

# Cassandra

Курс «Базы данных»

Цесько Вадим Александрович

<http://incubos.org>

@incubos

Computer Science Center

30 сентября 2013 г.

# Содержание

- 1 Введение
- 2 Модель данных
- 3 Архитектура
- 4 Детали реализации
- 5 Заключение
- 6 Домашнее задание

# History

- July 2008 — open-sourced by Facebook<sup>1</sup>
- March 2010 — graduated from the Apache Incubator
- Influenced by Amazon Dynamo
- Committers: Rackspace, Digg, Twitter, Amazon, Microsoft

---

<sup>1</sup>Facebook. Cassandra — A Decentralized Structured Storage System:  
<http://www.cs.cornell.edu/projects/ladis2009/papers/lakshman-ladis2009.pdf>

# Features

- Decentralized — no SPoF, every node is identical
- Multi data center replication
- Elastic Scalability
- High Availability and Fault Tolerance
- Tunable consistency
- CQL

# Data Model

- Based on Google Bigtable
- Hybrid: Key-value + Column-oriented
- Каждая строка содержит множество колонок
- Колонка (ячейка): имя + значение
- Разные строки могут иметь разный набор колонок
- Группировка колонок в семейства и суперсемейства

## Table

A distributed multi dimensional map indexed by a key

# ACID

- Никаких JOIN-ов и внешних ключей
- *Атомарность* на уровне строк, нет rollback
- Можно получить ошибку при записи, но запись состоится
- Отметки времени от клиента при конфликтах
- Настраиваемая *консистентность* (количество реплик)
- *Изоляция* на уровне строк
- *Durability*: commit log + replication

# CAP

- Tunable Consistency
- Tunable Availability
- Partition Tolerance

# Пользователи

Самое популярное колоночное хранилище<sup>2</sup>:

- Cisco
- CERN
- Digg
- Facebook
- IBM
- Netflix
- Reddit
- SoundCloud
- Twitter
- Odnoklassniki
- Yandex

<sup>2</sup><http://db-engines.com/en/ranking/wide+column+store>



# Мотивация

*If I had asked people what they wanted, they would have said faster horses.*

---

Henry Ford

## Что не так с RDBMS?

Очень рекомендую Chapter 1. "What's Wrong with Relational Databases?" in "Cassandra: The Definitive Guide"<sup>a</sup>

---

<sup>a</sup>Eben Hewitt. Cassandra: The Definitive Guide. 2010. ISBN 1449390412.

# У нас есть сервис на RDBMS

Всё работает?

Не трогай!

С чем могут быть проблемы:

- (Distributed) Transactions<sup>3</sup>
- (Distributed) JOINS
- (Distributed) Schema evolution
- Sharding<sup>4</sup>
  - Functional segmentation
  - Key-based sharding
  - Lookup table

<sup>3</sup>Gregor Hohpe. Starbucks Does Not Use Two-Phase Commit:  
[http://www.eaipatterns.com/ramblings/18\\_starbucks.html](http://www.eaipatterns.com/ramblings/18_starbucks.html)

<sup>4</sup>Michael Stonebraker. The Case for Shared Nothing. 1986

# Use Case: RDBMS Scaling

## LiveJournal

Brad Fitzpatrick. LiveJournal's Backend: A history of scaling<sup>a</sup>. August 2005.

---

<sup>a</sup>[http://www.danga.com/words/2005\\_oscon/oscon-2005.pdf](http://www.danga.com/words/2005_oscon/oscon-2005.pdf)

# Web Scale

IDC. The Expanding Digital Universe<sup>5</sup>. March 2007:

- YouTube: 2007 — 100 M videos/day, **2008** — 1 B videos/day
- Chevron: 2 TB/day
- Internet: 2006 — 166 EB, 2010 — 1 ZB
- Wall-Mart: 2000 — 110 TB, 2004 — 0.5 PB
- Email: 1998 — 253 M accounts, 2010 — 2 B accounts
- etc.<sup>6</sup>

