

# **Выпускная квалификационная работа Аналитика исходного кода ClickHouse с помощью ClickHouse**

---

Ширин Никита, гр. 155

Научный руководитель: Руководитель группы разработки ClickHouse в Яндексе,  
Миловидов Алексей Николаевич

# Предметная область

- Извлечение различных статистик из Git-репозитория
- Использование аналитической СУБД для хранения данных о репозитории

# Актуальность задачи

- Извлеченные из данных репозитория статистики полезны для аналитиков, исследующих различные аспекты работы над проектом, и разработчикам
- Существующие решения не обладают достаточной гибкостью и удобством использования
- Желаящим извлечь нетривиальную статистику приходится писать специальные программы

# Цель и задачи дипломной работы

Цель: получить инструмент, позволяющий извлекать нетривиальные статистики посредством SQL-запросов

Задачи:

- Обзор существующих решений - GitHub API и GitStats
- Разработать архитектуру базы данных под управлением СУБД ClickHouse
- Выбрать архитектуру и реализовать программу, наполняющую базу данных
- Выбрать архитектуру и реализовать web-приложение, обращающееся к базе данных
- Провести эксперименты подсчета нетривиальных статистик на примере репозитория ClickHouse

# Формальная постановка

Разработка приложения, получающего на вход git-репозиторий, сохраняющего данные о репозитории в базу данных под управлением СУБД ClickHouse и предоставляющего пользователю интерфейс, позволяющий задавать SQL-запросы к базе данных для подсчета нетривиальных статистик.

# Обзор существующих решений

## GitHub API

- REST API
- Кэширование ответов
- Небольшое число поддерживаемых запросов

Get the number of additions and deletions per week



```
GET /repos/:owner/:repo/stats/code_frequency
```

### Response

Returns a weekly aggregate of the number of additions and deletions pushed to a repository.

Status: 200 OK

```
[  
  [  
    1302998400,  
    1124,  
    -435  
  ]  
]
```

# Обзор существующих решений

## GitStats

- Приложение на python
- HTML-отчет
- Для расширения функционала требуется модификация исходного кода

Authors				
General	Activity	Authors	Files	Lines
List of Authors				
Author	Commits (%)	+ Files	- Files	First c
Alexandre Julliard	11630 (13.21%)	4063176	2925213	1994
Jack Caban	4268 (4.85%)	314444	124195	2003
Francois Gouget	3494 (3.97%)	287371	266141	1998
Henri Verbeet	3341 (3.79%)	186780	152227	2008
Juan Lang	2489 (2.83%)	178793	51327	2003
Michael Stefancu	2270 (2.58%)	72234	76521	2001
Eric Pouech	2035 (2.31%)	299284	189390	1998
Hans Leidekker	2012 (2.29%)	145810	53693	2003
Stefan Dörsinger	1926 (2.19%)	117234	67793	2004
Mike McCormack	1853 (2.10%)	129791	57853	2000
Rob Shearman	1808 (2.05%)	76147	30643	2004
Nikolay Sivov	1700 (1.93%)	82600	33977	2008
Paul Vriens	1696 (1.93%)	43681	13421	2004
Dmitry Timoshkov	1678 (1.91%)	109469	61747	1999
Hew Davies	1637 (1.86%)	128279	25880	2002
Marcus Meissner	1602 (1.82%)	79278	52573	1998
James Hawkins	1552 (1.76%)	102560	27200	2004
Robert Shearman	1371 (1.56%)	88607	27544	2002
Andrew Talbot	1353 (1.54%)	9533	10810	2006-05-23
Aric Stewart	1349 (1.53%)	112464	42711	2000-05-05

