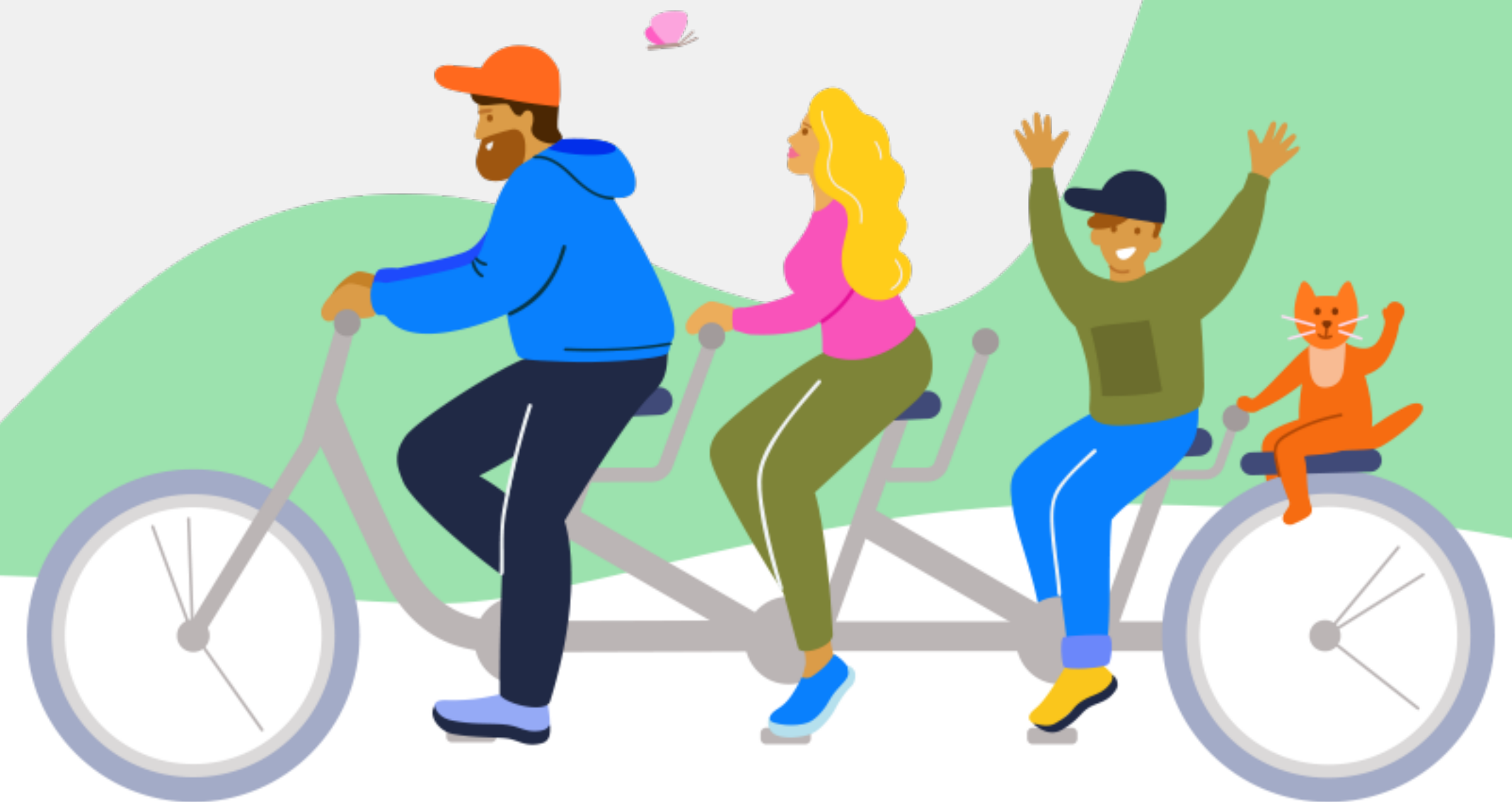


Эффективные

надежные

микросервисы



Одноклассники

6

ТЫС
машин

4

облака

15

ТЫС
задач

2000

сервисов

Общение - это ОК

5%

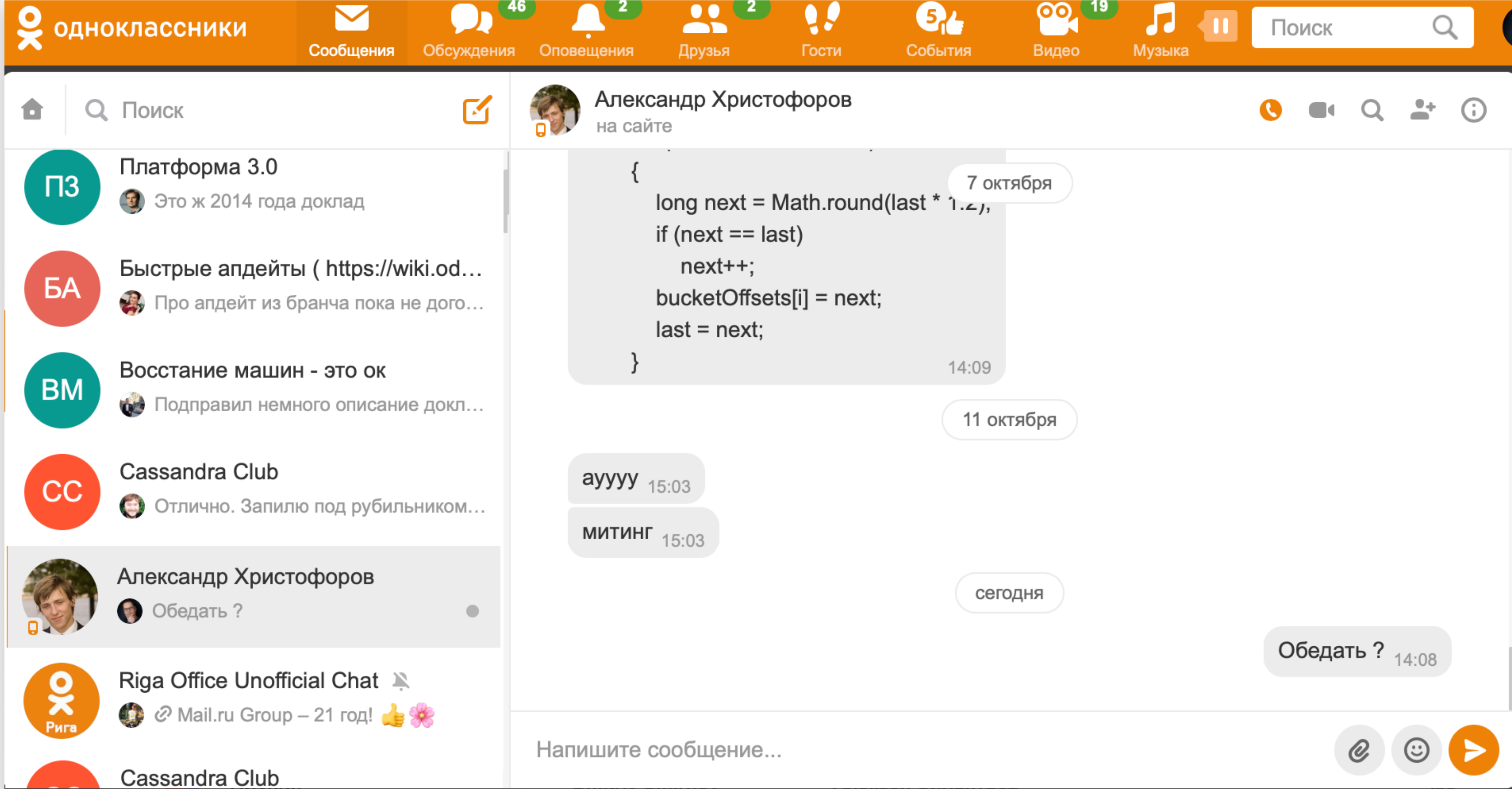
активных чатов

95%

запросов

80%

последние 13 сообщений



Сообщения в чатах