---

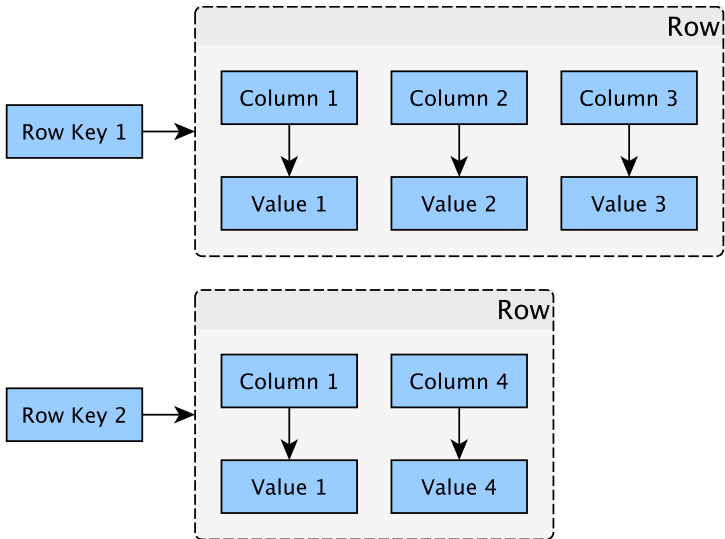
<sup>5</sup><http://www.emc.com/collateral/analyst-reports/expanding-digital-idc-white-paper.pdf>

<sup>6</sup><http://highscalability.com/>

# Модель данных

- Cassandra — партиционированное хранилище рядов
- Исходим из запросов, а не из сущностей и связей
- Таблица — коллекция упорядоченных колонок
- Keyspace — группа таблиц (обычно per application)

# Если глубже



# Интерфейсы

- Thrift API
  - Исторически первый
  - Низкоуровневый
  - Довольно многословный
- **CQL** (текущая версия 3.1)
  - SQL-подобный
  - (Почти насколько же) выразительный
  - Развивается

# Пример

```
1 CREATE TABLE songs (  
2   id uuid PRIMARY KEY,  
3   title text,  
4   album text,  
5   artist text,  
6   data blob  
7 );  
  
1 CREATE TABLE playlists (  
2   id uuid,  
3   song_order int,  
4   song_id uuid,  
5   title text,  
6   album text,  
7   artist text,  
8   PRIMARY KEY (id, song_order));
```



# Особенности

- Compound keys (clustering)
- Indexes
- Collection types<sup>7</sup>: set, list, map
- INSERT = UPDATE = UPSERT
- DELETE для строки или набора ячеек
- TTL
- Counters
- См. спеку<sup>8</sup>

---

<sup>7</sup>[http://www.datastax.com/dev/blog/cql3\\_collections](http://www.datastax.com/dev/blog/cql3_collections)

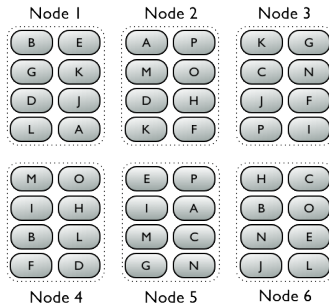
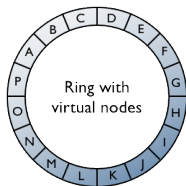
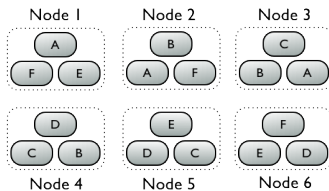
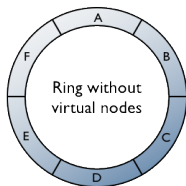
<sup>8</sup>[http:](http://www.datastax.com/documentation/cql/3.1/webhelp/index.html)

[//www.datastax.com/documentation/cql/3.1/webhelp/index.html](http://www.datastax.com/documentation/cql/3.1/webhelp/index.html)

# Компоненты

- **Consistent Hashing + VNodes** (см. лекцию 1)
- **Gossip** — состав и состояние узлов
- **Partitioner** — данные по узлам
- **Replica placement strategy** — реплики по узлам
- **Snitch** — топология (ДЦ и стойки)

# Virtual Nodes



# Gossip

- Узлы периодически обмениваются информацией о себе и других узлах
- Каждую секунду не больше чем с тремя узлами
- Все узлы быстро узнают информацию обо всех узлах
- Сообщения имеют версию — устаревшая информация затирается
- Seed nodes во избежания партиционирования
- Gossip-информация хранится на каждой ноде персистентно
- Обнаружение сбоев узлов (+ dynamic snitch)<sup>9</sup>