These didn't make it to the top: Piotr Caban, Vincent Povirk, André Hentschel, Maarten Lankhorst, Detlef Riekenberg, Christian Costa, Dimitrie O. Paun, Austin English, Vitaliy Margolen, Gerald Pfeifer, Roderick Colsonbrander, Dylan Smith, Alastair Leslie-Hughes, Robert Rief, H. Verbeet, Frédéric Delany, Andreas Mohr, Andrew Nguyen, Patrik Stridvall, Rico Schüller, Owe Kaaven, Mikolaj Zalewski, Andrew Eikam, Rein Kloos, Lionel Ulmer, Ulrich Weigand, Joerg Schmidt, Jon Griffiths, Stefan Leichter, Jason Edmendes, How D M Davies, Jörg Hühle, Alexander Nicolayson Sornes, Evan Stadel, Uwe Bonnes, Lei Zhang, Matteo Bruni, Misha Konshelov, Vincent Béron, Dan Hipschman, Thomas Mullaly, Ken Thomases, David Adam, David Hedberg, Damjan Jovanovic, Michael Jung, Alex Villacis Lasso, Mike Hearn, Ivan Gyurdiev, Ulrich Czekalla, Gerard Patel, Louis Lenders, Hwang YunSong, Aurimas Fierbas, Kai Blin, Owen Rudge, Luca Bannati, Dan Kegel, Raphael Jaquesira, Jeremy White, Chris Robinson, Ge van Geldorp, Martin Fuchs, Alexander Doroshev, Jukka Heino, Frank Richter, Hwang YunSong (황은성), Steven Edwards, Andrey Turkin, Oliver Stieber, Eric Kohl, Guy L. Albertelli, Hidenori Takekuma, Akhilesh Saurava, Bill Medland, Kevin Kellman, Rok Mandelc, Jonathan Ernst, Ricardo Filipe, Igor Palychuk, Vladimir Pankratov, Reif Kallenberg, Jeff Lattimer, Markus Ammer, Sasilm Krauss, Phil Krylov, Erich

## GitStats - wine

General

Activity

Authors

Files

Lines

Tags

### Project name:

wine

### Generated:

2012-05-26 05:48:03 (in 2944 seconds)

### Generator:

[GitStats](#) (version 81045f5), git version 1.7.8.247.g10f4e, gnuplot 4.4 patchlevel 0

### Report Period:

1993-06-29 19:33:12 to 2012-05-25 21:21:15

### Age:

6906 days, 4111 active days (59.53%)

### Total Files:

5885

### Total Lines of Code:

4225178 (9882574 added, 5657396 removed)

### Total Commits:

88038 (average 21.4 commits per active day, 12.7 per all days)

### Authors:

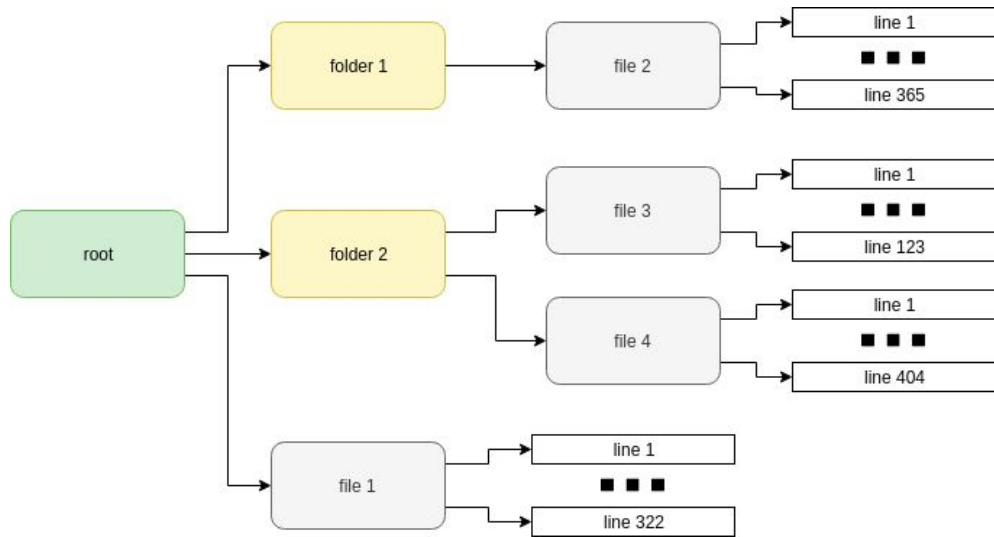
1343 (average 65.6 commits per author)

