

Выпускная квалификационная работа

Wait-free каталог баз данных в ClickHouse

Токмаков Александр Викторович, группа БПМИ165

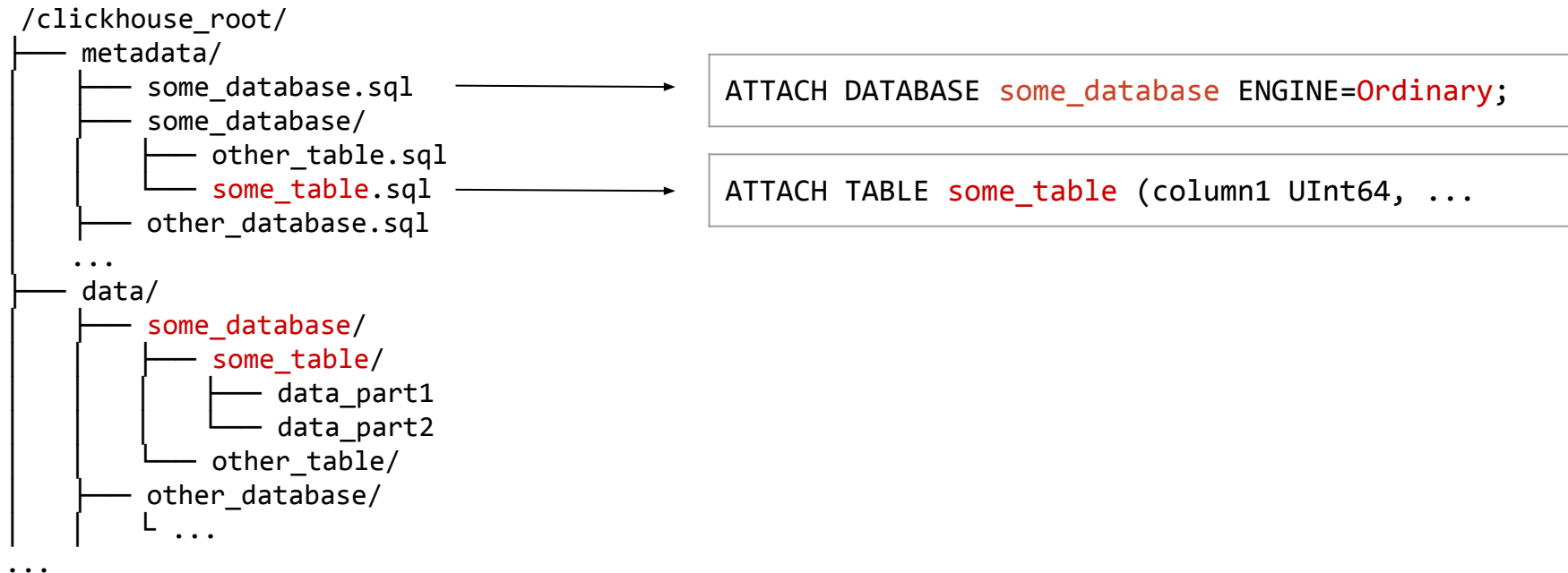
Научный руководитель: Миловидов Алексей Николаевич

Предметная область

- ClickHouse - аналитическая СУБД
- Каталог баз данных отвечает за хранение метаданных
- DDL-запросы требуют синхронизации

```
CREATE TABLE [IF NOT EXISTS] db.table ...;  
ATTACH TABLE [IF NOT EXISTS] db.table ...;  
DROP TABLE [IF EXISTS] db.table;  
DETACH TABLE [IF EXISTS] db.table;  
RENAME TABLE db.table1 TO db.table2;
```

Устройство каталога баз данных в ClickHouse



Синхронизация запросов

DDLGuard:

- С именем таблицы ассоциирован `std::mutex`
- DDL-запросы блокируют имя таблицы

Синхронизация запросов

DDLGuard:

- С именем таблицы ассоциирован `std::mutex`
- DDL-запросы блокируют имя таблицы

RWLock:

- Относится к объекту таблицы
- SELECT и INSERT блокируют таблицу на чтение
- DROP (DETACH) и RENAME блокируют таблицу на запись

