



# The handbook of gStore System v0.7.0

Edited by gStore team <sup>1</sup>

November 5, 2018

<sup>1</sup>The mailing list is given in Chapter 13.

# Contents

<b>Preface</b>	<b>7</b>
<b>I Start</b>	<b>9</b>
Chapter 00: A Quick Tour . . . . .	9
Getting Started . . . . .	9
Compile from Source . . . . .	9
Deploy via Docker . . . . .	10
Run . . . . .	10
Advanced Help . . . . .	10
Other Business . . . . .	11
Chapter 01: System Requirements . . . . .	13
Chapter 02: Basic Introduction . . . . .	18
What Is gStore . . . . .	18
Why gStore . . . . .	18
Open Source . . . . .	19
Chapter 03: Install Guide . . . . .	20
Install from Source . . . . .	20
Docker Deployment . . . . .	21
Prepare the Environment . . . . .	21
Build the Mirror via Dockerfile . . . . .	21
Pulling the Mirror Directly to Run . . . . .	21
Problems That Exist . . . . .	22
Follow-up . . . . .	22
Chapter 04: How To Use . . . . .	23
0. Format of data . . . . .	23

1. gbuild . . . . .	23
2. gquery . . . . .	23
3. ghttp . . . . .	25
4. gserver . . . . .	31
5. gclient . . . . .	31
6. gconsole . . . . .	32
7. gadd . . . . .	33
8. gsub . . . . .	33
9. gmonitor . . . . .	34
10. gshow . . . . .	34
11. shutdown . . . . .	34
12. ginit . . . . .	35
13. test utilities . . . . .	35

## II Advanced 37

Chapter 05: HTTP API Explanation . . . . .	37
Easy Examples . . . . .	37
API Structure . . . . .	38
C++ API . . . . .	39
Interface . . . . .	39
Java API . . . . .	41
Interface . . . . .	41
Python API . . . . .	42
Interface . . . . .	42
Php API . . . . .	43
Interface . . . . .	44
Nodejs API . . . . .	45

Interface . . . . .	45
Chapter 06: Socket API Explanation . . . . .	47
Easy Examples . . . . .	47
API structure . . . . .	48
C++ API . . . . .	49
Interface . . . . .	49
Compile . . . . .	51
Java API . . . . .	51
Interface . . . . .	51
Compile . . . . .	53
PHP API . . . . .	53
Interface . . . . .	53
Run . . . . .	55
Python API . . . . .	55
Interface . . . . .	55
Run . . . . .	57
Chapter 07: Use gStore in Web . . . . .	58
Example . . . . .	58
Chapter 08: Project Structure . . . . .	63
The core source codes are listed below: . . . . .	63
The parser part is listed below: . . . . .	66
The utilities are listed below: . . . . .	66
The interface part is listed below: . . . . .	67
More details . . . . .	68
Others . . . . .	68
Chapter 09: Publications . . . . .	73
Publications related with gStore are listed here: . . . . .	73

Chapter 10: Limitations . . . . .	74
Chapter 11: Frequently Asked Questions . . . . .	75
When I use the newer gStore system to query the original database, why error? . . . . .	75
Why error when I try to write programs based on gStore, just like the Main/gconsole.cpp? . . . . .	75
Why does gStore report “garbage collection failed” error when I use the Java API? . . . . .	75
When I compile the code in ArchLinux, why the error that “no -ltermcap” is reported? . . . . .	75
Why does gStore report errors that the format of some RDF datasets are not supported? . . . . .	75
When I read on GitHub, why are some documents unable to be opened? . . . . .	76
Why sometimes strange characters appear when I use gStore?	76
In centos7, if the watdiv.db(a generated database after gbuild) is copied or compressed/uncompressed, the size of watdiv.db will be different(generally increasing) if using <code>du -h</code> command to check?	76
In gclient console, a database is built, queried, and then I quit the console. Next time I enter the console, load the originally imported database, but no output for any queries(originally the output is not empty)? . . . . .	77
If query results contain null value, how can I use the full_test utility? Tab separated method will cause problem here because null value cannot be checked! . . . . .	77

When I compile and run the API examples, it reports the “unable to connect to server” error? . . . . .	77
When I use the Java API to write my own program, it reports “not found main class” error? . . . . .	77
Chapter 12: Recipe Book . . . . .	78
Config . . . . .	78
Backup . . . . .	78
Query . . . . .	78
KVstore . . . . .	78
String Buffer . . . . .	79
HTTP API . . . . .	79
<b>III Others</b>	<b>80</b>
Chapter 13: Contributors . . . . .	80
Faculty . . . . .	80
Students . . . . .	80
Alumni . . . . .	81
Chapter 14: Updated Logs . . . . .	82
Apr 24, 2018 . . . . .	82
Oct 2, 2017 . . . . .	82
Jan 10, 2017 . . . . .	83
Sep 15, 2016 . . . . .	83
Jun 20, 2016 . . . . .	83
Apr 01, 2016 . . . . .	84
Nov 06, 2015 . . . . .	84
Oct 20, 2015 . . . . .	85
Sep 25, 2015 . . . . .	85

Feb 2, 2015 . . . . .	85
Dec 11, 2014 . . . . .	86
Nov 20, 2014 . . . . .	86
Chapter 15: Test Result . . . . .	87
Preparation . . . . .	87
Result . . . . .	88
Chapter 16: Future Plan . . . . .	90
Improve The Core . . . . .	90
Better The Interface . . . . .	90
Idea Collection Box . . . . .	90
Chapter 17: Thanks List . . . . .	91
zhangxiaoyang . . . . .	91
Dingfeng Wang . . . . .	91
Libo Wang . . . . .	91
Xin Lv . . . . .	91
Zhiyuan Deng . . . . .	92
Hao Cui . . . . .	92
imbajin . . . . .	92
Chapter 18: Legal Issues . . . . .	93
<b>End</b>	<b>95</b>

## Preface

The RDF (*Resource Description Framework*) is a family of specifications proposed by W3C for modeling Web objects as part of developing the semantic web. In RDF model, each Web object is modeled as a uniquely named *resource* and denoted by a URI (*Uniform Resource Identifier*). RDF also uses URIs to name the properties of resources and the relationships between resources as well as the two ends of the link (this is usually referred to as a “triple”). Hence, an RDF dataset can be represented as a directed, labeled graph where resources are vertices, and triples are edges with property or relationship names as edge labels. For more details, please go to [RDF Introduction](#)

To retrieve and manipulate an RDF graph, W3C also proposes a structured query language, SPARQL (*Simple Protocol And RDF Query Language*), to access RDF repository. SPARQL contains capabilities for querying required and optional graph patterns along with their conjunctions and disjunctions. SPARQL also supports aggregation, subqueries, negation, creating values by expressions, extensible value testing, and constraining queries by source RDF graph. Similar to RDF graphs, a SPARQL query can also be modeled as a graph, which is a query graph with some variables. Then, evaluating a SPARQL query is equivalent to finding subgraph (homomorphism) matches of a query graph over an RDF graph. You can have a better understanding of SPARQL at [SPARQL Introduction](#).

Although there are some RDF data management systems (like Jena, Virtuoso, Sesame) that store the RDF data in relational systems, few existing systems exploit the native graph pattern matching semantics of SPARQL. **Here, we implement a graph-based RDF triple store named gStore, which is a joint research project by Peking University, University of Waterloo and Hong Kong University of Science and Technology. The system is developed and maintained by the database group in Institute of Computer Science and Technology, Peking University, China.** A detailed description of gStore can be found at our papers



[Zou et al., VLDB 11] and [Zou et al., VLDB Journal 14] in the **Publication** chapter. This HELP document includes system installment, usage, API, use cases and FAQ. gStore is a open-source project in github under the BSD license. You are welcome to use gStore, report bugs or suggestions, or join us to make gStore better. It is also allowed for you to build all kinds of applications based on gStore, while respecting our work.

**Please make sure that you have read **Legal Issues** before using gStore.**

## Part I

# Start

### Chapter 00: A Quick Tour

Gstore System(also called gStore) is a graph database engine for managing large graph-structured data, which is open-source and targets at Linux operation systems. The whole project is written in C++, with the help of some libraries such as readline, antlr, and so on.

#### Getting Started

**Compile from Source** This system is really user-friendly and you can pick it up in several minutes. Remember to check your platform where you want to run this system by viewing [System Requirements](#). After all are verified, please get this project's source code. There are several ways to do this:

- type `git clone https://github.com/pkumod/gStore.git` in your terminal or use git GUI to acquire it
- download the zip from this repository and extract it
- fork this repository in your github account

Then you need to compile the project, for the first time you need to type `make pre` to prepare the ANTLR library and some Lexer/Parser programs. Later you do not need to type this command again, just use the `make` command in the home directory of gStore, then all executables will be generated. (For faster compiling speed, use `make -j4` instead, using how many threads is up to your machine.) To check the correctness of the program, please type `make test` command.

The first strategy is suggested to get the source code because you can easily acquire the updates of the code by typing `git pull` in the home directory of gStore

repository. In addition, you can directly check the version of the code by typing `git log` to see the commit logs. If you want to use code from other branches instead of master branch, like 'dev' branch, then:

- clone the master branch and type `git checkout dev` in your terminal
- clone the dev branch directly by typing `git clone -b dev`

**Deploy via Docker** You can easily deploy gStore via Docker. We provide both of Dockerfile and docker image. Please see [Docker Deployment](#).

**Run** To run gStore, please type `bin/gbuild database_name dataset_path` to build a database named by yourself. And you can use `bin/gquery database_name` command to query a existing database. What is more, `bin/ghttp` is a wonderful tool designed for you, as a database server which can be accessed via HTTP protocol. Notice that all commands should be typed in the root directory of gStore.

*A detailed description can be found at Chapter 04 [How to use](#) in this document.*

## Advanced Help

If you want to understand the details of the gStore system, or you want to try some advanced operations(for example, using the API, server/client), please see the chapters below.

- [Basic Introduction](#): introduce the theory and features of gStore
- [Install Guide](#): instructions on how to install this system
- [How To Use](#): detailed information about using the gStore system
- [HTTP API Explanation](#): guide you to develop applications based on our HTTP API

- **Socket API Explanation**: guide you to develop applications based on our Socket API
- **Project Structure**: show the whole structure and sequence of this project
- **Publications**: contain essays and publications related with gStore
- **Update Logs**: keep the logs of the system updates
- **Test Results**: present the test results of a series of experiments

## Other Business

We have written a series of short essays addressing recurring challenges in using gStore to realize applications, which are placed in **Recipe Book**.

You are welcome to report any advice or errors in the github Issues part of this repository, if not requiring in-time reply. However, if you want to urgent on us to deal with your reports, please email bookug@qq.com to submit your suggestions and report bugs to us by emailing to gStoreDB@gmail.com. A full list of our whole team is in **Contributors**.

There are some restrictions when you use the current gStore project, you can see them on **Limitations**.

Sometimes you may find some strange phenomena(but not wrong case), or something hard to understand/solve(don't know how to do next), then do not hesitate to visit the **Frequently Asked Questions** page.

Graph database engine is a new area and we are still trying to go further. Things we plan to do next is in **Future Plan** chapter, and we hope more and more people will support or even join us. You can support in many ways:

- watch/star our project
- fork this repository and submit pull requests to us
- download and use this system, report bugs or suggestions

- ...

People who inspire us or contribute to this project will be listed in the [Thanks List](#) chapter.

## Chapter 01: System Requirements

*We have tested on linux server with CentOS 6.2 x86\_64 and CentOS 6.6 x86\_64.*

*The version of GCC should be 4.4.7 or later.*

Item	Requirement
operation system	Linux, such as CentOS, Ubuntu and so on
architecture	x86_64
disk size	according to size of dataset
memory size	according to size of dataset
glibc	version $\geq$ 2.14
gcc	version $\geq$ 4.8
g++	version $\geq$ 4.8
make	need to be installed
boost	version $\geq$ 1.54
readline	need to be installed
readline-devel	need to be installed
openjdk	needed if using Java api
openjdk-devel	needed if using Java api
requests	needed if using Python http api
pthread	needed if using php http api
curl-devel	needed if using php http api
node	needed if using Nodejs http api and version $\geq$ 10.9.0
realpath	needed if using gconsole
ccache	optional, used to speed up the compilation
libcurl-devel	needed to be installed

Table 1: software requirement

NOTICE:

**To help ease the burden of setting environments, several scripts are provided in `gstore/scripts/setup` for different Linux distributions. Please select**

**the setup scripts corresponding to your system and run it with root(or sudo) privilege. (As for CentOS system, you need to install boost-devel by yourselves.)**

1. The name of some packages may be different in different platforms, just install the corresponding one in your own operation system.
2. To install readline and readline-devel, just type `dnf install readline-devel` in Redhat/CentOS/Fedora, or `apt-get install libreadline-dev` in Debian/Ubuntu. Please use corresponding commands in other systems. If you use ArchLinux, just type `pacman -S readline` to install the readline and readline-devel.(so do other packages)
3. You do not have to install realpath to use gStore, but if you want to use the gconsole for its convenience, please do so by using `dnf install realpath` or `apt-get install realpath`.
4. Our programs use regEx functions, which are provided by GNU/Linux by default.
5. ANTLR3.4 is used in gStore to produce lexer and parser code for SPARQL query. However, you do not need to install the corresponding antlr libraries because we have merged the libantlr3.4 in our system.
6. When you type `make` in the root directory of the gStore project, the Java api will also be compiled. You can modify the makefile if you do not have JDK in your system. However, you are advised to install `openjdk-devel` in your Linux system.
7. To install ccache, you need to add epel repository if using CentOS, while in Ubuntu you can directly install it by `apt-get install ccache` comand. If you can not install ccahe(or maybe you do not want to), please go to modify the makefile(just change the CC variable to g++).