# Архитектура приложения

## База данных

### Первый вариант

- Независимое заполнение для различных коммитов
- Множество операций с жестким диском
- Оптимизация для хранения



Структура таблицы: commit\_hash (FixedString(20)), commit\_date (Date), commit\_time (DateTime), commit\_message (String), author\_name (String), author\_email (String), file\_name (String), file\_path (String), file\_extension (String), line\_num (UInt32), line (String)



# Архитектура приложения

## База данных

### Второй вариант

- Имитация репозитория python-структурами
- Хранение “blame”
- Две таблицы.
- Оптимизация для скорости запросов

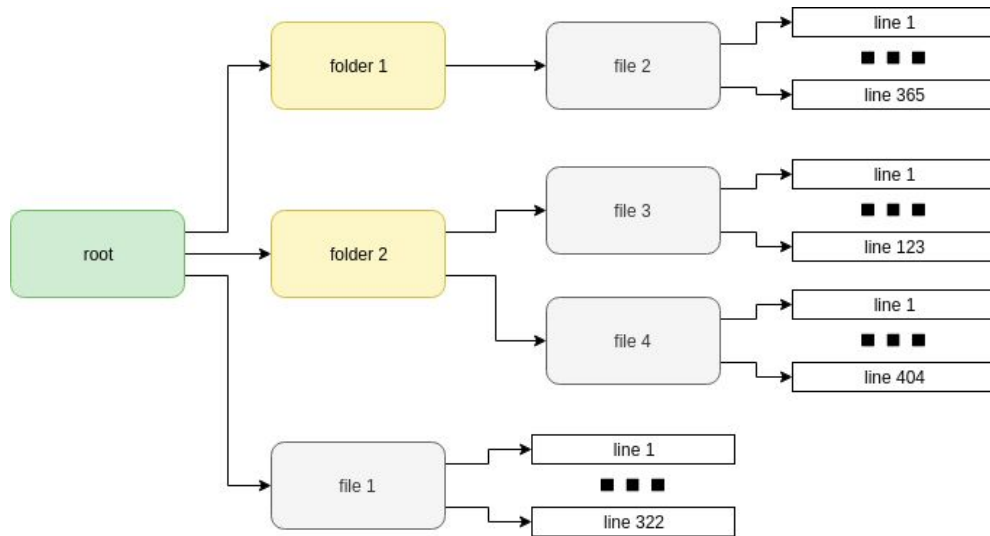


Таблица lines: commit\_hash (FixedString(40)), file\_name (String), file\_path (String), file\_id (UUID), line\_num (UInt32), line (String), last\_change (FixedString(40))

Таблица commits: commit\_hash (FixedString(40)), commit\_time (DateTime), commit\_message (String), author\_name (String), author\_email (String)

# Архитектура приложения

## База данных

## Результаты для репозитория ClickHouse

Вариант	Первый вариант	Второй вариант
Размер таблицы	600 Гб	200 Гб
Время заполнения	10 стуюк	6 суток
Время выполнения запросов	1000 секунд	Сильно быстрее

Характеристики компьютера: Intel Core i7-7700, 16 Гб RAM и накопитель типа HDD

Код написан на Python3 с использованием библиотек git, pydriller, clickhouse-driver

# Web-приложение

- Серверная часть написана на языке Python3 с использованием Flask
- PostgreSQL для хранения истории SQL-запросов
- Отдельная программа выполняет невыполненные запросы
- Графический интерфейс написан на TypeScript с использованием React