Синхронизация запросов

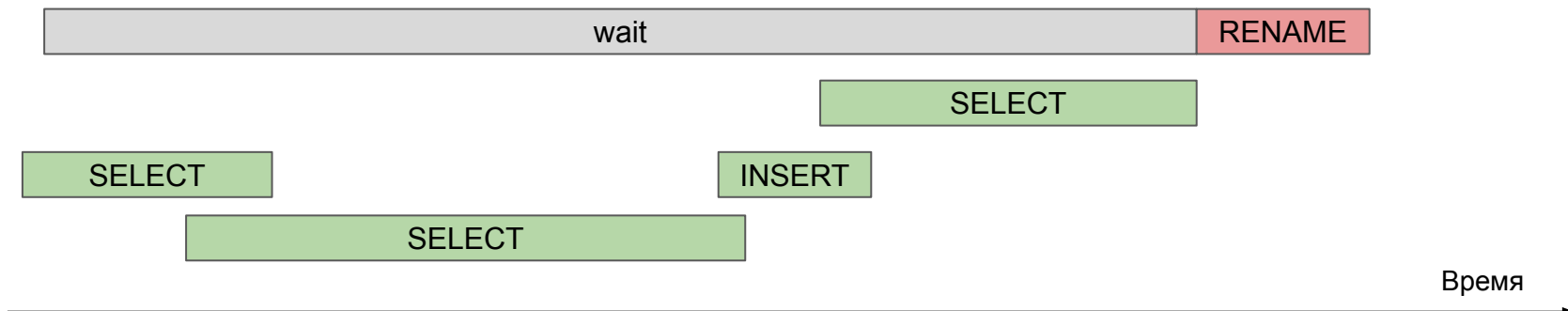
Специфическая реализация RWLock

- рекурсивный; блокируется запросом, а не потоком

Синхронизация запросов

Специфическая реализация RWLock

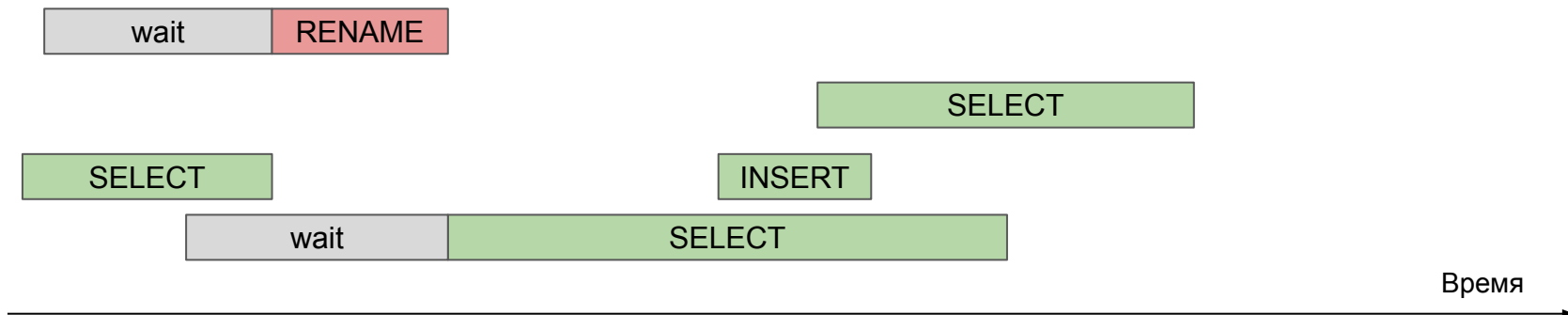
- рекурсивный; блокируется запросом, а не потоком



Синхронизация запросов

Специфическая реализация RWLock

- рекурсивный; блокируется запросом, а не потоком
- честный



Синхронизация запросов

Специфическая реализация RWLock

- рекурсивный; блокируется запросом, а не потоком
- честный
- механизм предотвращения взаимных блокировок