8. If you need to use the HTTP server in gStore, then Boost Library (like boost-devel, including boost headers for developing) must be installed and the version should not be less than 1.54. Remember to check the makefile for your installed path of Boost. To use Python api, you need to install requests by `pip install requests` in CentOS. To use php api, you need to install pthreads and curl in CentOS as follows:

```
1- get curl-devel
# yum install curl-devel

2- get php
# wget -c http://www.php.net/distributions/php-5.4.
  36.tar.gz

    get pthreads
# wget -c http://pecl.php.net/get/pthreads-1.0.0.tgz

3- extract both
# tar zxvf php-5.4.36.tar.gz
# tar zxvf pthreads-1.0.0.tgz

4- move pthreads to php/ext folder
# mv pthreads-1.0.0 php-5.4.36/ext/pthreads

5- reconfigure sources
# ./buildconf --force
# ./configure --help | grep pthreads
```

You have to see `--enable-pthreads` listed.

If do not, clear the builds with this commands:

```
# rm -rf aclocal.m4
```



```
# rm -rf autom4te.cache/  
# ./buildconf --force
```

6- Inside php folder run configure command  
to set what we need:

```
# ./configure --enable-debug --enable-maintainer-zts  
--enable-pthreads --prefix=/usr  
--with-config-file-path=/etc --with-curl
```

7- install php

We will run make clear just to be sure that no other  
crashed build will mess our new one.

```
# make clear  
# make  
# make install
```

8- Copy configuration file of PHP  
and add local lib to include path  
# cp php.ini-development /etc/php.ini

Edit php.ini and set Include\_path to be like this:

```
Include_path = "/usr/local/lib/php"
```

9- install node

```
# wget https://npm.taobao.org/mirrors/node/v10.9.0/node-v10.9.0.t  
# tar -xvf node-v10.9.0.tar.gz  
# cd node-v10.9.0  
# ./configure  
# make
```

```
# sudo make install
```

10- Check Modules

```
# php -m (check pthread loaded)
```

You have to see pthreads listed

11- If pthread is not listed, update php.ini

```
# echo "extension=pthreads.so" >> /etc/php.ini
```

9. Any other questions, please go to [FAQ](#) page.

## Chapter 02: Basic Introduction

*The first essay to come up with Gstore System is [gStore\\_VLDBJ](#), and you can find related publications in [Publications](#).*

### What Is gStore

gStore is a graph-based RDF data management system(or what is commonly called a “triple store”) that maintains the graph structure of the original [RDF](#) data. Its data model is a labeled, directed multi edge graph, where each vertex corresponds to a subject or an object.

We represent a given [SPARQL](#) query by a query graph Q. Query processing involves finding subgraph matches of Q over the RDF graph G, instead of joining tables in relational data management system. gStore incorporates an index over the RDF graph (called VS-tree) to speed up query processing. VS-tree is a height balanced tree with a number of associated pruning techniques to speed up subgraph matching.

**The gStore project is supported by the National Science Foundation of China (NSFC), Natural Sciences and Engineering Research Council (NSERC) of Canada, and Hong Kong RGC.**

### Why gStore

After a series of test, we analyse and keep the result in [Test Results](#). gStore runs faster to answer complicated queries(for example, contain circles) than other database systems. For simple queries, both gStore and other database systems work well.

In addition, now is the big data era and more and more structured data is coming, while the original relational database systems(or database systems based on relational tables) cannot deal with them efficiently. In contrast, gStore can utilize the features of graph data structures, and improve the performance.

What is more, gStore is a high-extensible project. Many new ideas of graph

database have be proposed, and most of them can be used in gStore. For example, our group is also designing a distributed gstore system.

### **Open Source**

The gStore source code is available as open-source code under the BSD license. You are welcome to use gStore, report bugs or suggestions, or join us to make gStore better. It is also allowed for you to build all kinds of applications based on gStore, while respecting our work.

## Chapter 03: Install Guide

### Install from Source

You are advised to read `init.conf` file, and modify it as you wish. (this file will configure the basic options of gStore system)

gStore is a green software, and you just need to compile it with three commands. Please run

```
$ sudo ./scripts/setup/setup_$(ARCH).sh
$ make pre
$ make
```

in the gStore home directory to prepare the dependency, link the ANTLR lib, compile the gStore code, and build executable “gbuild”, “gquery”, “ghttp”, “gserver”, “gclient”. (Please substitute the `$(ARCH)` with your system version, like `setup_archlinux.sh`, `setup_centos.sh` and `setup_ubuntu.sh`) What is more, the api of gStore is also built now.

Setup scripts and dependency preparation only need to be done once, later you can directly use `make` to compile the code.

(For faster compiling speed, use `make -j4` instead, using how many threads is up to your machine)

To check the correctness of the program, please type `make test` command.

Only if you use the `make dist` command, then you need to run `make pre` command again.

If you want to use API examples of gStore, please run `make APIexample` to compile example codes for both C++ API and Java API. For details of API, please visit [API](#) chapter.

Use `make clean` command to clean all objects, executables, and use `make dist` command to clean all objects, executables, libs, datasets, databases, debug logs, temp/text files in the gStore root directory.

You are free to modify the source code of gStore and create your own project while respecting our work, and type `make tarball` command to compress all useful files into a `.tar.gz` file, which is easy to carry.

Type `make gtest` to compile the gtest program if you want to use this test utility. You can see the [HOW TO USE](#) for details of gtest program.

## Docker Deployment

Roughly speaking, there are two ways to deploy gStore via Docker. The first one is using Dockerfile in the root directory of project to automatically build it. And then run the container. Another one is downloading the mirror which has been automatically built directly, then just run it.

**Prepare the Environment** Official doc of Docker has explained how to download and use it on common Linux release version in details. And here is the link: [Docker Docs](#).

It's worth noting that the Docker with too high version may lead to some problems. Please read the precautions carefully. The current version of test environment is Docker CE 17.06.1.

**Build the Mirror via Dockerfile** After having the correct Docker environment and network, use `git clone` to download the project firstly. After inputting command `docker build -t gstore`, it's available to start building. In the default case, it will use the Dockerfile in the root directory. More specific explanation has been written in the Dockerfile.

After the building, using `docker run -it gstore` directly to enter the container and execute other operations.

**Pulling the Mirror Directly to Run** Instead of downloading project or building on your own, input `docker pull suxunbin/auto_gstore:latest` to pull the mirror which has been automatically built well on the docker hub. Then

input `docker run -it suxunbin/auto_gstore:latest` to enter the container and execute other operations.

**Problems That Exist** Owing to the uncertain influence accompanying the containerization, including but not limited to some problems about network, lock, caching, rights and so on, it's quite difficult to locate the problem when debugging. It has been known that the following problems may exist: (The host is CentOS 7.4)

1. If add the environment `ENV CC="ccache g++"` in the process of Dockerfile building, it will lead to compile error. The reason is still unknown and it may affect the mirror cache layer, which results in repeated errors in the later process of building.
2. The case of shutdown automatically as soon as its reboot may occur after starting the container via `docker run`.

### **Follow-up**

**Performance Testing** The proportion of the loss of the performance of the container under the conditions of different file numbers/networks. The performance of running gStore in the native environment and in the container.

It's waiting for supplement.

**Test of Connecting Other Containers** Waiting for supplement.

**Simplification and Optimization of Building** The mirror of gcc:8 has conquered the space of 1.7G, bringing a lot of unnecessary things (including the environment of Go). Hoping to optimize the mirror source in the future excepting lowering the gcc's version.

## Chapter 04: How To Use

*gStore currently includes five executables and others.*

**All the commands of gStore should be used in the root directory of gStore like bin/ghttp, because executables are placed in bin/, and they may use some files whose paths are indicated in the code, not absolute paths. We will ensure that all paths are absolute later by asking users to give the absolute path in their own systems to really install/configure the gStore. However, you must do as we told now to avoid errors.**

**0. Format of data** The RDF data should be given in N-Triple format (XML is not supported by now) and queries must be given in SPARQL 1.1 syntax. Not all syntax in SPARQL 1.1 are parsed and answered in gStore, for example, property path is beyond the ability of gStore system. Tabs, '<' and '>' are not allowed to appear in entity, literal or predicates of the data and queries.

**1. gbuidl** gbuidl is used to build a new database from a RDF triple format file.

```
bin/gbuidl db_name rdf_triple_file_path
```

db\_name refers to the database name which you can setup as you wish, rdf\_triple\_file\_path stands for the path of data file . For example, we build a database from lubm.nt which can be found in data/lubm/ folder.

```
[bookug@localhost gStore]$ bin/gbuidl lubm ./data/lubm/
lubm.nt
gbuidl...
argc: 3 DB_store:lubm      RDF_data: ./lubm.nt
begin encode RDF from : ./data/lubm/lubm.nt ...
```

**2. gquery** gquery is used to query an existing database with files containing SPARQL queries.(each file contains exact one SPARQL query)



Type `bin/gquery db_name query_file` to execute the SPARQL query retrieved from `query_file` in the database named `db_name`.

Use `bin/gquery --help` for detail information of `gquery` usage.

To enter the `gquery` console, type `bin/gquery db_name`. The program shows a command prompt(`gsql>`), and you can type in a command here. Use `help` to see basic information of all commands, while `help command_t` shows details of a specified command.

Type `quit` to leave the `gquery` console.

For `sparql` command, input a file path which contains a single SPARQL query. (*answer redirecting to file is supported*)

When the program finish answering the query, it shows the command prompt again.

We also take `lumb.nt` as an example.

```
[bookug@localhost gStore]$ bin/gquery lubm
gquery...
argc: 2 DB_store:lubm/
loadTree...
LRUCache initial...
LRUCache initial finish
finish loadCache
finish loadEntityID2FileLineMap
open KVstore
finish load
finish loading
Type `help` for information of all commands
Type `help command_t` for detail of command_t
gsql>sparql ./data/lubm/lubm_q0.sql
... ..
Total time used: 4ms.
final result is :
<http://www.Department0.University0.edu/FullProfessor0>
```

```
<http://www.Department1.University0.edu/FullProfessor0>
<http://www.Department2.University0.edu/FullProfessor0>
<http://www.Department3.University0.edu/FullProfessor0>
<http://www.Department4.University0.edu/FullProfessor0>
<http://www.Department5.University0.edu/FullProfessor0>
<http://www.Department6.University0.edu/FullProfessor0>
<http://www.Department7.University0.edu/FullProfessor0>
<http://www.Department8.University0.edu/FullProfessor0>
<http://www.Department9.University0.edu/FullProfessor0>
<http://www.Department10.University0.edu/FullProfessor0>
<http://www.Department11.University0.edu/FullProfessor0>
<http://www.Department12.University0.edu/FullProfessor0>
<http://www.Department13.University0.edu/FullProfessor0>
<http://www.Department14.University0.edu/FullProfessor0>
```

Notice:

- “[empty result]” will be printed if no answer, and there is an empty line after all results.
- readline lib is used, so you can use arrow key in your keyboard to see command history, and use arrow key to move and modify your entire command.
- path completion is supported for utility. (not built-in command completion)

**3. ghttp** ghttp runs gStore like HTTP server with a specific port (You need to open this port in your environment, iptables tool is suggested). Visit from browser with prescriptive url, then gStore will execute corresponding operation.

Just type `bin/ghttp db_name serverPort` or `bin/ghttp serverPort db_name` to start server with serverPort and load database named db\_name initially.

Attention: the argument serverPort or db\_name can be left out. If you leave out the argument serverPort in the command, then the corresponding value will be set to

default as 9000. If you leave out the argument `db_name` in the command, then the server will start with no databases loaded.

operation: build, load, unload, query, monitor, show, checkpoint, checkall, user, drop

```
// build a new database by a RDF file.
gc.build("test", "data/lubm/lubm.nt", "root", "123456");

// drop a database already built but leave a backup.
gc.drop("test", "root", "123456");

// drop a database already built completely.
gc.drop_r("test", "root", "123456");

// load database
gc.load("test", "root", "123456");

// then you can execute SPARQL query on this database.
answer = gc.query("root", "123456", "test", sparql);

// output information of a database
cout << answer << std::endl;

// show all databases already built and if they are loaded
gc.show();

// show statistical information of a loaded database
gc.monitor("lubm");

// unload this database
gc.unload("lubm", "root", "123456");
```

```
//add a user (with username: Jack, password: 2)
answer = gc.user("add_user", "root", "123456", "Jack", "2");

//add privilege to user Jack (add_query, add_load, add_unload)
answer = gc.user("add_query", "root", "123456", "Jack", "lubm");

//delete privilege of a user Jack
//(delete_query, delete_load, delete_unload)
answer = gc.user("delete_query", "root", "123456", "Jack", "lubm");

//delete user (with username: Jack, password: 2)
answer = gc.user("delete_user", "root", "123456", "Jack", "2");
```

#### Parameters:

db\_name: the name of database, like lubm

format: html, json, txt, csv

sparql: select ?s where ?s ?p ?o .

ds\_path in the server: like /home/data/test.n3

operation: the type of operation: like load, unload, query ...

type: the type of operation that you execute on user, like: add\_user, delete\_user, add\_query, add\_load...

username: the username of the user that execute the operation

password: the password of the user that execute the operation

In order to unify the respond format and, at the same time, allow clients to obtain and handle the response information more efficiently, we change all the response format in ghttp.cpp to JSON.

