



Платформа инкрементальных вычислений на базе MapReduce

Федор Бочаров

Научный руководитель: Александр Юрьевич Юрченко, Яндекс
САНКТ-ПЕТЕРБУРГСКИЙ АКАДЕМИЧЕСКИЙ УНИВЕРСИТЕТ

13 июня 2017 г.



MapReduce¹ – парадигма программирования, которая позволяет удобно обрабатывать большие объемы данных.

Идея: отображаем данные в пары <ключ, значение> (map), затем агрегируем значения с одинаковыми ключами (reduce).

Особенности:

- Позволяет обрабатывать произвольно большие объемы данных
- Обладает большой латентностью

¹<http://bit.ly/2sfzbO8>



YТ² – платформа распределенных вычислений, поддерживающая парадигму MapReduce. Разработана в компании Яндекс.

Обеспечивает:

- Надежное выполнение операций на кластере
- Надежное хранение данных в *таблицах*
- Потабличные транзакции

²<http://bit.ly/2raQKvs>



Типичный MR-процесс при обновлениях обрабатывает весь объем данных, но со временем новые данные начинают поступать быстрее, чем успевают обрабатываться.

Вопрос: Как обрабатывать данные быстрее?

Ответ: Обрабатывать их *инкрементально*, т.е. не весь объем, а только изменившуюся часть.

Эта идея описана в статье о системе Percolator³.

³<http://bit.ly/2rSE3YL>



Идея: пользователь реализует *триггеры* (функции-обработчики), которые работают с большой таблицей, и запускаются только при обновлении данных.

Особенности:

- Данные хранятся в KV-базе⁴ и адресуются кортежем <ключ, колонка, время>
- Пользовательские триггеры могут транзакционно читать/писать любые наборы строк
- После записи новых данных в колонку срабатывают триггеры, которые «подписаны» на нее

⁴<http://bit.ly/2raH9ow>



Цель: сделать платформу, реализующую модель вычислений Percolator в модели MapReduce.

Задачи:

1. Разобраться с Percolator и понять, как реализовать его в MapReduce
2. Реализовать и протестировать платформу
3. Внедрить платформу в текущие бизнес-процессы и обеспечить ее поддержку



1. На MR-кластерах хранится большинство данных компании, а значит работать с ними лучше на этих же кластерах
2. Хранение данных в общем MR-кластере дает возможность обрабатывать их другими MR-приложениями
3. YТ – надежная протестированная система, использование которой позволит сэкономить большое количество ресурсов



Основные проблемы, возникающие при реализации модели Percolator на MapReduce:

- Отсутствие построчных транзакций.
- Операция работает только с данными одного ключа.



- Логически данные представляют из себя таблицу
- Триггеры выполняются в транзакциях с уровнем изоляции snapshot⁵ и работают с произвольными ключами базы
- Запуск exactly once на новых данных

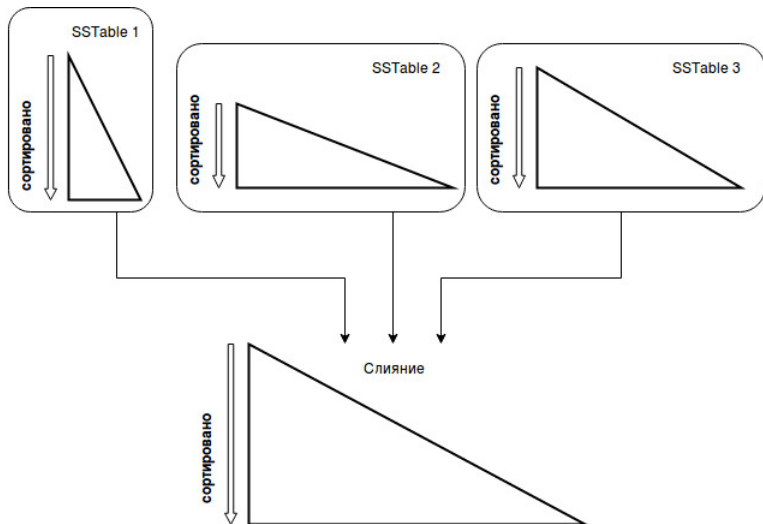
⁵<http://bit.ly/1SSxgID>



- Физически база представляет из себя LSM-дерево⁶, каждая SSTable которого – таблица на YT.
- Операций с базой – Reduce таблиц с данными и некоторых вспомогательных таблиц

Например: для запуска триггеров только на новых данных достаточно сделать Reduce таблиц с информацией о том, какие колонки изменились, и таблиц базы.

⁶<http://bit.ly/2rnGVJV>





Свой ключ – ключ, который передан в Reduce-операцию, запускающую триггер.

Чужие ключи – все остальные ключи в базе.

Предположение: обычно, триггеры читают небольшое количество чужих ключей, поэтому их чтение можно реализовать с помощью *зеркалирования* данных.



- Каждый ключ хранит «кэш» чужих ключей – дополнительная колонка в таблице
- При попытке прочитать ключ, который отсутствует в кэше, генерируется «запрос на чтение» – запись во вспомогательную таблицу
- Перед запуском триггеров происходит разрешение запросов на чтение
- При обновлении данных по ключу *Key*, обновление *зеркалируется* в кэши всех ключей, которые читают *Key*



```
def trigger(url, ctx):
    content = ctx.get_bytes(url, "content")
    rank = 0
    for link in extract_links(content):
        rank += ctx.get_int(link, "rank")
    ctx.set_int(url, "rank", rank)
```



- Типизированный доступ к данным из триггеров
- Импорт/экспорт данных с произвольными форматами
- Автоматический конфигурируемый compaction и сборка мусора
- Мониторинг системы (сколько времени исполнялись триггеры, сколько использовали памяти, ...)



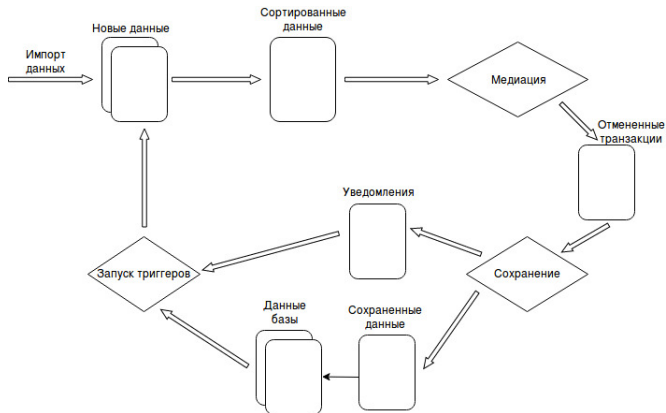
- Выявлены особенности модели Percolator в ограничениях модели MapReduce
- Платформа реализована и протестированная на MapReduce
- Внедрена в несколько production процессов



Спасибо за внимание!

Заключение

Общий pipeline



Заключение

Запуск триггеров