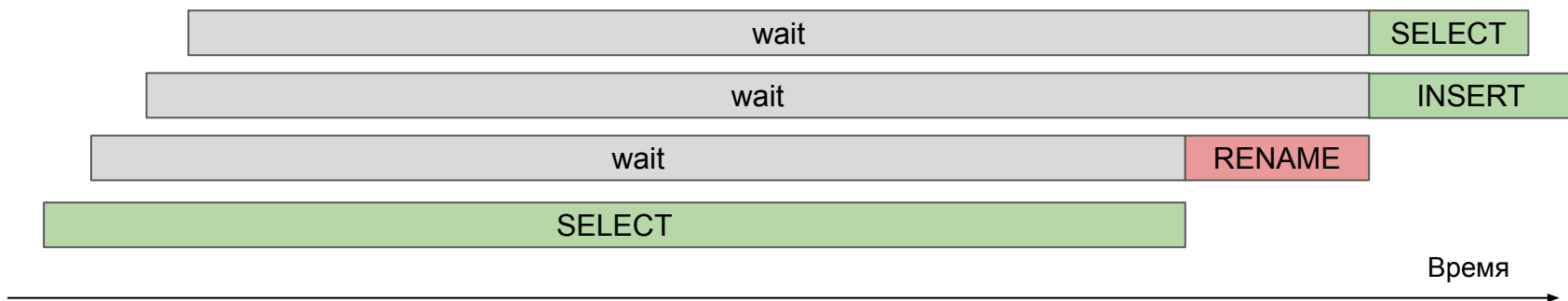
Синхронизация запросов

Специфическая реализация RWLock

- рекурсивный; блокируется запросом, а не потоком
- честный
- механизм предотвращения взаимных блокировок
- много ложноположительных срабатываний

Синхронизация запросов

Иногда честности блокировки недостаточно:



Недостатки существующего подхода

- `DB::Exception: Possible deadlock avoided. Client should retry.`
- Возможно длительное ожидание при выполнении DROP и RENAME
- Неатомарное переименование таблицы

Задачи работы

Реализовать движок баз данных Atomic, который:

- Может заменить движок Ordinary
- Не использует механизм табличных RWLock
- Выполняет DROP и RENAME не допуская длительного ожидания

Общее описание решения

В базе данных с движком Atomic:

- Таблицы имеют постоянный UUID
- Имена отображаются на идентификаторы
- Подсчёт ссылок на таблицы
- Фоновое удаление неиспользуемых таблиц

Изменения в структуре директорий



```
RENAME TABLE db.old_name TO db.new_name;
```

- Атомарно переименовывает файл метаданных
- Данные не перемещаются
- Нет длительных блокировок
- Устойчиво к `kill -9`

Множественное переименование

Частый случай использования:

```
RENAME TABLE table TO table_old, table_new TO table;
```

Альтернатива:

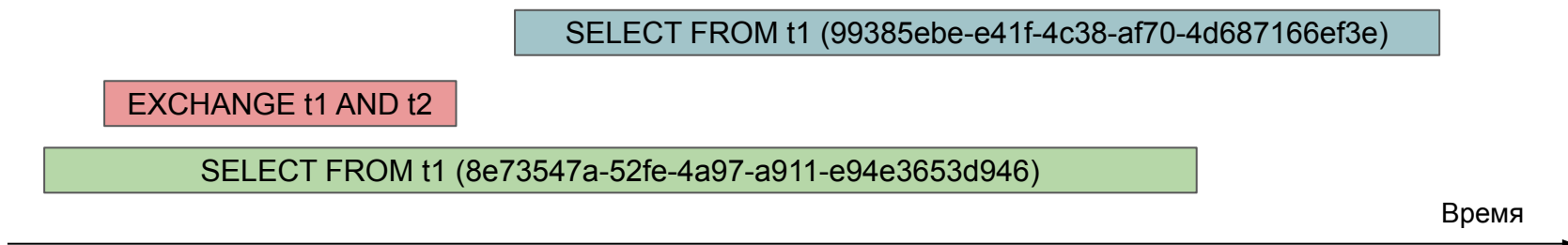
```
EXCHANGE TABLES table AND table_new;
```

Атомарно обменивает файлы метаданных (renameat2) и обновляет структуры данных в памяти

Множественное переименование

Альтернатива:

```
EXCHANGE TABLES table AND table_new;
```



```
DROP TABLE db.table_name;
```