Requests sent by clients will receive results encapsulated as JSON, which include

StatusCode( a number that simply represents the status of the response), StatusMsg( a string that describes the message from the server), ResponseBody( the body part of the response, if it's a query request, then the body part is the result of the query, and sometimes the body part can be empty).

We defined some StatusCode and it's corresponding message in the following table. For a client sending a request to ghttp, he can simply parse the JSON result to get the StatusCode, StatusMsg, or ResponseBody and then make use of them according to his need. If you are not clear with the meaning of one StatusCode, you can just look up in this table.

StatusCode	StatusMsg	Description
0	"success"	"operation success"
0	"file delete successfully"	
101	"could not open path"	"could not open file"
0	"import RDF file to database done"	
201	"database already built"	
202	"your db name to be built should not end with ".db""	
203	"database not built yet"	
204	"import RDF file to database failed"	
0	"database loaded successfully"	
301	"database already load"	
302	"no load privilege, operation failed"	
303	"unable to load due to loss of lock"	
304	"database not load yet"	
305	"failed to load the database"	
401	"empty SPARQL"	

StatusCode	StatusMsg	Description
402	"update query returns true"	
403	"search query returns false"	
404	"no query privilege, operation failed"	
501	"unable to monitor due to loss of lock"	
0	"database unloaded"	
601	"no unload privilege, operation failed"	
602	"unable to unload due to loss of lock"	
0	"database saved successfully"	
801	"no database"	
802	"no users"	
0	"check identity successfully"	
901	"wrong username"	
902	"wrong password"	
903	"username not found"	
904	"exactly 1 argument is required"	
905	"exactly 2 argument is required"	
0	"operation on users succeeded"	"add, delete users or modify the privilege of a user, operation succeeded"
907	"username already existed, add user failed"	
908	"you cannot delete root, delete user failed"	

StatusCode	StatusMsg	Description
909	"username not exist, delete user failed"	
910	"you can't add privilege to root user"	
911	"add privilege failed"	
912	"add query or load or unload privilege failed"	
913	"you can't delete privilege of root user"	
914	"delete privilege failed"	
915	"not root user, no privilege to perform this operation"	
916	"username not exist, change password failed"	

Table 2: StatusCode Definition

ghhttp supports concurrent read-only queries, but when queries containing updates come, the whole database will be locked. The number of concurrent running queries is suggested to be lower than 300 on a machine with dozens of kernel threads, though we can run 13000 queries concurrently in our experiments. To use the concurrency feature, you had better modify the system settings of 'open files' and 'maximum processes' to 65535 or larger. Three scripts are placed in setup folder to help you modify the settings in different Linux distributions.

**If queries containing updates are sent via ghhttp, a checkpoint command must be sent and done by ghhttp console before we shutdown the database server. Otherwise, the updates may not be synchronize to disk and will be lost if the ghhttp server is stopped.**

**Attention: you can not stop ghhttp by simply Ctrl+C, because this will cause the changes of databases lost. In order to stop the ghhttp server, you can type `bin/shutdown serverPort`.**

#### **4. gserver   This is not maintained now.**

gserver is a daemon. It should be launched first when accessing gStore by gclient or API. It communicates with client through socket.

```
[bookug@localhost gStore]$ bin/gserver -s  
Server started at port 3305
```

```
[bookug@localhost gStore]$ bin/gserver -t  
Server stopped at port 3305
```

You can also assign a custom port for listening.

```
[bookug@localhost gStore]$ bin/gserver -s -p 3307  
Server started at port 3307.
```

Notice: Multiple threads are not supported by gserver. If you start up gclient in more than one terminal in the same time, gserver will go down.

#### **5. gclient   This is not maintained now.**

gclient is designed as a client to send commands and receive feedbacks.

```
[bookug@localhost gStore]$ bin/gclient  
ip=127.0.0.1 port=3305  
gsql>
```

You can also assign gserver's ip and port.

```
[bookug@localhost gStore]$ bin/gclient 172.31.19.15 3307  
ip=172.31.19.15 port=3307  
gsql>
```

We can use these following commands now:



- `help` shows the information of all commands
- `import db_name rdf_triple_file_name` build a database from RDF triple file
- `load db_name` load an existing database
- `unload db_name` unload database, but will not delete it on disk, you can load it next time
- `sparql "query_string"` query the current database with a SPARQL query string(quoted by "")
- `show` displays the name of the current loaded database

Notice:

- at most one database can be loaded in the gclient console
- you can place ' ' or '\t' between different parts of command, but not use characters like ';'.
- you should not place any space or tab ahead of the start of any command

## **6. gconsole This is not maintained now.**

gconsole is the main console of gStore, which integrates with all functions to operate on gStore, as well as some system commands. Completion of commands name, line editing features and access to the history list are all provided. Feel free to try it, and you may have a wonderful tour!(spaces or tabs at the beginning or end is ok, and no need to type any special characters as separators)

Just type `bin/gconsole` in the root directory of gStore to use this console, and you will find a `gstore>` prompt, which indicates that you are in native mode and can type in native commands now. There are another mode of this console, which is called remote mode. Just type `connect` in the native mode to enter the remote mode, and type `disconnect` to exit to native mode.(the console connect to a

gStore server whose ip is '127.0.0.1' and port is 3305, you can specify them by  
`type connect gStore_server_ip gStore_server_port)`

You can use `help` or `?` either in native mode or remote mode to see the help information, or you can type `help command_name` or `? command_name` to see the information of a given command. Notice that there are some differences between the commands in native mode and commands in remote mode. For example, system commands like `ls`, `cd` and `pwd` are provided in native mode, but not in remote mode. Also take care that not all commands contained in the help page are totally achieved, and we may change some functions of the console in the future.

What we have done is enough to bring you much convenience to use gStore, just enjoy it!

**7. gadd** gadd is used to add triples in a file to an existing database.

Usage: `bin/gadd db_name rdf_triple_file_path.`

```
[bookug@localhost gStore]$ bin/gadd lubm ./data/lubm/lubm.nt
...
argc: 3 DB_store:lubm    insert file:./data/lubm/lubm.nt
get important pre ID
...
insert rdf triples done.
inserted triples num: 99550
```

**8. gsub** gsub is used to remove triples from an existing database.

Usage: `bin/gsub db_name rdf_triple_file_path.`

```
[bookug@localhost gStore]$ bin/gsub lubm ./data/lubm/lubm.nt
...
argc: 3 DB_store:lubm    remove file: ./data/lubm/lubm.nt
...
```

```
remove rdf triples done.  
removed triples num: 99550
```

**9. gmonitor** After starting ghttp, type `bin/gmonitor ip port` to check current status of gStore.

```
[bookug@localhost bin]$ bin/gmonitor 127.0.0.1 9000  
parameter: ?operation=monitor  
request: http://127.0.0.1:9000/%3Foperation%3Dmonitor  
null--->[HTTP/1.1 200 OK]  
Content-Length--->[127]  
database: lubm  
triple num: 99550  
entity num: 28413  
literal num: 0  
subject num: 14569  
predicate num: 17  
connection num: 7
```

**10. gshow** After starting ghttp, type `bin/gshow ip port` to check loaded database.

```
[bookug@localhost gStore]$ bin/gshow 127.0.0.1 9000  
parameter: ?operation=show  
request: http://127.0.0.1:9000/%3Foperation%3Dshow  
null--->[HTTP/1.1 200 OK]  
Content-Length--->[4]  
lubm
```

**11. shutdown** After starting ghttp, type `bin/shutdown port` to stop the server.

**12. ginit** If you want to restore the initial configuration of the ghttp server, type `bin/ginit` to rebuild the system.db.

**13. test utilities** A series of test program are placed in the `test/` folder, and we will introduce the two useful ones: `gtest.cpp` and `full_test.sh`

**gtest is used to test gStore with multiple datasets and queries.**

To use `gtest` utility, please type `make gtest` to compile the `gtest` program first. Program `gtest` is a test tool to generate structural logs for datasets. Please type `./gtest --help` in the working directory for details.

**Please change paths in the `test/gtest.cpp` if needed.**

You should place the datasets and queries in this way:

```
DIR/WatDiv/database/*.nt
```

```
DIR/WatDiv/query/*.sql
```

Notice that `DIR` is the root directory where you place all datasets waiting to be used by `gtest`. And `WatDiv` is a class of datasets, as well as `LUBM`. Inside `WatDiv`(or `LUBM`, etc. please place all datasets(named with `.nt`) in a `database/` folder, and place all queries(corresponding to datasets, named with `.sql`) in a `query` folder.

Then you can run the `gtest` program with specified parameters, and the output will be sorted into three logs in `gStore` root directory: `load.log`/(for database loading time and size), `time.log`/(for query time) and `result.log`/(for all query results, not the entire output strings, but the information to record the selected two database systems matched or not).

All logs produced by this program are in TSV format(separated with ‘\t’), you can load them into `Calc/Excel/Gnumeric` directly. Notice that time unit is ms, and space unit is kb.

**`full_test.sh` is used to compare the performance of `gStore` and other database systems on multiple datasets and queries.**

To use `full_test.sh` utility, please download the database system which you want to test and compare, and set the exact position of database systems and datasets in this script. The name strategy should be the same as the requirements of `gtest`, as well as the logs strategy.

Only `gStore` and `Jena` are tested and compared in this script, but it is easy to add other database systems, if you would like to spend some time on reading this script. You may go to [Test Report](#) or [Frequently Asked Questions](#) for help if you encounter a problem.

## Part II

# Advanced

### Chapter 05: HTTP API Explanation

**We provide HTTP api(suggested) and socket api, corresponding to ghttp and gserver respectively.**

This chapter provides API for ghttp. Compared with socket API, HTTP API is more stable and more standard, and can maintain connection. Socket API can not guarantee correct transmission, so the network transmission is faster.

#### Easy Examples

We provide JAVA, C++, Python, PHP and Nodejs API for ghttp now. Please see `api/http`. To use these examples, please make sure that executables have already been generated.

Next, **start up ghttp service by using `./ghttp` command**. It is ok if you know a running usable ghttp server and try to connect to it. (you don't need to change anything if using examples, just by default). Then, for Java and C++ code, you need to compile the example codes in the directory `gStore/api/http/`.

Finally, go to the example directory and run the corresponding executables. All these four executables will connect to a specified ghttp server and do some load or query operations. Be sure that you see the query results in the terminal where you run the examples, otherwise please go to [Frequently Asked Questions](#) for help or report it to us.(the report approach is described in [README](#))

You are advised to read the example code carefully, as well as the corresponding Makefile. This will help you to understand the API, specially if you want to write your own programs based on the API interface.

## API Structure

The HTTP API of gStore is placed in `api/http/` directory in the root directory of gStore, whose contents are listed below:

- `gStore/api/http/`
  - `cpp/` (C++ API)
    - \* `client.cpp` (source code of C++ API)
    - \* `client.h`
    - \* `example` (example program to show the basic idea of using the C++ API)
      - `example.cpp`
      - `Benchmark.cpp`
      - `Benchmark1.cpp`
      - `CppAPIExample.cpp`
      - `Makefile`
    - \* `Makefile` (compile and build lib)
  - `java/` (the Java API)
    - \* `client.java`
    - \* `lib/`
    - \* `src/`
      - `jgsc/GstoreConnector.java`
      - `Makefile`
    - \* `example/`
      - `Benchmark.java`
      - `Benchmark1.java`
      - `JavaAPIExample.java`
      - `Makefile`

- python/ (the Python API)
  - \* example/
    - Benchmark.py
    - PyAPIExample.py
  - \* src/
    - GstoreConnector.py
- php/ (the Php API)
  - \* example/
    - Benchmark.php phpAPIExample.php
  - \* src/
    - GstoreConnector.php
- nodejs/ (the Nodejs API)
  - \* example.js
  - \* index.js
  - \* package.json
  - \* README.md

## C++ API

**Interface** To use the C++ API, please place the phrase `#include "GstoreConnecotr.h"` in your `cpp` code. Functions in `GstoreConnector.h` should be called like below:

```
[language={ [Visual]C++}]
// initialize
GstoreConnector gc("172.31.222.93", 9016);

// build a new database by a RDF file.
// note that the relative path is related to gserver
gc.build("test", "data/lubm/LUBM_10.n3", "root", "123456");
```



```

// load the database that you built.
gc.load("test", "root", "123456");

// then you can execute SPARQL query on this database.
std::string answer = gc.query("root", "123456", "test", sparql);
std::cout << answer << std::endl;

// make a SPARQL query and save the result in ans.txt
gc.query("root", "123456", "test", sparql, "ans.txt");

// unload this database.
gc.unload("test", "root", "123456");

// you can load some exist database directly and then query.
gc.load("lubm", "root", "123456");
answer = gc.query("root", "123456", "lubm", sparql);

std::cout << answer << std::endl;
gc.unload("lubm", "root", "123456");

