8400 @ 3.00GHz \times 2 processor, 4GB(2x2GB) of 667Mhz DDR2 SDRAM, running Ubuntu Linux 12.10 Quantal Quetzal, as a point of reference, should be able to visualize about 50.000 data points on a scatterplot with ease.

3.3 Orgware

The server needs to be able to connect to the client with a reasonable latency. Also it has to be set up before the client is able to connect to it.

To fulfill optional requirements, the program musst be able to run queries against the origin database. //FIXME I DON'T KNOW ABOUT THAT

3.4 Product interfaces

To fulfill optional requirements .csv files have to be imported.

4 Functional requirements

4.1 Main functions

This functions are required to fulfill the core criteria.

General

- /F10/ Provide access via web page
- /F20/ Provide option for chart types
- /F30/ Show diagrams
- /F40/ Show histograms
- /F50/ Provide option to select and filter variables for histograms
- /F60/ Show scatter plots
- /F70/ Provide option to select and filter variables for scatter plots
- /F80/ Show bubble charts
- /F90/ Provide option to select and filter variables for bubble charts
- /F100/ Show information about chart types
- /F110/ Show information about selectable variables
- /F120/ Provide easy navigation on the web page

Administrator specific

This functions are required for the administrative business.

- /F130/ Provide access via web page with administrator rights
- /F140/ Provide the opportunity to pass log-files to the parser

Parser specific

The following functions specify the parses functionality. Thereby /F160/ to /F180/ specify the function /F150/.

- /F150/ Extract specific data from the log-files
- /F160/ Extract access database, access time and user information (dimensions)
- /F170/ Extract number of rows, elapsed time, busy time (measures)
- /F180/ Extract type of data requested from the where-part

- /F190/ Transform the data to fit into the data warehouse schema
- /F200/ Load the data into the data warehouse

Analyzer specific

- /F210/ Run specific queries against data warehouse
- /F220/ Transform received data serving the web page

4.2 Extending functions

To fulfill the optional goals the following functions are required.

General

- /F230/ Select language on web page
- /F240/ Show bubble map
- /F250/ Show combination of histogram and scattered plot
- /F260/ Show other diagrams and charts
- /F270/ Show history of the 5 last requested charts
- /F280/ Show interesting attribute combinations
- /F290/ Show clusters in the data
- /F300/ Add background clustering
- /F310/ Turn background clustering on and off //to make the program faster

Administrator specific

- /F320/ Provide the opportunity to initialize the data warehouse
- /F330/ Provide the opportunity to request new log-files for specific time intervals from the database
- /F340/ Provide the opportunity to clean the data warehouse

5 Data

5.1 Static data

/D350/ Language files

/D360/ Manual

/D370/ Source Code

/D380/ Documentation

/D390/ Graphics for GUI

/D400/ HTML backbone

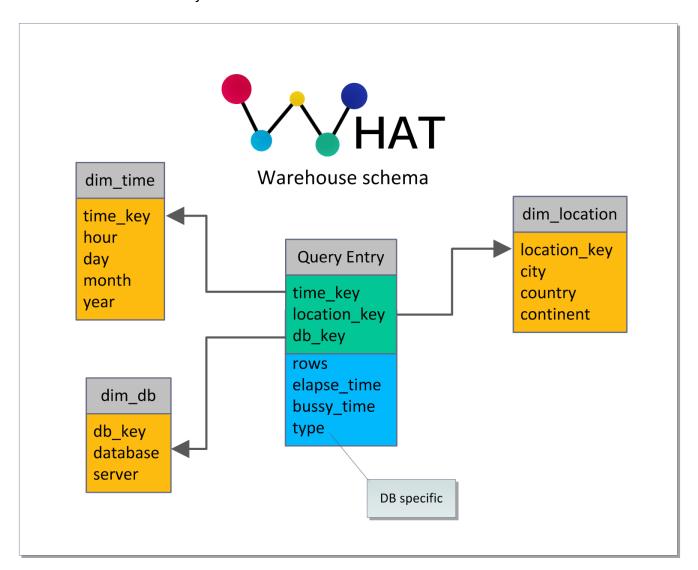
/D410/ Javascript files

/D420/ Stylesheet files

5.2 Data-Warehouse data

The data warehouse stores all the data, coming from from the database (Skyserver) server logs. They were subject of the ETL-process in the parser. If optional criteria are fulfilled they can be replaced by clearing the data warehouse and loading new data via web page as administrator.

The data warehouse may use a star schema.



See rows, elapse time and bussy time in the glossary (10) for descriptions.

6 Nonfunctional requirements

6.1 First Whatever

/NF10/ Parse 1000 log rows in few seconds

/NF20/ the data warehouse can store more than 1000000 entries

/NF30/ Visualize scatterplot of 50000 points in a few seconds

/NF40/ minimalistic UI design

/NF50/ responsive UI

/NF60/ UI that is easy to get used to

/NF70/ easy translating in other languages

/NF80/ reasonably commented, clean code

6.2 Next Whatever

/NF90/ 30!

/NF100/ 40!

6.3 Optional

/NF10/ Show a nice loading screen while visualizing a diagram.

7 Test Cases

7.1 Global test cases

7.1.1 General

Those test-cases are necessary for the webpage to work properly and succeed in serving basic requests.