- Атомарно перемещает файл метаданных в `metadata_dropped/`
- Отвязывает UUID от имени таблицы

```
DROP TABLE db.table_name;
```

- Атомарно перемещает файл метаданных в `metadata_dropped/`
- Отвязывает UUID от имени таблицы
- Данные не удаляются
- Запросы продолжают выполняться

A diagram illustrating the execution of a DROP statement. A horizontal timeline is shown with an arrow pointing to the right, labeled 'Время' (Time). Two rectangular boxes are positioned above the timeline. The first box is red and contains the text 'DROP t1'. The second box is green and contains the text 'SELECT FROM t1 (8e73547a-52fe-4a97-a911-e94e3653d946)'. The green box starts after the red box and extends further to the right, indicating that the SELECT query continues to execute even after the DROP statement has been processed.

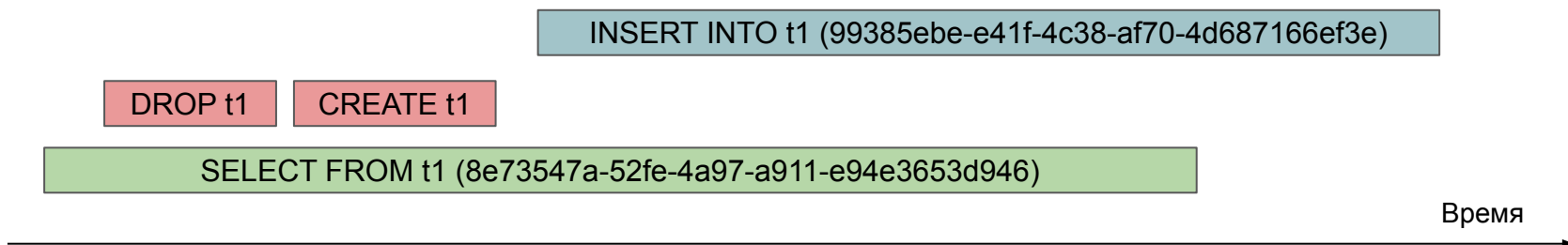
DROP t1

SELECT FROM t1 (8e73547a-52fe-4a97-a911-e94e3653d946)

Время

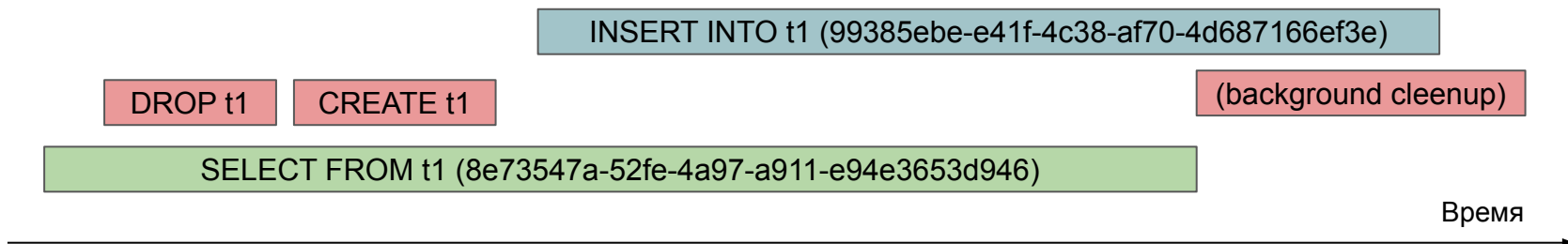
```
DROP TABLE db.table_name;
```

- Атомарно перемещает файл метаданных в `metadata_dropped/`
- Отвязывает UUID от имени таблицы
- Данные не удаляются
- Запросы продолжают выполняться
- Новая таблица может быть создана сразу



```
DROP TABLE db.table_name;
```

- Атомарно перемещает файл метаданных в `metadata_dropped/`
- Отвязывает UUID от имени таблицы
- Данные не удаляются
- Запросы продолжают выполняться
- Новая таблица может быть создана сразу
- Фоновый поток удаляет неиспользуемые данные



Результаты

- В движке Atomic устранены существенные недостатки движка Ordinary
- Доступно в ClickHouse с версии 20.4 (в экспериментальном режиме)

Вопросы