```

The original declaration of these functions are as below:

```

GstoreConnector();

bool build(std::string _db_name, std::string _rdf_file_path,
           std::string username, std::string password);

bool load(std::string _db_name, std::string username,
           std::string password);

bool unload(std::string _db_name, std::string username,

```

```

        std::string password);

void query(std::string username, std::string password,
          std::string db_name, std::string sparql,
          std::string filename);

std::string query(std::string username,
                 std::string password, std::string db_name,
                 std::string sparql);

```

## Java API

**Interface** To use the Java API, please see [gStore/api/http/java/src/jgsc/Gstore-Connector.java](http://java/src/jgsc/Gstore-Connector.java). Functions should be called like below:

```

// initialize
GstoreConnector gc = new GstoreConnector("172.31.222.78", 3305);

// build the database
gc.build("LUBM10", "data/lubm/lubm.nt", "root", "123456");

// load the database you built
gc.load("LUBM10", "root", "123456");

// you can load some exist database directly and then query.
gc.load("LUBM10", "root", "123456");

// make a SPARQL query
answer = gc.query("root", "123456", "LUBM10", sparql);
System.out.println(answer);

// make a SPARQL query and save the result in ans.txt

```

```
gc.query("root", "123456", "LUBM10", sparql, "ans.txt");
```

```
// unload the database
```

```
gc.unload("LUBM10", "root", "123456");
```

The original declaration of these functions are as below:

```
public class GstoreConnector();
```

```
public boolean build(String _db_name, String _rdf_file_path,  
    String _username, String _password);
```

```
public boolean load(String _db_name,  
    String _username, String _password);
```

```
public boolean unload(String _db_name,  
    String _username, String _password);
```

```
public String query(String _username, String _password,  
    String _db_name, String _sparql);
```

```
public void query(String _username, String _password,  
    String _db_name, String _sparql, String _filename);
```

## Python API

**Interface** To use Python API, please see `gStore/api/http/python/src/GstoreConnector.py`. Functions should be called like following:

```
[language={python}]
```

```
# start a gc with given IP and Port
```

```
gc = GstoreConnector.GstoreConnector("172.31.222.78", 3305)
```

```

# build database with a RDF graph
ret = gc.build("test", "data/lubm/lubm.nt", username, password)

# load the database
ret = gc.load("test", username, password)

# query
print (gc.query(username, password, "test", sparql))

# query and save the result in a file
gc.fquery(username, password, "test", sparql, filename)

# unload the database
ret = gc.unload("test", username, password)

```

The original declaration of these functions are as below:

```

public class GstoreConnector()

def build(self, db_name, rdf_file_path, username, password):

def load(self, db_name, username, password):

def unload(self, db_name, username, password):

def query(self, username, password, db_name, sparql):

def fquery(self, username, password, db_name, sparql, filename):

```

## **Php API**

**Interface** To use Php API, please see `gStore/api/http/php/src/GstoreConnector.php`. Functions should be called like following:

```
// start a gc
$gc = new GstoreConnector("172.31.222.78", 3305);

// build database
$ret = $gc->build("test", "data/lubm/lubm.nt",
    $username, $password);
echo $ret . PHP_EOL;

// load rdf
$ret = $gc->load("test", $username, $password);
echo $ret . PHP_EOL;

// query
echo $gc->query($username, $password, "test", $sparql) . PHP_EOL;

// fquery--make a SPARQL query and save the result in the file
$gc->fquery($username, $password, "test", $sparql, $filename);

// unload rdf
$ret = $gc->unload("test", $username, $password);
echo $ret . PHP_EOL;
```

The original declaration of these functions are as below:

```
class GstoreConnector

function build($db_name, $rdf_file_path, $username, $password)
```

```

function load($db_name, $username, $password)

function unload($db_name, $username, $password)

function query($username, $password, $db_name, $sparql)

function fquery($username, $password, $db_name, $sparql, $filename)

```

## Nodejs API

**Interface** Before using Nodejs API, type `npm install request` and `npm install request-promise` under the `nodejs` folder to add the required module.

To use Nodejs API, please see `gStore/api/http/nodejs/index.js`. Functions should be called like following:

```

# start a gc with given IP and Port
gc = new GStoreClient(username, password, "http://127.0.0.1:9000")

# build database with a RDF graph
ret = gc.build("test", "data/lubm/lubm.nt")

# load the database
ret = gc.load("test")

# query
print (gc.query("test", sparql))

# unload the database
ret = gc.unload("test")

```

The original declaration of these functions are as below:

```
class GStoreClient

    async build(db = '', dataPath = '')

    async load(db = '')

    async unload(db = '')

    async query(db = '', sparql = '')
```

## Chapter 06: Socket API Explanation

**This API is not maintained now.**

### Easy Examples

We provide JAVA, C++, PHP and Python API for gStore now. Please refer to example codes in `api/socket/cpp/example`, `api/socket/java/example`, `api/socket/php` and `api/socket/python/example`. To use the four examples to have a try, please ensure that executables have already been generated. Otherwise, for Java and C++, just type `make APIexample` in the root directory of gStore to compile the codes, as well as API.

Next, **start up a gStore server by using `./gserver` command**. It is ok if you know a running usable gStore server and try to connect to it, but notice that **the server ip and port of server and client must be matched**. (you don't need to change any thing if using examples, just by default) Then, for Java and C++ code, you need to compile the example codes in the directory `gStore/api/socket/`. We provide a utility to do this, and you just need to type `make APIexample` in the root directory of gStore. Or you can compile the codes by yourself, in this case please go to `gStore/api/socket/cpp/example/` and `gStore/api/socket/java/example/`, respectively.

Finally, go to the example directory and run the corresponding executables. For C++, just use `./example` command to run it. And for Java, use `make run` command or `java -cp ../lib/GstoreJavaAPI.jar:. JavaAPIExample` to run it. For PHP, use `php ./PHPAPIExample`. For python, use `python ./PythonAPIExample`. All these four executables will connect to a specified gStore server and do some load or query operations. Be sure that you see the query results in the terminal where you run the examples, otherwise please go to [Frequently Asked Questions](#) for help or report it to us. (the report approach is described in [README](#))

You are advised to read the example code carefully, as well as the corresponding Makefile. This will help you to understand the API, specially if you want to write



your own programs based on the API interface.

## API structure

The API of gStore is placed in `api/socket/` directory in the root directory of gStore, whose contents are listed below:

- `gStore/api/socket/`
  - `cpp/` (the C++ API)
    - \* `src/` (source code of C++ API, used to build the `lib/libgstoreconnector.a`)
      - `GstoreConnector.cpp` (interfaces to interact with gStore server)
      - `GstoreConnector.h`
      - `Makefile` (compile and build lib)
    - \* `lib/` (where the static lib lies in)
      - `.gitignore`
      - `libgstoreconnector.a` (only exist after compiled, you need to link this lib when you use the C++ API)
    - \* `example/` (small example program to show the basic idea of using the C++ API)
      - `CppAPIExample.cpp`
      - `Makefile`
  - `java/` (the Java API)
    - \* `src/` (source code of Java API, used to build the `lib/Gstore-JavaAPI.jar`)
      - `jgsc/GstoreConnector.java` (the package which you need to import when you use the Java API)
      - `Makefile` (compile and build lib)
    - \* `lib/`

- .gitignore
- GstoreJavaAPI.jar (only exist after compiled, you need to include this JAR in your class path)
- \* example/ (small example program to show the basic idea of using the Java API)
  - JavaAPIExample.cpp
  - Makefile
- php/ (the PHP API)
  - \* GstoreConnector.php (source code of PHP API, you need to include this file when you use the PHP API)
  - \* PHPAPIExample.php (small example program to show the basic idea of using the PHP API)
- python/ (the Python API)
  - \* src/ (source code of Python API)
    - GstoreConnector.py (the package which you need to import when you use the Python API)
  - \* example/ (small example program to show the basic idea of using the Python API)
    - PythonAPIExample.py

## C++ API

**Interface** To use the C++ API, please place the phrase `#include "GstoreConnector.h"` in your cpp code. Functions in `GstoreConnector.h` should be called like below:

```
// initialize the Gstore server's IP address and port.
GstoreConnector gc("127.0.0.1", 3305);
// build a new database by a RDF file.
// note that the relative path is related to gserver.
```

```

gc.build("LUBM10", "example/LUBM_10.n3");
// then you can execute SPARQL query on this database.
std::string sparql = "select ?x where \
{ \
?x    <rdf:type>      <ub:UndergraduateStudent>. \
?y    <ub:name> <Course1>. \
?x    <ub:takesCourse> ?y. \
?z    <ub:teacherOf>   ?y. \
?z    <ub:name> <FullProfessor1>. \
?z    <ub:worksFor>    ?w. \
?w    <ub:name>      <Department0>. \
}";
std::string answer = gc.query(sparql);
// unload this database.
gc.unload("LUBM10");
// also, you can load some exist database directly and then query.
gc.load("LUBM10");
// query a SPARQL in current database
answer = gc.query(sparql);

```

The original declaration of these functions are as below:

```

GstoreConnector();
GstoreConnector(string _ip, unsigned short _port);
GstoreConnector(unsigned short _port);
bool load(string _db_name);
bool unload(string _db_name);
bool build(string _db_name, string _rdf_file_path);
string query(string _sparql);

```

Notice:

1. When using GstoreConnector(), the default value for ip and port is 127.0.0.1 and 3305, respectively.

2. When using `build()`, the `rdf_file_path`(the second parameter) should be related to the position where `gserver` lies in.
3. Please remember to unload the database you have loaded, otherwise things may go wrong.(the errors may not be reported!)

**Compile** You are advised to see `gStore/api/socket/cpp/example/Makefile` for instructions on how to compile your code with the C++ API. Generally, what you must do is compile your own code to object with header in the C++ API, and link the object with static lib in the C++ API.

Let us assume that your source code is placed in `test.cpp`, whose position is `${GSTORE}/gStore/`.(if using `devGstore` as name instead of `gStore`, then the path is `${GSTORE}/devGstore/` directory first:

```
Use g++ -c -I${GSTORE}/gStore/api/socket/cpp/src/
test.cpp -o test.o to compile your test.cpp into test.o, rela-
tive API header is placed in api/socket/cpp/src/.
```

```
Use g++ -o test test.o -L${GSTORE}/gStore/api/socket/cpp/lib/
-lgstoreconnector to link your test.o with the libgstoreconnec-
tor.a(a static lib) in api/socket/cpp/lib/.
```

Then you can type `./test` to execute your own program, which uses our C++ API. It is also advised for you to place relative compile commands in a `Makefile`, as well as other commands if you like.

## Java API

**Interface** To use the Java API, please place the phrase `import jgsc.GstoreConnector;` in your java code. Functions in `GstoreConnector.java` should be called like below:

```
// initialize the Gstore server's IP address and port.
GstoreConnector gc = new GstoreConnector("127.0.0.1", 3305);
```

```

// build a new database by a RDF file.
// note that the relative path is related to gserver.
gc.build("LUBM10", "example/LUBM_10.n3");
// then you can execute SPARQL query on this database.
String sparql = "select ?x where " + "{" +
"?x    <rdf:type>      <ub:UndergraduateStudent>. " +
"?y    <ub:name> <Course1>. " +
"?x    <ub:takesCourse> ?y. " +
"?z    <ub:teacherOf>   ?y. " +
"?z    <ub:name> <FullProfessor1>. " +
"?z    <ub:worksFor>    ?w. " +
"?w    <ub:name>      <Department0>. " +
"}";
String answer = gc.query(sparql);
//unload this database.
gc.unload("LUBM10");
//also, you can load some exist database directly and then query.
gc.load("LUBM10");// query a SPARQL in current database
answer = gc.query(sparql);

```

The original declaration of these functions are as below:

```

GstoreConnector();
GstoreConnector(string _ip, unsigned short _port);
GstoreConnector(unsigned short _port);
bool load(string _db_name);
bool unload(string _db_name);
bool build(string _db_name, string _rdf_file_path);
string query(string _sparql);

```

Notice:

1. When using GstoreConnector(), the default value for ip and port is 127.0.0.1 and 3305, respectively.

2. When using build(), the rdf\_file\_path(the second parameter) should be related to the position where gserver lies in.
3. Please remember to unload the database you have loaded, otherwise things may go wrong.(the errors may not be reported!)

**Compile** You are advised to see gStore/api/socket/java/example/Makefile for instructions on how to compile your code with the Java API. Generally, what you must do is compile your own code to object with jar file in the Java API.

Let us assume that your source code is placed in test.java, whose position is \${GSTORE}/gStore/.(if using devGstore as name instead of gStore, then the path is \${GSTORE}/devGstore/ directory first:

```
Use javac -cp ${GSTORE}/gStore/api/socket/java/lib/GstoreJavaAPI.jar
test.java to compile your test.java into test.class with the Gstore-
JavaAPI.jar(a jar package used in Java) in api/socket/java/lib/.
```

Then you can type java -cp \${GSTORE}/gStore/api/socket/java/lib/GstoreJavaAPI.jar test to execute your own program(notice that the “:.” in command cannot be neglected), which uses our Java API. It is also advised for you to place relative compile commands in a Makefile, as well as other commands if you like.