Analytical queries

QUERY LIST

Description

Query

SUBMIT

Filter

Description	Query	Executed	Added
<a href="#">Codebase size</a>	<pre>select * from (select commit_hash, commit_time from commits order by commit_hash limit 100) commits inner join (select commit_hash, count() as n_lines from lines group by commit_hash) lines on commits.commit_hash=lines.commit_hash order by commit_time;</pre>	True	28/05/2019, 00:18:18
<a href="#">Codebase size</a>	<pre>select * from (select commit_hash, commit_time from commits order by commit_hash limit 100) commits inner join (select commit_hash, count() as n_lines from lines group by commit_hash) lines on commits.commit_hash=lines.commit_hash;</pre>	True	28/05/2019, 00:16:29
<a href="#">Codebase size</a>	<pre>select * from (select commit_hash, commit_time from commits order by commit_hash limit 100) commits inner join (select commit_hash, count() as n_lines from lines group by commit_hash) lines;</pre>	Fail	28/05/2019, 00:15:44
<a href="#">Dangerous files</a>	<pre>select fpath, argMax(author_name, n_lines) as author, max(n_lines) / sum(n_lines) as fraction from (select concat(lines.file_path, '/') as fpath, lines.file_name) as fpath, commits.author_name, count(*) as n_lines from (select * from lines where commit_hash =( select commit_hash from commits where commit_time =( select max(commit_time) from commits limit 1))) lines inner join commits on lines.last_change = commits.commit_hash group by fpath, commits.author_name ) group by fpath order by fraction desc;</pre>	True	28/05/2019, 00:07:17

# Эксперименты

Analytical queries QUERY LIST

## Commits by int weekday

```
select weekday, count() as n_commits from (select cast(formatDateTime(commit_time, '%w') as Int8) as weekday from commits) group by weekday order by weekday;
```

Elapsed 1 secs.

X

Y

GET GRAPH

weekday (Int8)

n\_commits (UInt64)

0	1753
1	4163
2	3851
3	4363
4	4215
5	4307
6	1648

Analytical queries QUERY LIST

## Commits by int weekday

```
select weekday, count() as n_commits from (select cast(formatDateTime(commit_time, '%w') as Int8) as weekday from commits) group by weekday order by weekday;
```

Elapsed 1 secs.

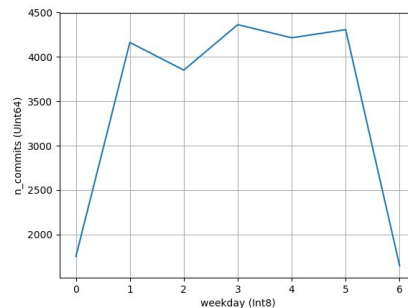
X

weekday

Y

n\_commits

GET GRAPH



# Эксперименты

Analytical queries QUERY LIST

## Codebase size

```
select * from (select commit_hash, commit_time from commits order by commit_hash limit 100) commits inner join (select commit_hash, count() as n_lines from lines group by commit_hash) lines on commits.commit_hash=lines.commit_hash order by commit_time;
```

Elapsed 46 secs.