<sup>9</sup>Naohiro Hayashibara, Xavier Défago, Rami Yared and Takuya Katayama. The  $\phi$  Accrual Failure Detector. 2008

# Partitioner

- Consistent Hashing<sup>10</sup>
- Новый Murmur3Partitioner ( $-2^{63}$  to  $+2^{63}$ ) и старый RandomPartitioner ( $0$  to  $2^{127} - 1$ )
- ByteOrderedPartitioner
  - Лексикографически по байтам ключа
  - Сложная балансировка нагрузки (расчёт диапазонов партиций вручную)
  - Неравномерная балансировка для разных таблиц
  - Последовательная запись упирается в один узел

---

<sup>10</sup>[http://www.datastax.com/documentation/cassandra/2.0/webhelp/cassandra/architecture/architectureDataDistributeAbout\\_c.html](http://www.datastax.com/documentation/cassandra/2.0/webhelp/cassandra/architecture/architectureDataDistributeAbout_c.html)

# Replica Placement Strategy

- Replication Factor — количество реплик на кластер
- Каждая реплика на отдельном узле
- Все реплики равноправны
- SimpleStrategy:
  - 1 ДЦ (если планируете расширяться, не используйте)
  - Первая реплика на узел от Partitioner
  - Остальные — по часовой стрелке по кольцу
- NetworkTopologyStrategy:
  - Определяет количество реплик в каждом ДЦ
  - Узлы для реплик — идём по кольцу до следующей стойки
  - Часто стойки вырубаются целиком (питание, охлаждение, сеть, ...)

# NetworkTopologyStrategy

## Факторы

- Желательно чтение внутри ДЦ
- Учитываем типичные сбои

## По 2 реплики в каждом ДЦ

Сбой одного узла — `ConsistencyLevel.ONE`

## По 3 реплики в каждом ДЦ

- Сбой одного узла —  
`ConsistencyLevel.LOCAL_QUORUM`
- Сбой двух узлов — `ConsistencyLevel.ONE`

# Snitch

- **Одинаковая конфигурация** на узлах
- Dynamic snitching wrapper<sup>11</sup>
  - По умолчанию и только на чтение
  - Данные с самой быстрой реплики, с остальных — контрольные суммы
  - Статистика Read Latency и текущие задачи узлов
- SimpleSnitch: никаких ДЦ и стоек
- RackInferringSnitch: (110.dc.rack.node)
- PropertyFileSnitch: адреса в ДЦ/стойки в файле
- GossipingPropertyFileSnitch: локальный ДЦ и стойка в файле на каждом узле + Gossip

<sup>11</sup><http://www.datastax.com/dev/blog/>



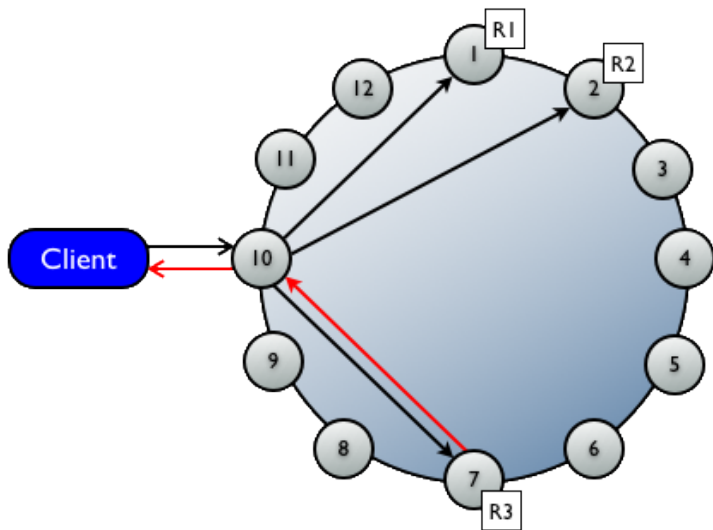
# Клиентские запросы

- Можно обращаться к любому узлу
- Узел становится координатором для текущего запроса
- Координатор проксирует с учётом Partitioner и Replica Placement Strategy