## PHP API

**Interface** To use the PHP API, please place the phrase include('GstoreConnector.php'); in your php code. Functions in GstoreConnector.php should be called like below:

```
// initialize the Gstore server's IP address and port.
$gc = new Connector("127.0.0.1", 3305);
// build a new database by a RDF file.
// note that the relative path is related to gserver.
$gc->build("LUBM10", "example/LUBM_10.n3");
// then you can execute SPARQL query on this database.
```

```

$sparql = "select ?x where " + "{" +
"?x    <rdf:type>      <ub:UndergraduateStudent>. " +
"?y    <ub:name> <Course1>. " +
"?x    <ub:takesCourse> ?y. " +
"?z    <ub:teacherOf>   ?y. " +
"?z    <ub:name> <FullProfessor1>. " +
"?z    <ub:worksFor>    ?w. " +
"?w    <ub:name>      <Department0>. " +
"}";
$answer = gc->query($sparql);
//unload this database.
$gc->unload("LUBM10");
//also, you can load some exist database directly and then query.
$gc->load("LUBM10");// query a SPARQL in current database
$answer = gc->query(sparql);

```

The original declaration of these functions are as below:

```

class Connector {
    public function __construct($host, $port);
    public function send($data);
    public function recv();
    public function build($db_name, $rdf_file_path);
    public function load($db_name);
    public function unload($db_name);
    public function query($sparql);
    public function __destruct();
}

```

Notice:

1. When using Connector(), the default value for ip and port is 127.0.0.1 and 3305, respectively.

2. When using build(), the rdf\_file\_path(the second parameter) should be related to the position where gserver lies in.
3. Please remember to unload the database you have loaded, otherwise things may go wrong.(the errors may not be reported!)

**Run** You can see gStore/api/socket/php/PHPAPIExample for instructions on how to use PHP API. PHP script doesn't need compiling. You can run PHP file directly or use it in your web project.

## Python API

**Interface** To use the Python API, please place the phrase from GstoreConnector import GstoreConnector in your python code. Functions in GstoreConnector.py should be called like below:

```
// initialize the Gstore server's IP address and port.
gc = GstoreConnector('127.0.0.1', 3305)
// build a new database by a RDF file.
// note that the relative path is related to gserver.
gc.build('LUBM10', 'data/LUBM_10.n3')
// then you can execute SPARQL query on this database.
$sparql = "select ?x where " + "{" +
"?x    <rdf:type>      <ub:UndergraduateStudent>. " +
"?y    <ub:name> <Course1>. " +
"?x    <ub:takesCourse>  ?y. " +
"?z    <ub:teacherOf>    ?y. " +
"?z    <ub:name> <FullProfessor1>. " +

"?z    <ub:worksFor>      ?w. " +
"?w    <ub:name>      <Department0>. " +
"}";
```



```

answer = gc.query(sparql)
//unload this database.
gc.unload('LUBM10')
//also, you can load some exist database directly and then query.
gc.load('LUBM10')// query a SPARQL in current database
answer = gc.query(sparql)

```

The original declaration of these functions are as below:

```

class GstoreConnector {
    def _connect(self)
    def _disconnect(self)
    def _send(self, msg):
    def _recv(self)
    def _pack(self, msg):
    def _communicate(f):
    def __init__(self, ip='127.0.0.1', port=3305):
    @_communicate
    def test(self)
    @_communicate
    def load(self, db_name)
    @_communicate
    def unload(self, db_name)
    @_communicate
    def build(self, db_name, rdf_file_path)
    @_communicate
    def drop(self, db_name)
    @_communicate
    def stop(self)
    @_communicate
    def query(self, sparql)
    @_communicate

```

```
def show(self, _type=False)
}
```

Notice:

1. When using `GstoreConnector()`, the default value for ip and port is 127.0.0.1 and 3305, respectively.
2. When using `build()`, the `rdf_file_path`(the second parameter) should be related to the position where gserver lies in.
3. Please remember to unload the database you have loaded, otherwise things may go wrong.(the errors may not be reported!)

**Run** You are advised to see `gStore/api/socket/python/example/PythonAPIExample` for examples on how to use python API. Python file doesn't need compiling, and you can run it directly.

## Chapter 07: Use gStore in Web

**This Chapter provides a specific example on how to use our API in a web project.**

### Example

Now you have the basic idea on how to use our APIs to connect gStore. Yet you might be still a little confused. Here we provide a simple demo to show you what to do explicitly.

Let's say, you need to use gStore in a web project. PHP is a popular general-purpose scripting language that is especially suited to web development. So, using our PHP API can meet your requirements. Here is what we implement: [demo](#).

First, get your web server ready so it can run PHP files. We won't give detailed instructions on this step here. You can easily google it according to your web server(for example, Apache or Nginx, etc.)

Next, go to your web document root(usually in /var/www/html or apache/htdocs, you can check it in config file), and create a folder named "Gstore". Then copy the GstoreConnector.php file into it. Create a "PHPAPI.php" file. Edit it like below:

```
<?php
//require "../src/GstoreConnector.php";
include( 'gStoreConnector.php' );

$username = "root";
$password = "123456";
$gc = new GstoreConnector("127.0.0.1", 9001);

$dbname = $_POST["databasename"];
$query = $_POST["query"];
$format = $_POST["format"];
$result = $gc->query($username, $password, $dbname, $query);
```

```

$js = json_decode($result);
switch ($format) {
    case 1:
        $vars = $js->head->vars;
        $rows = $js->results->bindings;

        $html = '<html><table class="sparql" border="1"><tr>';
        for ($i = 0; $i < count($vars); $i++) {
            $html.= '<th>' . $vars[$i] . '</th>';
        }
        $html.= '</tr>';
        for ($i = 0; $i < count($rows); $i++) {
            $html.= '<tr>';
            $row = $rows[$i];
            for ($j = 0; $j < count($vars); $j++)
                $html.= '<td>' . htmlspecialchars($row->
                    $vars[$j]->value) . '</td>';
            $html.= '</tr>';
        }
        $html.= '</table></html>';
        echo $html;
        exit;

    case 2:
        $filename = 'sparql.txt';
        header("Content-Type: application/octet-stream");
        header('Content-Disposition: attachment;
            filename="' . $filename . '"');

        $rows = $js->results->bindings;

```

```

for ($i = 0; $i < count($rows); $i++)
    echo $rows[$i] . "\n";
exit;

case 3:
$filename = 'sparql.csv';
header("Content-Type: application/octet-stream");
header('Content-Disposition: attachment;
        filename="' . $filename . '"');

$vars = $js->head->vars;
echo implode(", ", $vars);
echo "\n";
$rows = $js->results->bindings;
for($i = 0; $i < count($rows); $i++) {
    $row = $rows[$i];
    for($j = 0; $j < count($vars); $j++) {
        echo ($row->$vars[$j]->value);
        echo ",";
    }
    echo "\n";
}

exit;
}
?>

```

This PHP file get three parametres from a website, including databasename, sparql and output format. Then it use our PHP API to connect gStore and run the query. Finally, the "switch" part gives the output.

After that, we need a website to collect those imformation(databasename, sparql and output format). We create a html file and use a form to do it, just like below:

```

1 <form id="form_1145884" class="appnitro" method="post" action="
    PHPAPI.php">
2 <div class="form_description">
3   <h2>Gstore SPARQL Query Editor</h2>
4   <p></p>
5 </div>
6 <ul>
7   <li id="li_1" >
8     <label class="description" for="element_1">
9       Database Name
10    </label>
11    <div>
12      <input id="element_1" name="databasename" class="element
        text medium"
13        type="text" maxlength="255" value="dbpedia_2014_reduce">
14      </input>
15    </div>
16  </li>
17
18  <li id="li_3">
19    <label class="description" for="element_3">Query Text </label>
20    <div>
21      <textarea id="element_3" name="sparql" class="element
        textarea large">
22        SELECT DISTINCT ?uri
23        WHERE {
24          ?uri <type> <Astronaut> .
25          { ?uri < nationality > <Russia> . }
26          UNION
27          { ?uri < nationality > <Soviet_Union> . }
28          }
29      </textarea>
30    </div>
31  </li>
32
33  <li id="li_5" >
34    <label class="description" for="element_5">

```

```

35     Results Format
36 </label>
37 <div>
38     <select class="element select medium" id="element_5" name
        ="format">
39         <option value="1" selected="ture">HTML</option>
40         <option value="2" >Text</option>
41         <option value="3" >CSV</option>
42     </select>
43 </div>
44
45 <li class="buttons">
46     <input type="hidden" name="form_id" value="1145884" />
47     <input id="saveForm" class="button_text" type="submit"
48         name="submit" value="Run Query" />
49 </li>
50 </ul>
51 </form>

```

As you can see in the code, we use a <input> element to get the databasename, and <texarea> for sparql, <select> for output format. <form> lable has an attribute "action" which specifies which file to execute. So, when you click the "submit" button, it will call PHPAPI.php file and post the values from the form.

Finally, don't forget to start ghttp on your server.

## Chapter 08: Project Structure

This chapter introduce the whole structure of the gStore system project.

The core source codes are listed below:

- Database/ (calling other core parts to deal with requests from interface part)
  - Database.cpp (achieve functions)
  - Database.h (class, members and functions definitions)
  - Join.cpp (join the node candidates to get results)
  - Join.h (class, members and functions definitions)
  - Strategy.cpp
  - Strategy.h
- KVstore/ (a key-value store to swap between memory and disk)
  - KVstore.cpp (interact with upper layers)
  - KVstore.h
  - ISArray/
    - \* ISArray.cpp
    - \* ISArray.h
    - \* ISBlockManager.cpp
    - \* ISBlockManager.h
    - \* ISEntry.cpp
    - \* ISEntry.h
  - ISTree/
    - \* ISTree.cpp
    - \* ISTree.h
    - \* heap/



- ISHeap.cpp
  - ISHeap.h
  - \* node/
    - ISIntlNode.cpp
    - ISIntlNode.h
    - ISLeafNode.cpp
    - ISLeafNode.h
    - ISNode.cpp
    - ISNode.h
  - \* storage/
    - ISSStorage.cpp
    - ISSStorage.h
  - IArray/
  - IVTree/
  - SITree/
- Query/ (needed to answer SPARQL query)
  - BasicQuery.cpp (basic type of queries without aggregate operations)
  - BasicQuery.h
  - IDList.cpp (candidate list of a node/variable in query)
  - IDList.h
  - ResultFilter.cpp
  - ResultFilter.h
  - ResultSet.cpp (keep the result set corresponding to a query)
  - ResultSet.h
  - SPARQLquery.cpp (deal with a entire SPARQL query)
  - SPARQLquery.h

- Varset.cpp
  - Varset.h
  - QueryCache.cpp
  - QueryCache.h
  - QueryTree.cpp
  - QueryTree.h
  - GeneralEvaluation.cpp
  - GeneralEvaluation.h
  - TempResult.cpp
  - TempResult.h
  - RegexExpression.h
- Signature/ (assign signatures for nodes and edges, but not for literals)
    - SigEntry.cpp
    - SigEntry.h
    - Signature.cpp
    - Signature.h
- VSTree/ (an tree index to prune more efficiently)
    - EntryBuffer.cpp
    - EntryBuffer.h
    - LRUCache.cpp
    - LRUCache.h
    - VNode.cpp
    - VNode.h
    - VSTree.cpp
    - VSTree.h

**The parser part is listed below:**

- Parser/
  - DBParser.cpp
  - DBParser.h
  - RDFParser.cpp
  - RDFParser.h
  - SparqlParser.c (auto-generated, subtle modified manually, compressed)
  - SparqlParser.h (auto-generated, subtle modified manually, compressed)
  - SparqlLexer.c (auto-generated, subtle modified manually, compressed)
  - SparqlLexer.h (auto-generated, subtle modified manually, compressed)
  - TurtleParser.cpp
  - TurtleParser.h
  - Type.h
  - QueryParser.cpp
  - QueryParser.h

**The utilities are listed below:**

- Util/
  - Util.cpp (headers, macros, typedefs, functions. . .)
  - Util.h
  - Bstr.cpp (represent strings of arbitrary length)
  - Bstr.h (class, members and functions definitions)
  - Stream.cpp (store and use temp results, which may be very large)
  - Stream.h

- Triple.cpp (deal with triples, a triple can be divided as subject(entity), predicate(entity), object(entity or literal))
- Triple.h
- BloomFilter.cpp
- BloomFilter.h
- ClassForVListCache.h
- VList.cpp
- VList.h

**The interface part is listed below:**

- Server/ (client and server mode to use gStore)
  - Client.cpp
  - Client.h
  - Operation.cpp
  - Operation.h
  - Server.cpp
  - Server.h
  - Socket.cpp
  - Socket.h
  - client\_http.hpp
  - server\_http.hpp
- web/ (a series of applications/main-program to operate on gStore)
  - ...

**More details** It is really not strange to see something different with the original design in the source code. And some designed functions may have not be achieved so far.

**Others** The `api/` folder in gStore is used to store API program, libs and examples, please go to [API](#) for details. And `test/` is used to store a series test programs or utilities, such as `gtest`, `full_test` and so on. Chapters related with `test/` are [How To Use](#) and [Test Result](#). This project need an ANTLR lib to parse the SPARQL query, whose code is placed in `tools/`(also archived here) and the compiled `libantlr.a` is placed in `lib/` directory.