/T10/ Access web page (local or internet) (/F10/)

/T12/ via Google Chrome

/T14/ via Mozilla Firefox

- /T20/ Navigate switfly through the web page. (/F120/)
- /T30/ Open page for scatter plot-creation and create an example with x-Axis date and y-Axis number of lines per request. (/F20/, /F30/, /F60/, /F70/)
- /T40/ Open page for histogram-creation and create an example with x-Axis year and y-Axis number of requests. (/F20/, /F30/, /F40/, /F50/)
- /T50/ Open page for bubble chart-creation and create an example with x-Axis year, y-Axis Country and size number of requests. (/F20/, /F30/, /F80/, /F90/)
- /T60/ Call the information about chart types. (/F100/)
- /T70/ Request information about selectable variables. (/F110/)

7.1.2 Administrator specific

Those test cases check the correct admin-login.

- /T80/ Enter loginname and password for admin-login. (/F130/)
- /T90/ Log-in and get administrative rights. (/F130/)
- /T100/ Try to log in with wrong loginname. (/F130/)
- /T110/ Try to log in with wrong password. (/F130/)
- /T120/ Pass example log-file to the parser. (/F140/)

7.1.3 Parser specific

The following test cases test the correct work of our parser.

/T130/ Parse example log-file and extract valuable information correctly. (/F150/) More specifically:

/T132/ Extract correct database, time, user information and type (/F160/)

/T134/ Extract correct number of rows, elapsed time, busy time from the log-file (/F170/)

/T136/ Extract the type requested (/F180/)

/T140/ Load extracted information into the data warehouse correctly. (/F190/, /F200/)

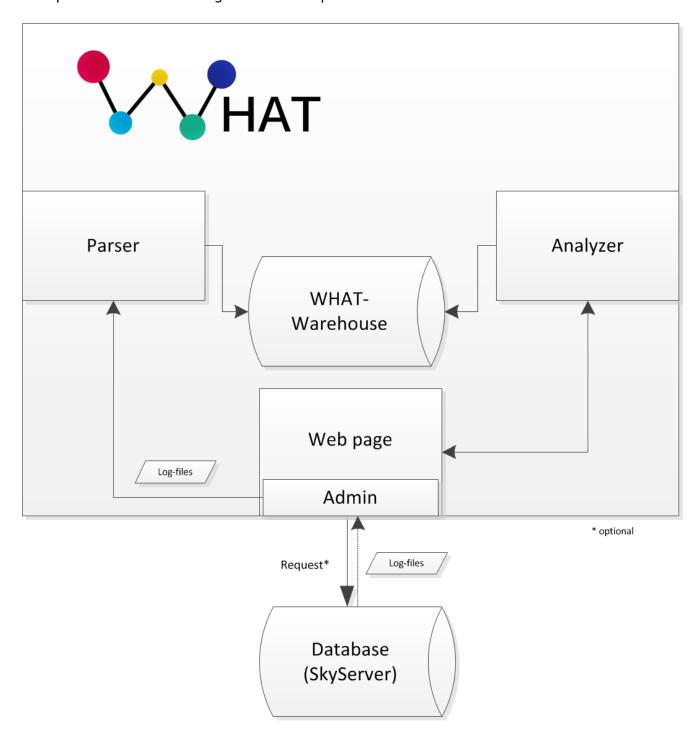
7.1.4 Analyzer specific

The functions /F210/ and /F220/ of the analyzer are already tested in the general test-cases /T30/ to /T50/.

8 Models

8.1 Overview

This picture describes the generell concept of WHAT.



8.1.1 Parser

Parser blabla

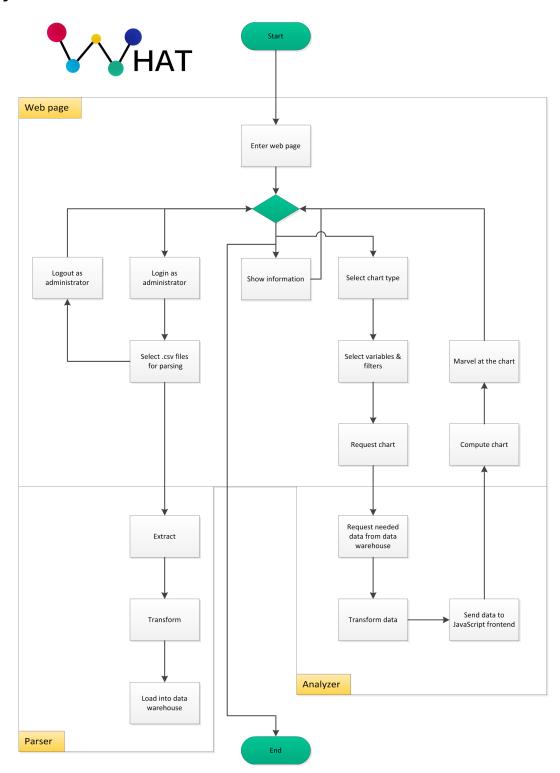
8.1.2 Analyzer

analyzer blabla

8.1.3 Web page

web page blabla

8.2 Dynamic models



8.3 Web interfaces

9 Development Environment

Operating System

- Windows 7
- Fedora 17
- Whatever Niko uses EDIT THIS, NIKO!

Object Models

• Microsoft Visio

Version Control

• Git

Miscellaneous

- Github (Hosting of Repository, Issue Tracking)
- Travis (Integration)
- LATEX (Documents)
- Gradle (Building)

10 Glossary

.. ..

KIT is a technological university in Karlsruhe, Baden-Württemberg, Germany (Karlsruhe Institute of Technology)

SkyServer is a huge database for astronomical data. It contains pictures and other information concerning astronomy and tries to form a 'map of the universe'. Web page ofSkyServer SkyServer.