# Write

- Запрос всем репликам независимо от `ConsistencyLevel`
- `ConsistencyLevel` определяет, сколько реплик должно ответить для успеха

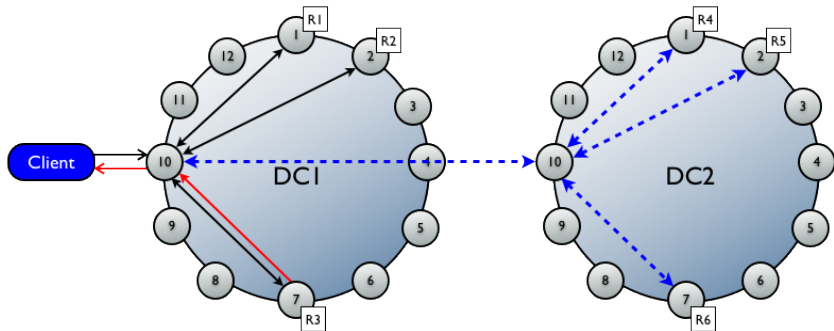
# Write: Пример



# MultiDC Write

- Оптимизация — один координатор в каждом удалённом ДЦ
- `ConsistencyLevel.ONE` или `ConsistencyLevel.LOCAL_QUORUM` — обязаны ответить только локальные узлы (география не влияет на задержку)

# MultiDC Write: Пример



# Write: ConsistencyLevel

- ANY — всегда успех (hinted handoff)
- ONE — в commit-log одного узла
- TWO
- THREE
- QUORUM — inter DC + нужен запас прочности
- LOCAL\_QUORUM — быстрее, чем QUORUM
- EACH\_QUORUM — выше консистентность
- ALL — WAT?

## Не забываем

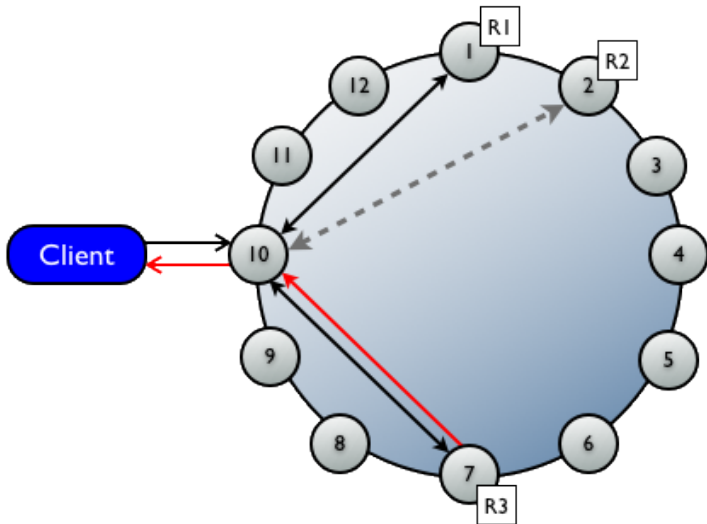
Даже при ONE и LOCAL\_QUORUM запись пошлют всем репликам в т. ч. в других ДЦ

# Read

Read-запросы от координатора репликам:

- Прямой запрос на чтение (в соответствии с `ConsistencyLevel`)
  - Самым «быстрым» репликам
  - Сравниваем ответы
  - Если не совпадают, то самая свежая (по `timestamp`) — клиенту
- Фоновый `read repair` (всем остальным репликам)
  - Для синхронизации «горячих» данных
  - `read_repair_chance` по умолчанию 0.1
  - Если не совпадают хэши, то перезаписываем

# Read: Пример





# Read ConsistencyLevel

- ONE — ближайшая реплика, возможно, устаревшие данные
- TWO
- THREE
- QUORUM — inter DC + нужен запас прочности
- LOCAL\_QUORUM — быстрее, чем QUORUM
- EACH\_QUORUM — выше консистентность
- ALL — WAT?