X

Y

GET GRAPH

commit_hash (FixedString(40))	commit_time (DateTime)	lines.commit_hash (FixedString(40))	n_lines (UInt64)
00f23ee493b547c2f62922491b404ed61aceab4c	2012-03-05 08:29:16	00f23ee493b547c2f62922491b404ed61aceab4c	34337
002d6a617850a2e1ca3c154df94403d52b83f66e	2013-01-07 12:27:39	002d6a617850a2e1ca3c154df94403d52b83f66e	59949
00e0307de17c74f3ead7cc948b43e70dd18f5925	2013-03-14 17:04:50	00e0307de17c74f3ead7cc948b43e70dd18f5925	64864
00f2ccea0218e3b08a60b17182e2d6f0fa7a050	2013-07-12 21:47:19	00f2ccea0218e3b08a60b17182e2d6f0fa7a050	70912
00b10d30df41b5b1f79ad65e607c30ba9e712293	2014-03-11 11:32:47	00b10d30df41b5b1f79ad65e607c30ba9e712293	93943
00d9c285719981a7fcd1dc5ee7d820cb59674771	2014-03-13 21:44:00	00d9c285719981a7fcd1dc5ee7d820cb59674771	94071
003c7b30cdf118e4b079134273c4762929bb7274	2014-03-14 19:42:30	003c7b30cdf118e4b079134273c4762929bb7274	95670
00d5cfbe45acd80ae90f9ef687059b644dac206	2014-04-03 19:57:36	00d5cfbe45acd80ae90f9ef687059b644dac206	95997
00074ffaa5d2282e824efec2fa0882dade66f717	2014-07-03 15:22:12	00074ffaa5d2282e824efec2fa0882dade66f717	110977

Analytical queries QUERY LIST

## Codebase size

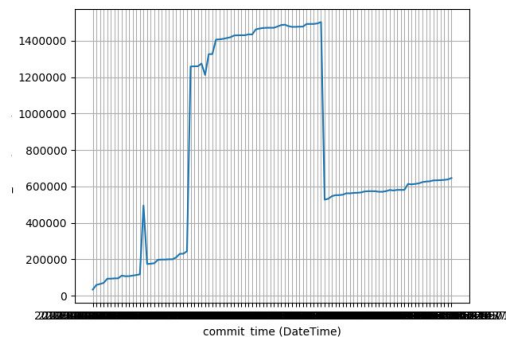
```
select * from (select commit_hash, commit_time from commits order by commit_hash limit 100) commits inner join (select commit_hash, count() as n_lines from lines group by commit_hash) lines on commits.commit_hash=lines.commit_hash order by commit_time;
```

Elapsed 46 secs.

X  
time

Y  
n\_lines

GET GRAPH



# Эксперименты

Analytical queries QUERY LIST

### Top 10 renames

```
select file_id, length(groupArray(fpath)) as n_renames, arrayStringConcat(groupArray(fpath), ', ') as names from
(select distinct file_id, concat(file_path, '/', file_name) as fpath from lines) group by file_id order by n_renames
desc limit 10;
```

Elapsed 515 secs.

X

Y

GET GRAPH

file_id (UUID)	n_renames (UInt64)	names (String)
d0d1166b-69fe-4c93-9701-4900b5573e9	5	dbms/tests/queries/0_stateless/00534_functions_bad_arguments12.reference, dbms/tests/queries/0_stateless/00534_long_functions_bad_arguments12.reference, dbms/tests/queries/0_stateless/00534_long_functions_bad_arguments5.reference, dbms/tests/queries/0_stateless/00534_long_functions_bad_arguments4.reference, dbms/tests/queries/0_stateless/00534_functions_bad_arguments4.reference
bb2461d6-d618-4e46-9b69-17e8314d04b7	5	dbms/tests/queries/0_stateless/00534_functions_bad_arguments2.reference, dbms/tests/queries/0_stateless/00534_long_functions_bad_arguments2.reference, dbms/tests/queries/0_stateless/00534_long_functions_bad_arguments7.reference, dbms/tests/queries/0_stateless/00534_long_functions_bad_arguments6.reference, dbms/tests/queries/0_stateless/00534_functions_bad_arguments6.reference
fa3718c9-4464-4025-911a-61138358f3a	5	dbms/tests/queries/0_stateless/00534_functions_bad_arguments1.reference, dbms/tests/queries/0_stateless/00534_long_functions_bad_arguments1.reference, dbms/tests/queries/0_stateless/00534_functions_bad_arguments.reference, dbms/tests/queries/0_stateless/00534_long_functions_bad_arguments2.reference
02467954-c2fa-47b2-810b-cc0271f8437	5	dbms/tests/queries/0_stateless/00534_functions_bad_arguments13.reference, dbms/tests/queries/0_stateless/00534_long_functions_bad_arguments13.reference, dbms/tests/queries/0_stateless/00534_long_functions_bad_arguments6.reference, dbms/tests/queries/0_stateless/00534_long_functions_bad_arguments5.reference, dbms/tests/queries/0_stateless/00534_functions_bad_arguments5.reference