We place some datasets and queries in `data/` directory as examples, and you can try them to see how gStore works. Related instructions are in [How To Use](#). The `docs/` directory contains all kinds of documents of gStore, including a series of markdown files and two folders, `pdf/` and `jpg/`. Files whose type is pdf are placed in `pdf/` folder, while files with jpg type are placed in `jpg/` folder.

You are advised to start from the [README](#) in the gStore root directory, and visit other chapters only when needed. At last, you will see all documents from link to link if you are really interested in gStore.

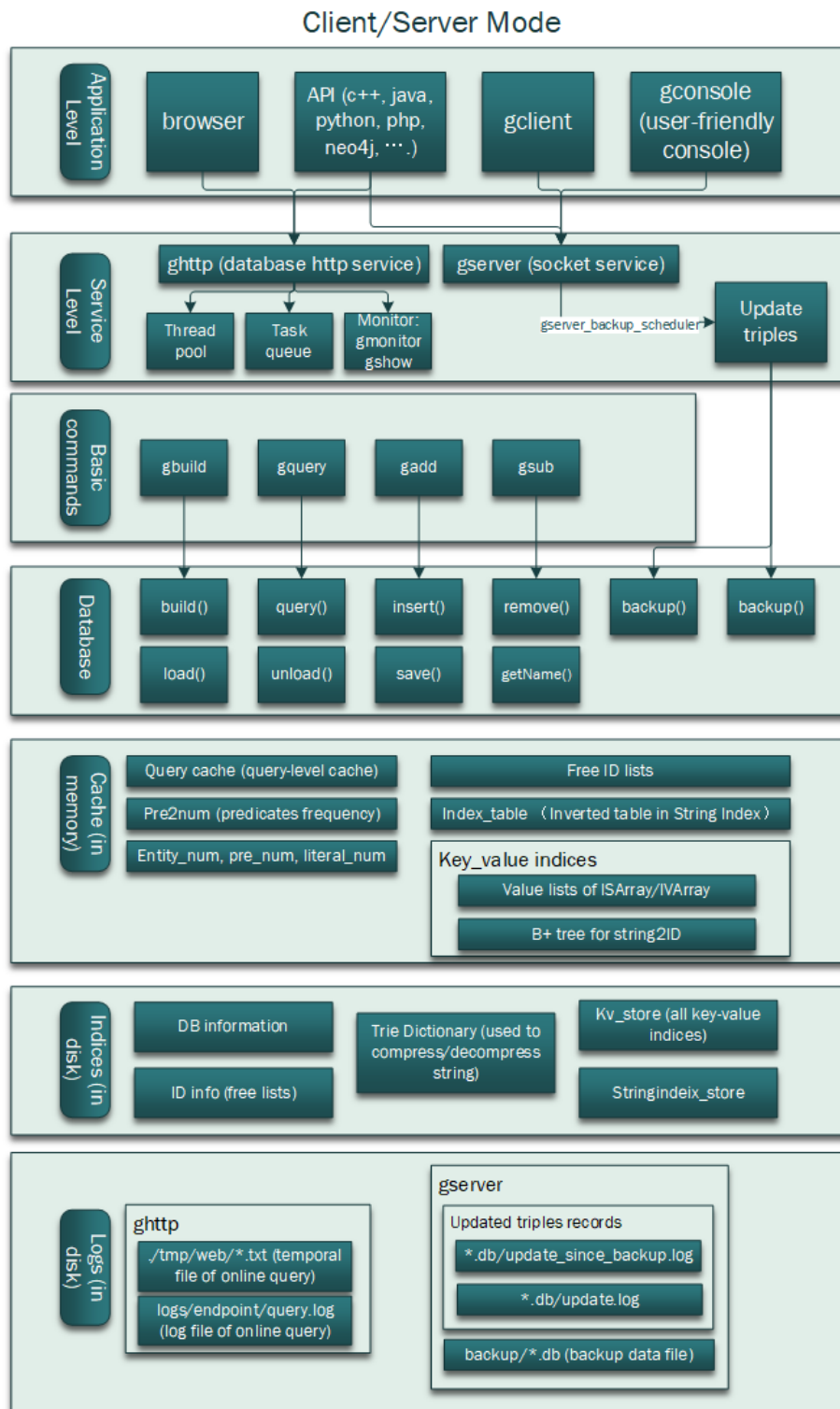
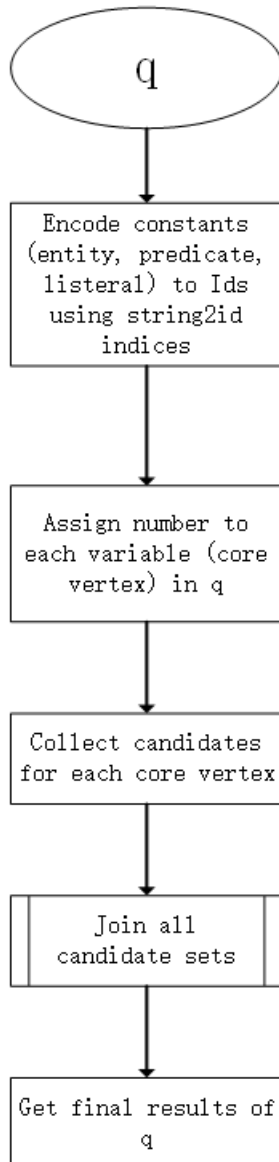


Figure 1: framework

## BGP (Basic Graph Pattern) processing

```
Select ?s where {
  <xxx> <ppp> ?s.
  ?s ?p <yyy>
}
```

Core vertex (degree>1)  
Satellite vertex (select and degree==1)