## Интерпретация

Возвращаем самое свежее значение из требуемого множества реплик

# Quorum

## Quorum

$$Quorum = \lfloor \frac{RF}{2} + 1 \rfloor$$

- $RF = 2 \Rightarrow Quorum = 2$  — нет запаса
- $RF = 3 \Rightarrow Quorum = 2$  — в запасе 1 узел
- $RF = 4 \Rightarrow Quorum = 3$  — в запасе 1 узел
- $RF = 5 \Rightarrow Quorum = 3$  — в запасе 2 узла

## Consistency

$$R + W > RF$$

- $R = 2 + W = 2 > RF = 3$  — в запасе 1 узел

# Клиентские запросы

## Проблема

- Узел сбойнул: железо или (чаще) сеть
- А мы хотим на него записать
- Что делать?

## Hinted Handoff

- Пусть известно (Gossip), что узел лежит, или узел не отвечает
- Координатор запоминает hint локально
- Когда обнаружится, что узел поднялся (Gossip), ему перешлют накопленные hints

# Hint

## Содержимое

- Кому предназначается
- Значение ключа
- Данные

## Внимание

- Hints хранятся ограниченное время (по умолчанию 3 часа)
- Официально рекомендуется периодический запуск `repair`

# Ограничения

- Hint засчитывается не на всех уровнях консистентности<sup>12</sup>: 2 узла (1 мёртв) + RF=1 + ConsistencyLevel.ONE
- Но **всегда работает**, если не выключен
- ConsistencyLevel.ANY — extreme write availability
- Если умрёт машина с hints, то мы их потеряем

---

<sup>12</sup><http://www.datastax.com/dev/blog/modern-hinted-handoff>

# Anti-entropy

## Проблема

- Узел умер — пропустил удаление данных
- Вернулся к жизни — данные возродились

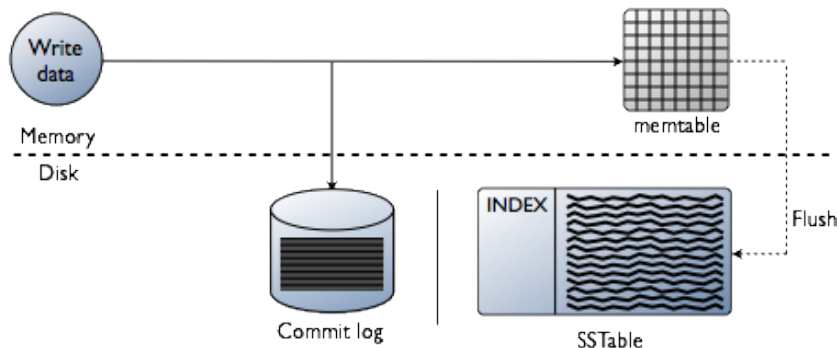
## Anti-entropy

- Иницилируем с помощью `nodetool repair`
- Запускаем `readonly major compaction`
- Строим Merkle Tree<sup>a</sup>
- Обмениваемся деревьями и ищем отличия
- Обмениваемся отличающимися сегментами

---

<sup>a</sup>[http://en.wikipedia.org/wiki/Merkle\\_tree](http://en.wikipedia.org/wiki/Merkle_tree)

# Write Path



# Пояснения

- Memtables and SSTables per table
- Memtable — отсортирована
- **SSTable** — неизменяемая отсортированная ассоциативная таблица
- Обычно строка распределена по нескольким SSTable



# Пример

Команды:

```
write (k1, c1:v1)
write (k2, c1:v1 c2:v2)
write (k1, c1:v4 c3:v3 c2:v2)
```

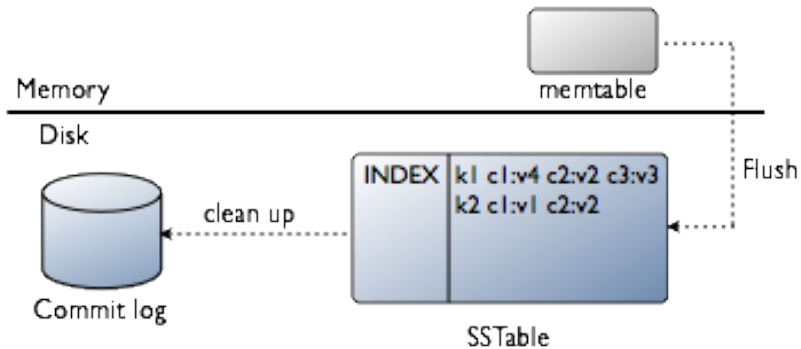
Commit-log:

```
k1, c1:v1
k2, c1:v1 c2:v2
k1, c1:v4 c3:v3 c2:v2
```