# Эксперименты

Analytical queries QUERY LIST

## Dangerous files

```
select fpath, argMax(author_name, n_lines) as author, max(n_lines) / sum(n_lines) as fraction from (select concat(lines.file_path, '/',  
lines.file_name) as fpath, commits.author_name, count(*) as n_lines from (select * from lines where commit_hash =( select commit_hash from  
commits where commit_time =( select max(commit_time) from commits limit 1))) lines inner join commits on lines.last_change =  
commits.commit_hash group by fpath, commits.author_name ) group by fpath order by fraction desc;
```

Elapsed 0 secs.

X

Y

GET GRAPH

fpath (String)	author (String)	fraction (Float64)
dbms/tests/queries/0_stateless/00586_removing_unused_columns_from_subquery.sql	alexey-milovidov	1
dbms/tests/queries/0_stateless/00711_array_enumerate_variants.sql	Alexey Milovidov	1
dbms/src/Functions/bitXor.cpp	chertus	1
dbms/tests/queries/0_stateless/00743_limit_by_not_found_column.reference	Alexey Zatelepin	1
dbms/tests/queries/0_stateless/00721_force_by_identical_result_after_merge_zookeeper.reference	Alexey Zatelepin	1
dbms/programs/performance-test/clickhouse-performance-test.cpp	Anastasiya Tsarkova	1
dbms/tests/queries/0_stateless/00333_parser_number_bug.sql	proller	1
dbms/src/Functions/bitTest.cpp	chertus	1
dbms/src/Functions/bitTest.cpp	Alexey Milovidov	1

# Эксперименты

Analytical queries QUERY LIST

Rewrite frags 2019 having at least 5 commits with filenames 2nd try

```
select fpath, mean_frac from (select file_id, AVG(frac) as mean_frac from (select commit_hash, file_id, countIf(commit_hash=last_change) / count() as frac from lines where commit_hash in (select commit_hash from commits where toYear(commit_time) >= 2019) group by commit_hash, file_id) group by file_id having count(commit_hash)>= 5 order by mean_frac desc) inner join (select distinct file_id, concat(file_path, '/', file_name) as fpath from lines where commit_hash in (select commit_hash from commits where commit_time=(select max(commit_time) from commits))) on file_id=file_id;
```

Elapsed 61 secs.

X

Y

GET GRAPH

fpath (String)

mean\_frac (Float64)

dbms/tests/queries/0_stateless/00933_ttl_simple.sql	0.11192450447507236
cmake/find_ifalloc.cmake	0.10869565217391307
dbms/tests/queries/0_stateless/00933_reserved_word.sql	0.0971428571428571
utils/list_backports.sh	0.08967391304347826
dbms/tests/queries/0_stateless/00933_ttl_with_default.sql	0.08948106591865358
dbms/tests/queries/0_stateless/00931_low_cardinality_nullable_aggregate_function_type.sql	0.08583333333333333
utils/github/query.py	0.08466249212350348
dbms/src/Interpreters/RowRefs.cpp	0.06500206691711478



# Результаты работы

- Проведен обзор существующих решений - GitHub API и GitStats
- Выбранная архитектура базы данных поддерживает функционал GitHub API и GitStats, а также предоставляет возможность считать пользовательские статистики посредством SQL-запросов
- Программа, наполняющая базу данных написана на языке Python
- Web-приложение написано на языке TypeScript с React, серверная часть - на языке Python с Flask
- В экспериментах продемонстрированы нетривиальные запросы, реализуемые с помощью приложения, на примере репозитория ClickHouse

Спасибо за внимание