## Join all candidate sets

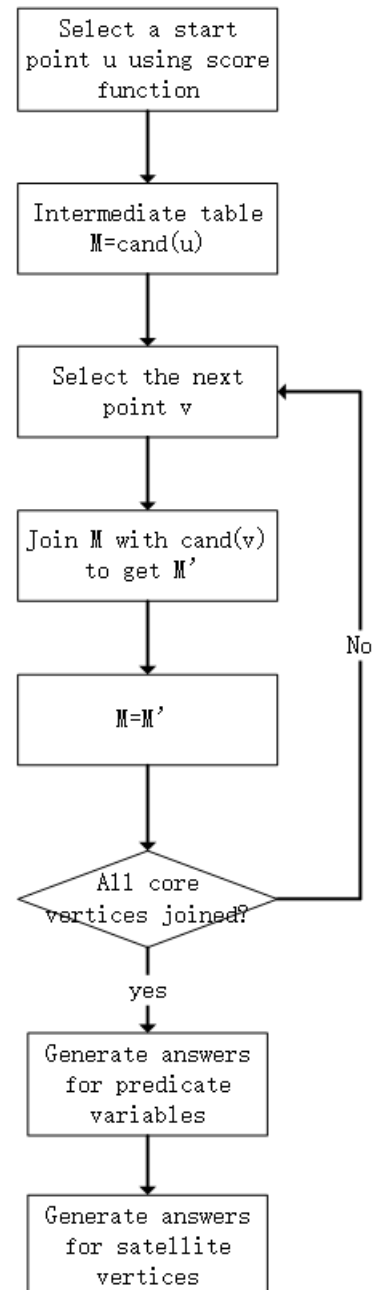


Figure 2: Basic Graph Pattern Processing

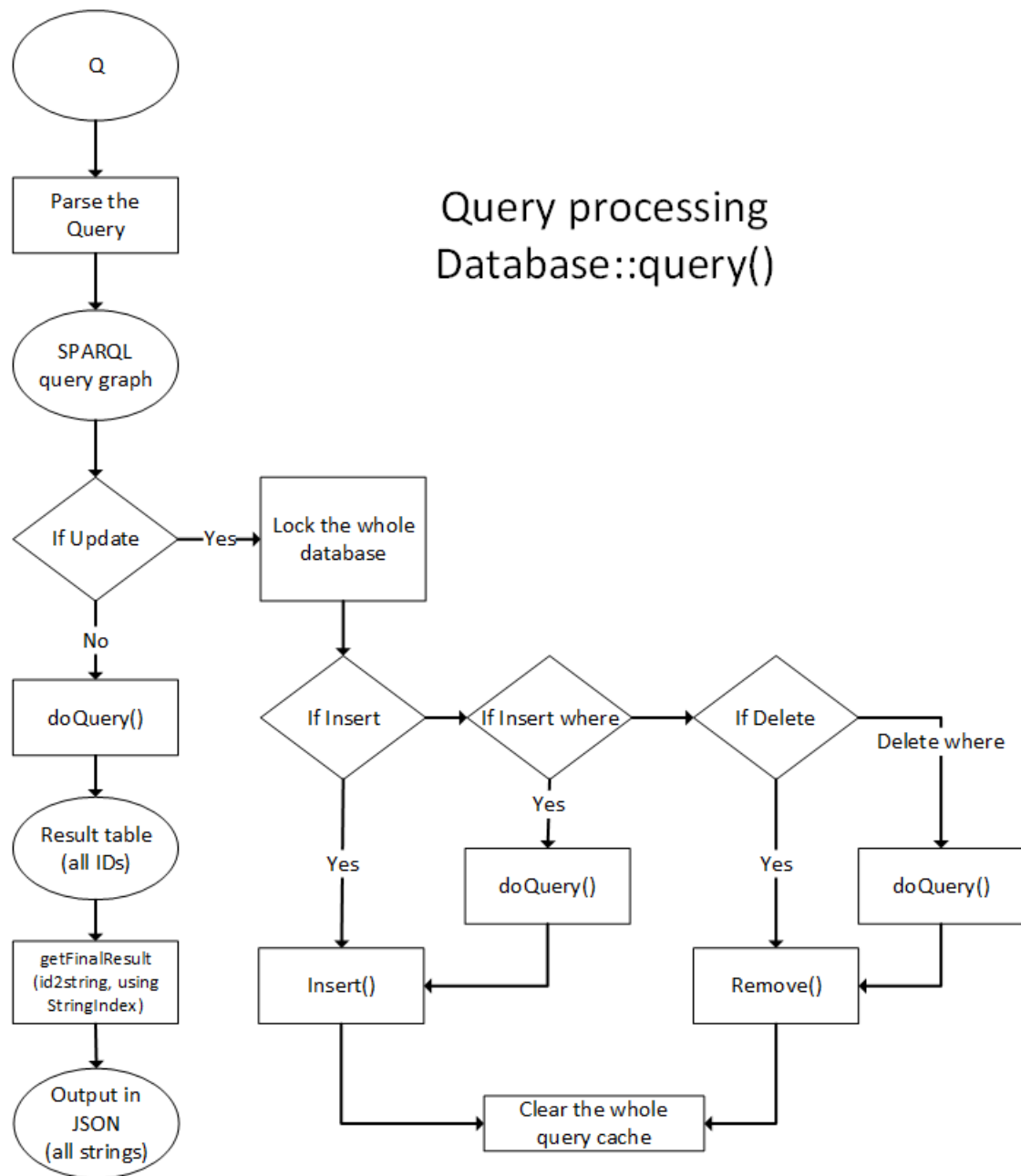


Figure 3: Query Processing



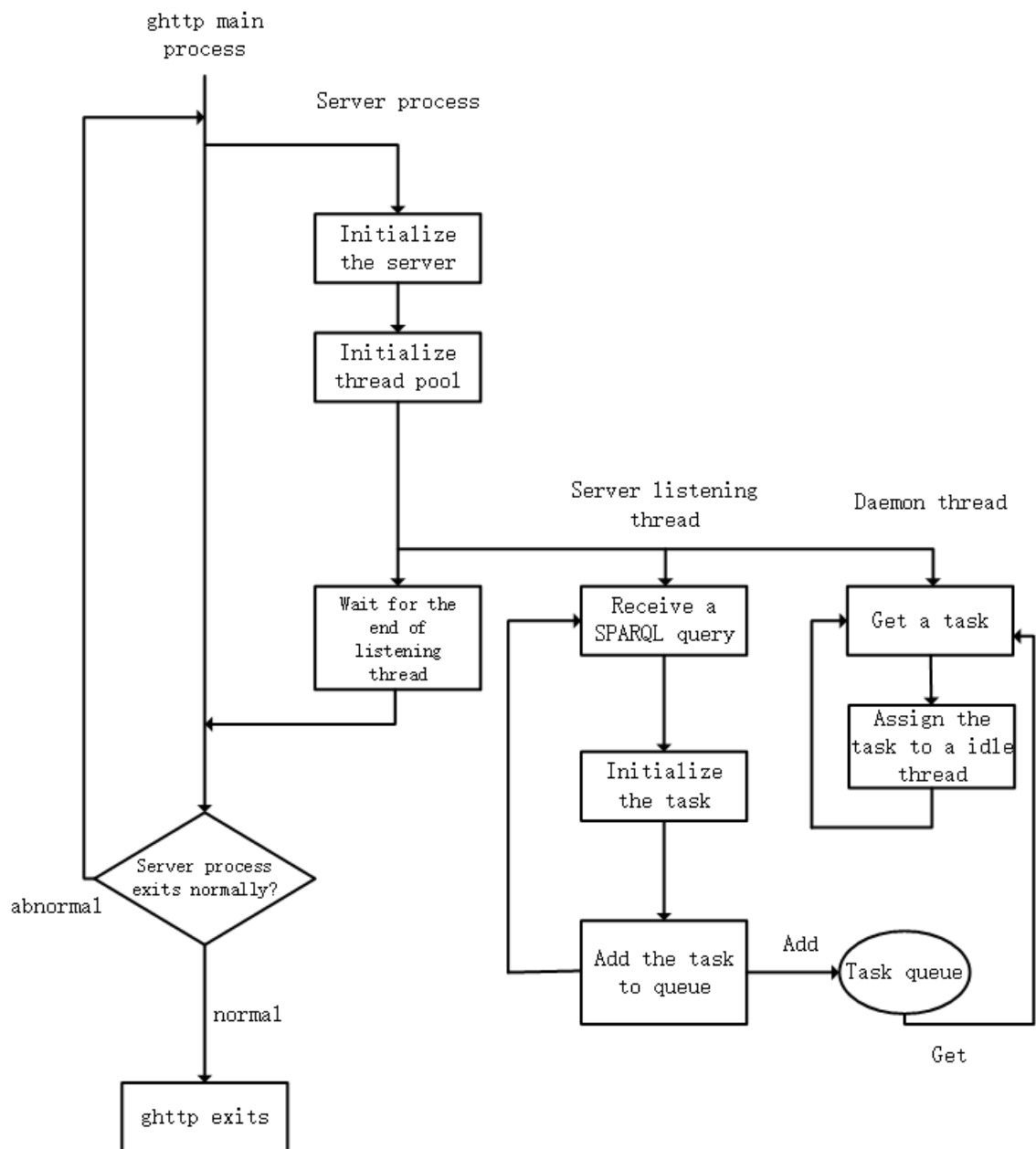


Figure 4: ghttp thread

## Chapter 09: Publications

**Publications related with gStore are listed here:**

- Lei Zou, M. Tamer Özsu, Lei Chen, Xuchuan Shen, Ruizhe Huang, Dongyan Zhao, [gStore: A Graph-based SPARQL Query Engine](#), VLDB Journal , 23(4): 565-590, 2014.
- Lei Zou, Jinghui Mo, Lei Chen, M. Tamer Özsu, Dongyan Zhao, [gStore: Answering SPARQL Queries Via Subgraph Matching](#), Proc. VLDB 4(8): 482-493, 2011.
- Xuchuan Shen, Lei Zou, M. Tamer Özsu, Lei Chen, Youhuan Li, Shuo Han, Dongyan Zhao, [A Graph-based RDF Triple Store](#), ICDE 2015: 1508-1511.
- Peng Peng, Lei Zou, M. Tamer Özsu, Lei Chen, Dongyan Zhao: [Processing SPARQL queries over distributed RDF graphs](#). VLDB Journal 25(2): 243-268 (2016).
- Dong Wang, Lei Zou, Yansong Feng, Xuchuan Shen, Jilei Tian, and Dongyan Zhao, [S-store: An Engine for Large RDF Graph Integrating Spatial Information](#), in Proc. 18th International Conference on Database Systems for Advanced Applications (DASFAA), pages 31-47, 2013.
- Dong Wang, Lei Zou and Dongyan Zhao, [gst-Store: An Engine for Large RDF Graph Integrating Spatiotemporal Information](#), in Proc. 17th International Conference on Extending Database Technology (EDBT), pages 652-655, 2014 (demo).
- Lei Zou, Yueguo Chen, [A Survey of Large-Scale RDF Data Management](#), Communications of CCCF Vol.8(11): 32-43, 2012 (Invited Paper, in Chinese).

## **Chapter 10: Limitations**

1. only support RDF datasets in N-Triples format
2. the filter effects in query propcessing are not good(less than 25%).

## Chapter 11: Frequently Asked Questions

### **When I use the newer gStore system to query the original database, why error?**

The database produced by gStore contains several indexes, whose structures may have been changed in the new gStore version. So, please rebuild your dataset just in case.

### **Why error when I try to write programs based on gStore, just like the Main/gconsole.cpp?**

You need to add these phrases at the beginning of your main program, otherwise gStore will not run correctly:

```
//NOTICE:this is needed to ensure the file path is the work path
chdir(dirname(argv[0]));
//NOTICE:this is needed to set several debug files
Util util;
```

### **Why does gStore report “garbage collection failed” error when I use the Java API?**

You need to adjust the parameters of jvm, see [url1](#) and [url2](#) for details.

### **When I compile the code in ArchLinux, why the error that “no -ltermcap” is reported?**

In ArchLinux, you only need to use `-lreadline` to link the readline library. Please remove the `-ltermcap` in the makefile which is located in the root of the gStore project if you would like to use ArchLinux.

### **Why does gStore report errors that the format of some RDF datasets are not supported?**

gStore does not support all RDF formats currently, please see [formats](#) for details.

### **When I read on GitHub, why are some documents unable to be opened?**

Codes, markdowns or other text files, and pictures can be read directly on GitHub. However, if you are using some light weight browsers like midori, for files in pdf type, please download them and read on your computer or other devices.

### **Why sometimes strange characters appear when I use gStore?**

There are some documents's names are in Chinese, and you don't need to worry about it.

### **In centos7, if the watdiv.db(a generated database after gbuild) is copied or compressed/uncompressed, the size of watdiv.db will be different(generally increasing) if using `du -h` command to check?**

It's the change of B+-trees' size in watdiv/kv\_store/ that causes the change of the whole database's size. The reason is that in storage/Storage.cpp, many operations use fseek to move file pointer. As everyone knows, file is organized in blocks, and if we request for new block, file pointer may be moved beyond the end of this file(file operations are all achieved by C in gStore, no errors are reported), then contents will be written in the new position!

In **Advanced Programming In The Unix Environment**, "file hole" is used to describe this phenomenon. "file hole" will be filled with 0, and it's also one part of the file. You can use `ls -l` to see the size of file(computing the size of holes), while `du -h` command shows the size of blocks that directory/file occupies in system. Generally, the output of `du -h` is large than that of `ls -l`, but if "file hole" exists, the opposite is the case because the size of holes are neglected.

The actual size of files containing holes are fixed, while in some operation systems, holes will be transformed to contents(also 0) when copied. Operation `mv` will not affect the size if not across different devices.(only need to adjust the file tree index) However, `cp` and all kinds of compress methods need to scan the file and transfer data.(there are two ways to achieve `cp` command, neglect holes or not, while the output size of `ls -l` not varies)

It is valid to use "file hole" in C, and this is not an error, which means you can go

on using gStore. We achieve [a small program](#) to describe the “file holes”, you can download and try it yourself.

**In gclient console, a database is built, queried, and then I quit the console. Next time I enter the console, load the originally imported database, but no output for any queries(originally the output is not empty)?**

You need to unload the using database before quitting the gclient console, otherwise errors come.

**If query results contain null value, how can I use the [full\\_test](#) utility? Tab separated method will cause problem here because null value cannot be checked!**

You may use other programming language(for example, Python) to deal with the null value cases. For example, you can change null value in output to special character like ‘,’, later you can use the [full\\_test](#) utility.

**When I compile and run the API examples, it reports the “unable to connect to server” error?**

Please use `./gserver` command to start up a gStore server first, and notice that the server ip and port must be matched.

**When I use the Java API to write my own program, it reports “not found main class” error?**

Please ensure that you include the position of your own program in class path of java. The whole command should be something like `java -cp /home/bookug/project/devGstore/api/java/lib/GstoreJavaAPI.jar:. JavaAPIExample`, and the “:.” in this command cannot be neglected.

## **Chapter 12: Recipe Book**

**This chapter introduces some useful tricks if you are using gStore to implement applications.**

### **Config**

If you are using gStore to serve as a SPARQL Endpoint, you had better set the SPARQL\_ENDPOINT macro in Util.h. In addition, all DEBUG macros in Util.h should not be used. What is more, if you do not use gStore as endpoint, but you do not need to update the database, ONLY\_READ macro should be set in Util.h.

### **Backup**

When running as a HTTP server, gStore has provided a backup function, and you can modify the time interval of backup procedure in Util. However, this backup utility can not take effect if the whole disk is down. As a result, you had better use multiple machines or clouds to do other backup procedures by yourself, if you are requiring a high security.

### **Query**

When running as a HTTP server, gStore has set a time limit for query processing, i.e. 1 hour. You can modify this parameter in Util as you wish, but we suggest that the lowest bound is 1 minute.

### **KVstore**

The efficiency of KVstore has huge impact on the performance of the whole system, and you can modify related parameters in KVstore.h, according to your demand and memory capacity.

## **String Buffer**

To speed up the process of reading Strings from disk when running into `getFinalResult()` function of the query processing, gStore has provided String Buffers for entities and literals. You can set this parameters in `setStringBuffer()` of `Database.h` according to the memory capacity of your machine.

## **HTTP API**

If you are using HTTP API, and use Java to visit it, then you must be care for the efficiency(due to the Garbage Collection process in jvm). We strongly suggest that you read codes in `api/http/java` carefully, before you write your own program to use the REST API providing by HTTP protocol.



## Part III

# Others

### Chapter 13: Contributors

Please contact with Lei Zou(zoulei@pku.edu.cn), Li Zeng(zengli-bookug@pku.edu.cn), Jiaqi Chen(chenjiaqi93@pku.edu.cn) and Peng Peng(pku09pp@pku.edu.cn) if you have suggestions or comments about gStore or you need help when using gStore.

#### Faculty

- Lei Zou (Peking University) Project Leader
- M. Tamer Özsu (University of Waterloo)
- Lei Chen (Hong Kong University of Science and Technology)
- Dongyan Zhao (Peking Univeristy)
- Zhiyuan Deng (Wuhan University)

#### Students

*Li Zeng and Jiaqi Chen are responsible for the gStore system optimization. Peng Peng is responsible for the distributed version of gStore, which is expected to be released before October.*

- Shuo Han (Peking University) (PhD student)
- Li Zeng (Peking University) (PhD student)
- Jiaqi Chen (Peking University) (Master student)
- Lin Hu (Peking University) (PhD student)

- Xunbin Su (Peking University) (Master student)
- Jing Li (Peking University) (Master student)

### **Alumni**

- Xuchuan Shen (Peking University) (Master's student, graduated)
- Dong Wang (Peking University) (PhD student, graduated)
- Ruizhe Huang (Peking University) (Undergraduate intern, graduated)
- Jinhui Mo (Peking University) (Master's, graduated)
- Peng Peng (Peking University) (PhD student)
- Youhuan Li (Peking University) (PhD student)
- Libo Wang (Peking University) (Intern student)

## **Chapter 14: Updated Logs**

**Apr 24, 2018**

Multithreading is enabled by Li Zeng in ghttp, to improve the performance of this HTTP web server.

In addition, openmp is added by Xunbin Su for sort and qsort to excavate hidden performance. However, this optimization does not yield good result. As a result, we reserve the code of openmp, but still use standard sort and qsort functions.

Jing Li adds support for multiple users and databases in ghttp, and improve the functions of web server, as well as the SPARQL query endpoint.

New key-value indices are designed by Zongyue Qin to take place of the original B+ trees. In detail, we think the original ISTree and IVTree are not efficient enough, so we choose to implement array+hash method instead of B+ tree. What needs to be noticed is that all these indices does not support parallism now, so when multiple queries are running concurrently, we must add locks to ensure that the indices are visited in sequence. Furthermore, Zongyue Qin designs a new method to compress the original string in RDF dataset. For example, the prefix can be extracted and compressed using some special characters.

Lin Hu fixes the bugs in preFilter and Join functions, which has no impact on the performance of answering SPARQL queries.

**Oct 2, 2017**

Bind and GroupBy are supported in SPARQL queries now, and the Join module has been optimized to harvest a big improvement.

In addition, VSTree module is deprecated to save memory and support larger datasets like freebase, which has almost 2.5B triples. In fact, we have finish experiments on a microbiology dataset with 3.5B triples, and the time of query answering is almost in a scale of 10s.

What's more, we redesign the Database Server using HTTP1.1 protocol and provide REST interface and Java API example now. Users can visit the server by URL

in a browser directly, and we have built several SPARQL endpoints based on it(including freebase, dbpedia, and openKG).

What is not mentioned above is the robustness of Database Server, which supports restart function and query timeout handler(1 hour by default). Backup function is also included in Database Server now, as well as the REDO function to finish a update query which was interrupted by a system crash.

### **Jan 10, 2017**

preFilter() function in Join module is optimized using the pre2num structure, as well as the choose\_next\_node() function. A global string buffer is added to lower the cost of getFinalResult(), and the time of answering queries is reduced greatly.

In addition, we assign buffers of different size for all B+ trees.(some of them are more important and more frequently used) WangLibo merges several B+ trees into one, and the num of all B+ trees are reduced to 9 from 17. This strategy not only reduces the space cost, but also reduces the memory cost, meanwhile speeding up the build process and query process.

What is more, ChenJiaqi has done a lot of work to optimize the SPARQL query. For example, some unconnected SPARQL query graphs are dealt specially.

### **Sep 15, 2016**

ZengLi splits the KVstore into 3 parts according to the types of key and value, i.e. int2string, string2int and string2string. In addition, updates are supported now. You can insert, delete or modify some triples in the gStore database. In fact, only insert() and remove() are implemented, while the modify() are supported by removing first and insert again.