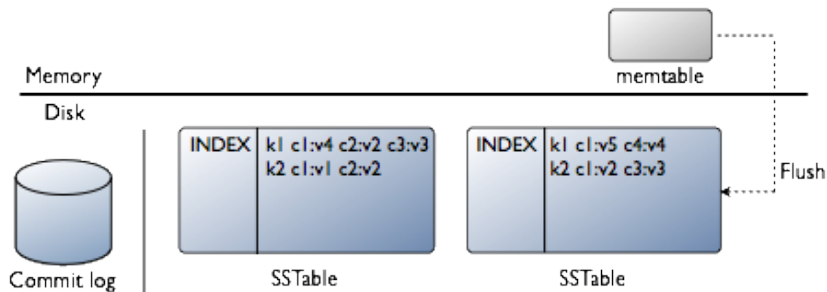
Memtable/SSTable:

```
k1 c1:v4 c2:v2 c3:v3
k2 c1:v1 c2:v2
```

# Flush



# Update Path



# Delete

- SSTable неизменяема
- Delete — tombstone marker
- Реальное удаление по истечении `gc_grace_seconds` во время compaction
- Удалённые данные могут возродиться (см. Anti-entropy)

# Compaction

- Объединяет SSTable-файлы
  - Фрагменты строк
  - Протухшие tombstones
  - Перестраивает индексы
- SSTable отсортирована  $\Rightarrow$  последовательный проход
- `SizeTieredCompactionStrategy` и `LeveledCompactionStrategy`<sup>1314</sup>

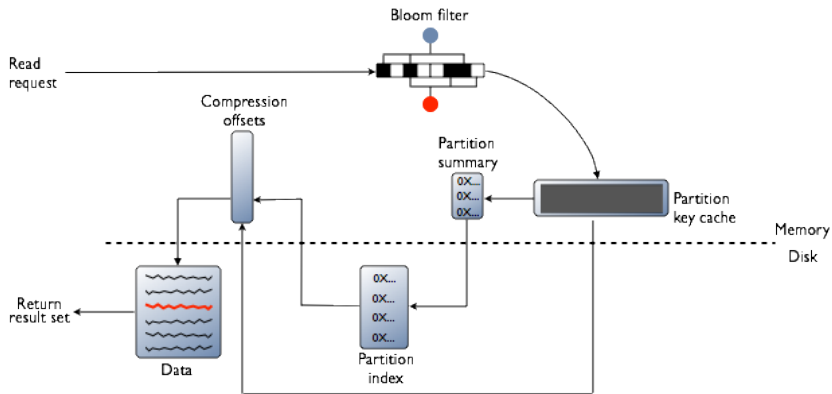
---

<sup>13</sup>[http:](http://www.datastax.com/dev/blog/when-to-use-leveled-compaction)

[//www.datastax.com/dev/blog/when-to-use-leveled-compaction](http://www.datastax.com/dev/blog/when-to-use-leveled-compaction)

<sup>14</sup><http://www.datastax.com/dev/blog/leveled-compaction-in-apache-cassandra>

# Read



# Комментарии

- Каждая SSTable имеет Bloom filter<sup>15</sup> — вероятность нахождения ключа в файле
- Если вероятность отлична от 0, идём в partition key cache
- Если нашли ключ в кэше, идём по смещению, находим сжатый блок и достаём данные
- Если не нашли ключ в кэше:
  - В partition summary примерно находим смещение на диске
  - Читаем последовательно блок с диска
  - Из compression offset map вынимаем индекс блока
  - Читаем сжатые данные и возвращаем клиенту

---

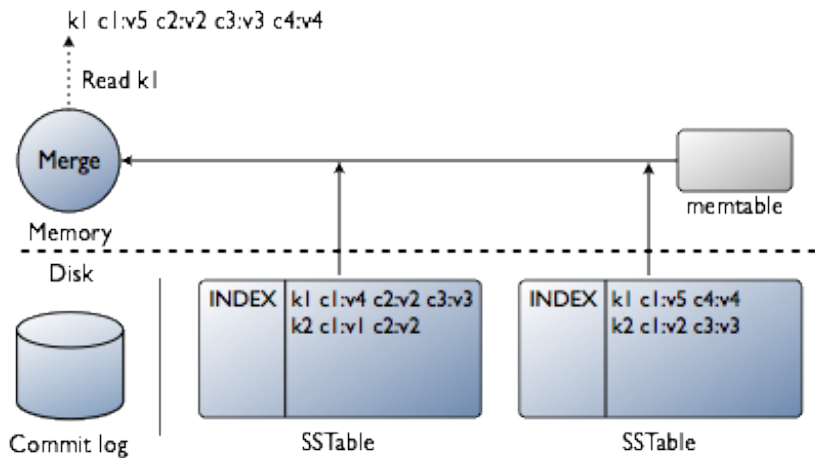
<sup>15</sup>[http://en.wikipedia.org/wiki/Bloom\\_filter](http://en.wikipedia.org/wiki/Bloom_filter)

# Дополнительные структуры

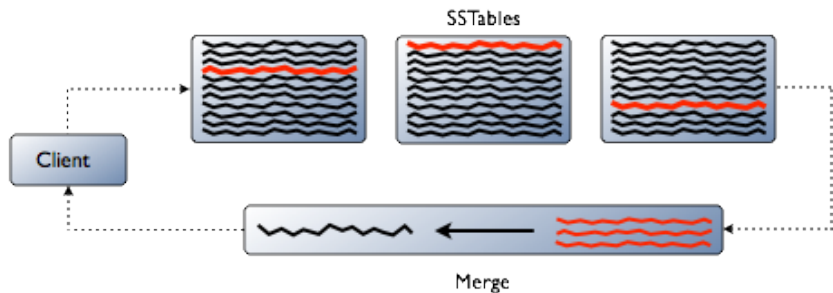
- Bloom filter — 1-2 ГБ / млрд. ключей
- Partition summary — по умолчанию шаг 128
- Compression offset map — 1-3 ГБ / 1 ТБ сжатых данных



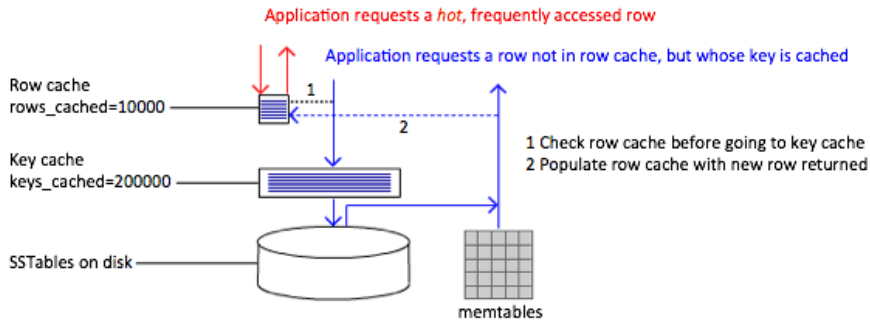
# Read Path



# Write-through Row Cache



# Key Cache + Row Cache



См. подробности<sup>16</sup>.

<sup>16</sup><http://www.datastax.com/docs/1.1/operations/tuning>

# CQL

```
1 CREATE TABLE comments (  
2     article_id uuid,  
3     posted_at timestamp,  
4     author text,  
5     karma int,  
6     content text,  
7     PRIMARY KEY (article_id, posted_at))
```

# Колонки

620e8...	1340000000:	
	1340000000:author	Alice
	1340000000:content	Nice article, thanks
	1340000000:karma	4
	1340130000:	
	1340130000:author	Bob
	1340130000:content	I agree with Alice, +1
	1340130000:karma	0

См. подробности<sup>17</sup>.