### **Jun 20, 2016**

ZengLi has enabled the gStore to answer queries with predicate variables. In addition, the structures of many queries have been studied to speed up the query

processing. ChenJiaqi rewrites the sparql query plan to acquire a more efficient one, which brings many benefits to us.

### **Apr 01, 2016**

The structure of this project has changed a lot now. A new join method has been achieved and we use it to replace the old one. The test result shows that speed is improved and the memory cost is lower. We also do some change to Parser/Sparql\*, which are all generated by ANTLR. They must be modified because the code is in C, which brings several multiple definition problems, and its size is too large.

There is a bug in the original Stream module, which brings some control characters to the output, such as ^C, ^V and so on. We have fixed it now and enabled the Stream to sort the output strings(both internal and external). In addition, SPARQL queries which are not BGP(Basic Graph Pattern) are also supported now, using the naive method.

A powerful interactive console, which is named `gconsole` now, is achieved to bring convenience to users. What is more, we use `valgrind` tools to test our system, and deal with several memory leaks.

The docs and API have also changed, but this is of little importance.

### **Nov 06, 2015**

We merge several classes(like Bstr) and adjust the project structure, as well as the debug system.

In addition, most warnings are removed, except for warnings in Parser module, which is due to the use of ANTLR.

What is more, we change RangeValue module to Stream, and add Stream for ResultSet. We also better the gquery console, so now you can redirect query results to a specified file in the gsql console.

Unable to add Stream for IDlist due to complex operations, but this is not necessary. Realpath is used to supported soft links in the gquery console, but it not works in

Gstore.(though works if not in Gstore)

### **Oct 20, 2015**

We add a gtest tool for utility, you can use it to query several datasets with their own queries.

In addition, gquery console is improved. Readline lib is used for input instead of fgets, and the gquery console can support commands history, modifying command and commands completion now.

What is more, we found and fix a bug in Database/(a pointer for debugging log is not set to NULL after fclose operation, so if you close one database and open another, the system will fail entirely because the system think that the debugging log is still open)

### **Sep 25, 2015**

We implement the version of B+Tree, and replace the old one.

After testing on DBpedia, LUBM, and WatDiv benchmark, we conclude that the new BTree performs more efficient than the old version. For the same triple file, the new version spends shorter time on executing gload command.

Besides, the new version can handle the long literal objects efficiently, while triples whose object's length exceeds 4096 bytes result in frequent inefficient split operations on the old version BTree.

### **Feb 2, 2015**

We modify the RDF parser and SPARQL parser.

Under the new RDF parser, we also redesign the encode strategy, which reduces RDF file scanning times.

Now we can parse the standard SPARQL v1.1 grammar correctly, and can support basic graph pattern(BGP) SPARQL queries written by this standard grammar.

**Dec 11, 2014**

We add API for C/CPP and JAVA.

**Nov 20, 2014**

We share our gStore2.0 code as an open-source project under BSD license on github.

## Chapter 15: Test Result

### Preparation

We have compared the performance of gStore with several other database systems, such as [Jena](#), [Sesame](#), [Virtuoso](#) and so on. Contents to be compared are the time to build database, the size of the built database, the time to answer single SPARQL query and the matching case of single query's results. In addition, if the memory cost is very large(>20G), we will record the memory cost when running these database systems.(not accurate, just for your reference)

To ensure all database systems can run correctly on all datasets and queries, the format of datasets must be supported by all database systems and the queries should not contain update operations, aggregate operations and operations related with uncertain predicates. Notice that when measuring the time to answer queries, the time of loading database index should not be included. To ensure this principle, we load the database index first for some database systems, and warm up several times for others.

Datasets used here are WatDiv, Lubm, Bsbm and DBpedia. Some of them are provided by websites, and others are generated by algorithms. Queries are generated by algorithms or written by us. Table 3 summarizes the statistics of these datasets. The experiment environment is a CentOS server, whose memory size is 82G and disk size is 7T. We use [full\\_test](#) to do this test.

Dataset	Number of Triples	RDF N3 File Size(B)	Number of Entities
WatDiv 300M	329,539,576	47,670,221,085	15,636,385
LUBM 5000	66718642	8134671485	16437950
DBpedia 2014	170784508	23844158944	7123915
Bsbm 10000	34872182	912646084	526590

Table 3: Datasets



## Result

The performance of different database management systems is shown in Figures 5, 6, 7 and 8.

Notice that Sesame and Virtuoso are unable to operate on DBpedia 2014 and WatDiv 300M, because the size is too large. In addition, we do not use Sesame and Virtuoso to test on the LUBM 5000 due to format questions. Generally speaking, Virtuoso is not scalable, and Sesame is so weak.

This program produces many logs placed in result.log/, load.log/ and time.log/. You can see that all results of all queries are matched by viewing files in result.log/, and the time cost and space cost of gStore to build database are larger than others by viewing files in load.log/. More precisely, there is an order of magnitude difference between gStore and others in the time/space cost of building database.

Through analysing time.log/, we can find that gStore behave better than others on very complicated queries(many variables, circles, etc). For other simple queries, there is not much difference between the time of these database systems.

Generally speaking, the memory cost of gStore when answering queries is higher than others. More complicated the query is and more large the dataset is, more apparent the phenomenon is.

You can find more detailed information in [original test report](#). Notice that some questions in the test report have already be solved now. The latest test report is [formal experiment](#).

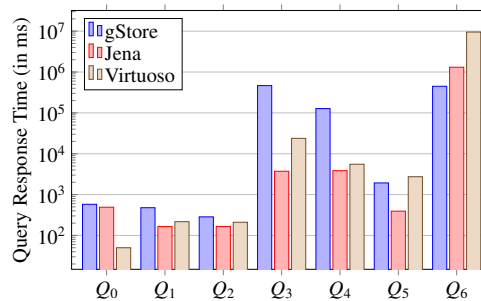


Figure 5: Query Performance over DBpedia 2014

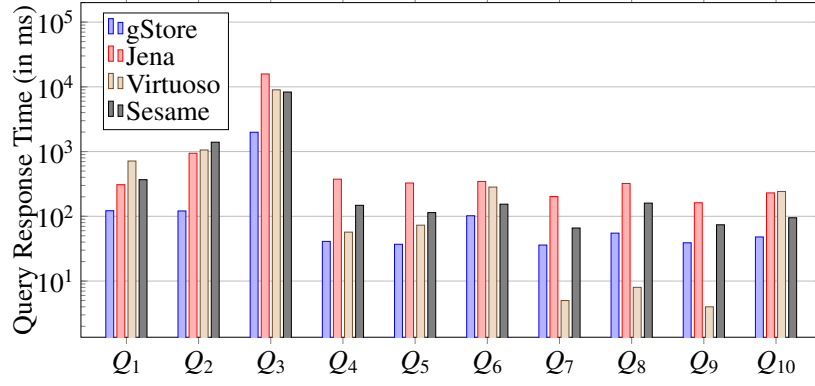
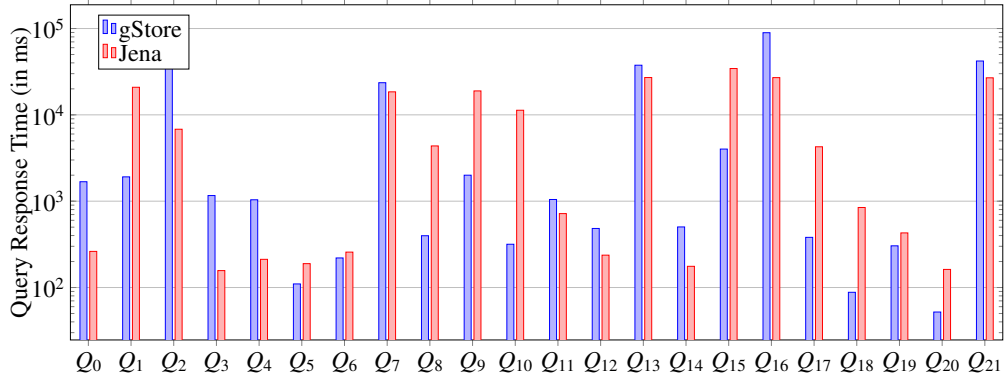
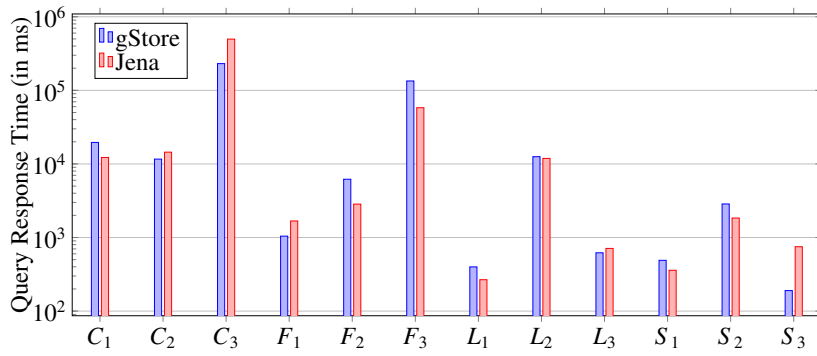


Figure 6: Query Performance over Bsbm 10000



(a) LUBM 5000

Figure 7: Query Performance over LUBM



(a) WatDiv 300M

Figure 8: Query Performance over WatDiv

## **Chapter 16: Future Plan**

### **Improve The Core**

- speed up the join process and postprocessing of SPARQL using GPU or FPGA
- improve the indices and support concurrent reads
- add numeric value query function. need to answer numeric range query efficiently and space consume cannot be too large
- typedef all frequently used types, to avoid inconsistency and high modify cost

### **Better The Interface**

- the usability of ghttp(ERROR\_CODE, API ...)
- improve socket interface
- docker settings

### **Idea Collection Box**

- warnings remain in using Parser/(antlr)!(modify sparql.g 1.1 and regenerate). change name to avoid redefine problem, or go to use executable to parse
- mmap to speedup KVstore?
- the strategy for Stream:is 85% valid? consider sampling, analyse the size of result set and decide strategy? how to support order by: sort in memory if not put in file; otherwise, partial sort in memory, then put into file, then proceed external sorting

## Chapter 17: Thanks List

*This chapter lists people who inspire us or contribute to this project.*

### **zhangxiaoyang**

- add python socket API (api/socket/python)

### **Dingfeng Wang**

fei123581321@qq.com

- support python3 in socket API (api/socket/python3)

### **Libo Wang**

wlbqe@pku.edu.cn

- suppress the num of B+ Tree
- add backup function to gserver and ghttp
- add REDO function to gserver and ghttp (in case of interrupt during a update query)
- restart automatically if database server is down
- kill a query answering procedure if it runs more than 1 hour

### **Xin Lv**

lvxin1204@163.com

- provide a HTTP server for database

**Zhiyuan Deng**

331563360@qq.com

- provide a parallelable version of B+ Tree

**Hao Cui**

prospace@bupt.edu.cn

- provide a query-level cache for database server

**imbajin**

- provide the support of docker deployment

## **Chapter 18: Legal Issues**

Copyright (c) 2016 gStore team

All rights reserved.

Redistribution and use in source and binary forms, with or without modification, are permitted provided that the following conditions are met:

Redistributions of source code must retain the above copyright notice, this list of conditions and the following disclaimer.

Redistributions in binary form must reproduce the above copyright notice, this list of conditions and the following disclaimer in the documentation and/or other materials provided with the distribution.

Neither the name of the Peking University nor the names of its contributors may be used to endorse or promote products derived from this software without specific prior written permission.

THIS SOFTWARE IS PROVIDED BY THE COPYRIGHT HOLDERS AND CONTRIBUTORS "AS IS" AND ANY EXPRESS OR IMPLIED WARRANTIES, INCLUDING, BUT NOT LIMITED TO, THE IMPLIED WARRANTIES OF MERCHANTABILITY AND FITNESS FOR A PARTICULAR PURPOSE ARE DISCLAIMED. IN NO EVENT SHALL THE COPYRIGHT OWNER OR CONTRIBUTORS BE LIABLE FOR ANY DIRECT, INDIRECT, INCIDENTAL, SPECIAL, EXEMPLARY, OR CONSEQUENTIAL DAMAGES (INCLUDING, BUT NOT LIMITED TO, PROCUREMENT OF SUBSTITUTE GOODS OR SERVICES; LOSS OF USE, DATA, OR PROFITS; OR BUSINESS INTERRUPTION) HOWEVER CAUSED AND ON ANY THEORY OF LIABILITY, WHETHER IN CONTRACT, STRICT LIABILITY, OR TORT (INCLUDING NEGLIGENCE OR OTHERWISE) ARISING IN ANY WAY OUT OF THE USE OF THIS SOFTWARE, EVEN IF ADVISED OF THE POSSIBILITY OF SUCH DAMAGE.

What's more, you need to include the label "powered by gStore", as well as the

logo of gStore, in your software product which is using gStore.

We would be very grateful if you are willing to tell us about your name, institution, purpose and email. Such information can be sent to us by emailing to [gStoreDB@gmail.com](mailto:gStoreDB@gmail.com), and we promise not to reveal privacy.

**End**

**Thank you for reading this document. If any question or advice, or you have interests in this project, please don't hesitate to get in touch with us.**