<sup>17</sup><http://www.datastax.com/dev/blog/thrift-to-cal3>

# Wide vs Skinny Rows

- Wide Rows
  - (+) Колонки отсортированы — возможны трюки<sup>18</sup>
  - (+/-) Атомарность
  - (-) Row cache (скорее всего) не работает
- Skinny Rows
  - (+) Row cache
  - (+/-) Неатомарно, но меньше contention
  - (-) Больше записей — больше индексы, хуже работает Bloom Filter

---

<sup>18</sup>Классная Cassandra: <http://javapoint.ru/talks/15/>

# Опыт использования

- Из Hadoop можно «положить» кластер — ждём `CqlOutputFormat`
- Довольно много «ручек» для подкручивания
- Внимание к GC
- Развитые средства мониторинга
- Периодический `nodetool repair`

# Осталось за кадром

- SEDA<sup>19</sup>
- Эксплуатация и тюнинг (JMX, nodetool)
- Масштабные примеры использования<sup>20</sup>
- См. блог разработчиков<sup>21</sup>

---

<sup>19</sup>[http:](http://en.wikipedia.org/wiki/Staged_event-driven_architecture)

[//en.wikipedia.org/wiki/Staged\\_event-driven\\_architecture](http://en.wikipedia.org/wiki/Staged_event-driven_architecture)

<sup>20</sup><http://www.planetcassandra.org/blog/post/slideshare-presentations---cassandra-summit-2013>

<sup>21</sup><http://www.datastax.com/dev/blog>



# Следующая лекция

- MongoDB (приглашённый доклад)

# Требования (ещё раз)

- Проект на <http://bitbucket.org> (Hg) / <http://github.com> (Git) + Issue Tracking
  - Открытый проект — ссылку мне на почту
  - Закрытый проект — мне админские права (incubos)
- **Обязательно:** Scala/Java + ScalaTest/JUnit + Maven/SBT
- **Обязательно:** README + INSTALL
- Команда 1-3 человека

# Feature Request 0

- «Телефонная книга» (см. лекцию 01-intro<sup>22</sup> слайд #26)
- **Обязательный**
- Выдан 2013-09-16
- Сгорает **сегодня**

---

<sup>22</sup><http://incubos.org/2013/edu/db/01-intro.pdf>

# Зачёт ХОТЯТ

- [https://github.com/Elizaveta239/Storage\\_system\\_csc](https://github.com/Elizaveta239/Storage_system_csc)
- <https://github.com/Feodorov/DB>
- <https://github.com/medvector/cscenter-databases13>
- <https://github.com/Noxoomo/DBCSC>

WAT

?!

# Feature Request 1

- «Big data» per node (см. лекцию 04-commit<sup>23</sup> слайд #64)
- Выдан 2013-09-23
- Сгорает **2013-10-14**

---

<sup>23</sup><http://incubos.org/2013/edu/db/04-commit.pdf>

# Кто всё-таки хочет участвовать

- FR1 заменяет FR0
- Становится **обязательным**
- Не засчитывается как дополнительный (не приносит балла)
- Всё так же сгорает **2013-10-14**

Не забывайте

Можно участвовать командой до 3 человек

# Система оценивания

- Будет выдано не меньше 6 Feature Requests (FR)
- Для каждого FR — жёсткий срок (но не меньше 2 недель)
- За каждый FR — 0, 1 или 2 балла (за особое мастерство и креативность)
- Для зачёта нужно FR0/FR1 + 4 балла + отчёт
- Проверяю каждые выходные (по возможности чаще)
- Ценю равномерный прогресс
- Ведите стабильную и нестабильную ветки
- Обратная связь в виде issues
- Не все issues закрыты — не зачтено

# FR2: Static Sharding

- Равноправные узлы
- Данные распределены по кластеру
- Статическая конфигурация (узлы и диапазоны ключей)
- Клиент знает конфигурацию и куда ходить
- **Тесты** + Testability
- Сгорает — **2013-10-14**



# FR3: Master-Slave Replication

- Commit-log от master к slaves
- Статическая конфигурация
- Клиент знает конфигурацию — пишет в master, читает из slaves
- **Тесты** + Testability
- Сгорает — **2013-10-21**

# FR4: Scalability

- Динамическая конфигурация (добавляем и удаляем ноды на лету)
- В простейшем случае 2PC новой конфигурации
- Graceful Degradation
- Consistent Hashing
- Миграция партиций при добавлении/удалении узлов
- **Тесты** + Testability
- Сгорает — **2013-10-28**

# Вопросы?

- <http://incubos.org/contacts/>
- Общие вопросы — в Twitter: @incubos
- Вопросы по лекциям — в комментариях:  
<http://incubos.org/blog/>
- Частные вопросы — в почту  
[vadim.tsesko@gmail.com](mailto:vadim.tsesko@gmail.com)