600 млрд
сообщений

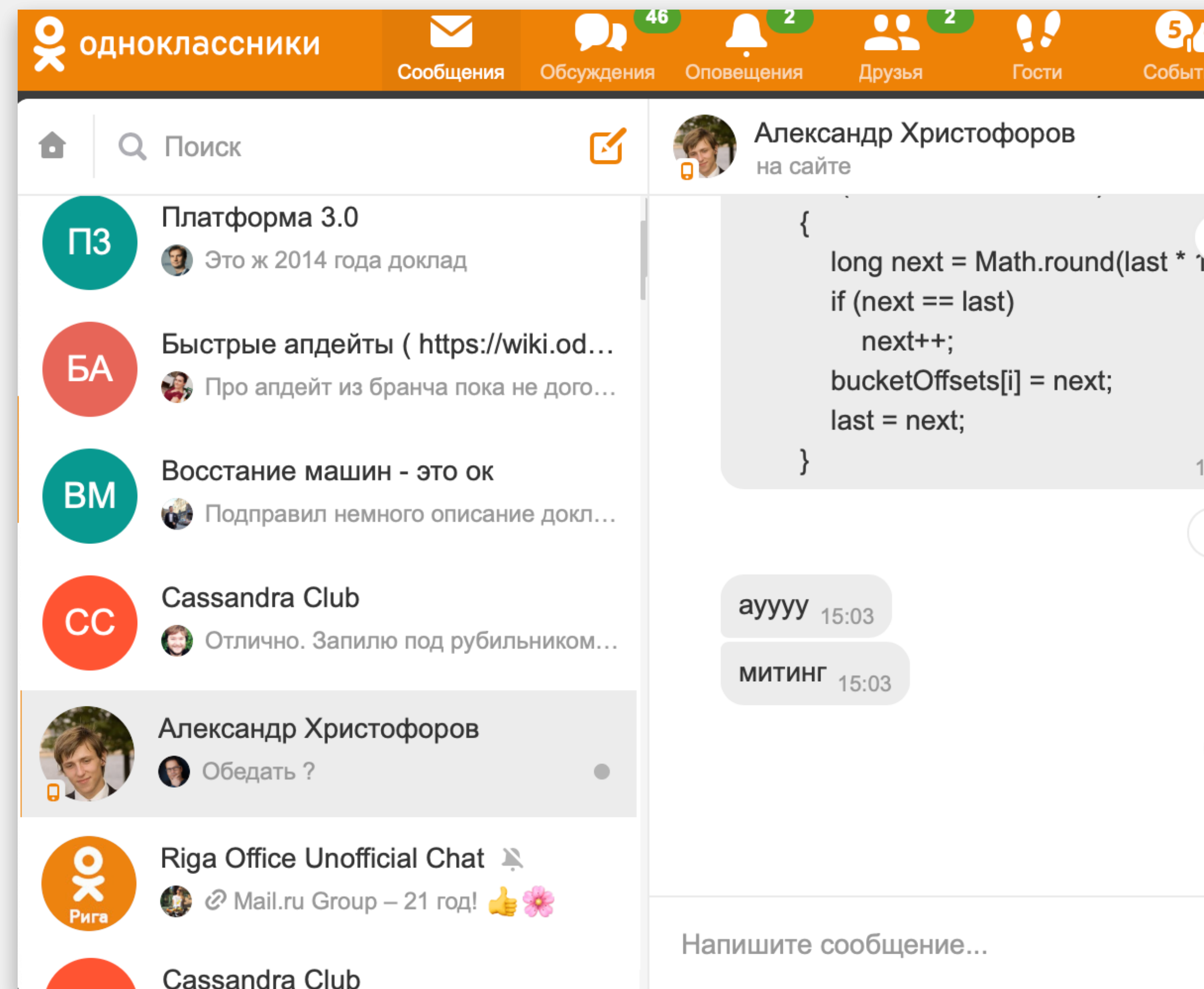
5 млрд
чатов

100 ТБ

ЧИСТЫХ ДАННЫХ

120 тыс/сек
чтений

8 тыс/сек
записей



Сообщения в чатах: поиск

600 млрд
сообщений

5 млрд
чатов

100 ТБ
ЧИСТЫХ ДАННЫХ

120 тыс/сек
чтений

8 тыс/сек
записей

The screenshot shows the VK chat search interface. At the top, there is a navigation bar with icons for messages, discussions, notifications, friends, guests, and events. The search bar contains the text 'joker'. Below the search bar, there is a list of search results under the heading 'ТЕКСТ В СООБЩЕНИЯХ'. The results include:

- ХабраЧат**: А кто два поста про игру на Joker себе п...
- Александр Анисимов**: Привет! А ты придешь обсуждать игру н...
- ВМ**: Восстание машин - это ок
сопру с хабра: Одноклассники — крупне...
- ODKL:all IT**: сюда публикуется инт...
Наша игрушка с JPoint2019 : https://h...
- CU**: Cassandra unconf/meetup
Вот такой текст в чат. Спасибо, за интер...
- Cassandra Club**

On the right side, there is a chat window with the contact 'Александр Христофоров на сайте'. The chat history shows a code block:

```
{  
  long next = Math.round(last * 1  
  if (next == last)  
    next++;  
  bucketOffsets[i] = next;  
  last = next;  
}
```

Below the code, there are two messages: 'ауууу 15:03' and 'МИТИНГ 15:03'. At the bottom of the chat window, there is a text input field with the placeholder 'Напишите сообщение...'.

Сервис сообщений в чатах

операции:

- `getMessages(viewer, chat, from, to)`
- `getLastMessages(viewer, chats)`
- `add(chat, message)`
- `search(viewer, text)`
- `indexMessages()`

	ИД	Автор	Тип	Текст	Аттачи[]
	1	Олег	ТХТ	Леха, привет! Звоню ?	
	1	Олег	ВЗ	callto: Леха, miss	
	1	Проктор	СПАМ	Чистите зубы нашими памперсами!	зубы.гиф
	1	Леха		Ты кто ?	

```
CREATE TABLE Messages (  
  chatId, msgId  
  
  user, type, text, attachments[], terminal, deletedBy[], replyTo,...  
  
  PRIMARY KEY ( chatId, msgId )  
)
```

Современные микросервисы



Прикладная логика



Состояние (данные)

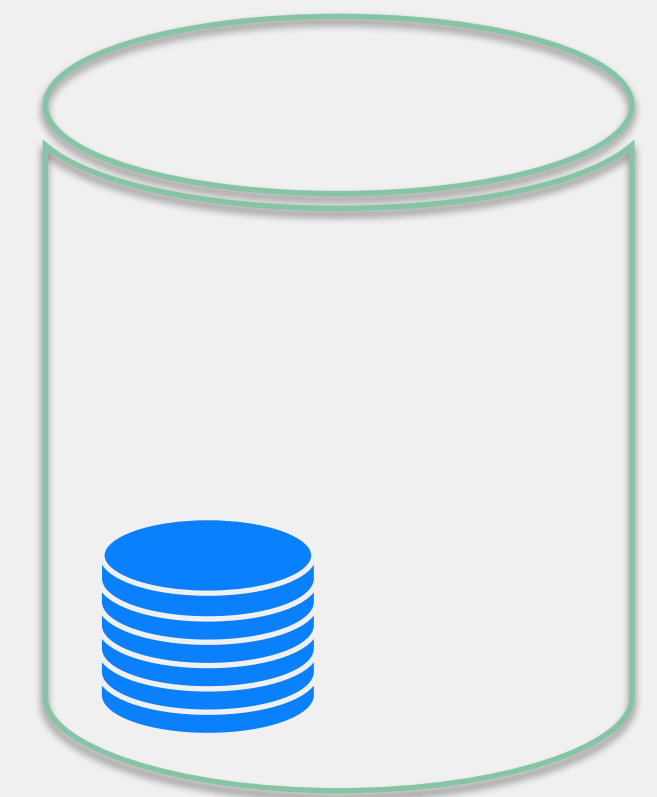
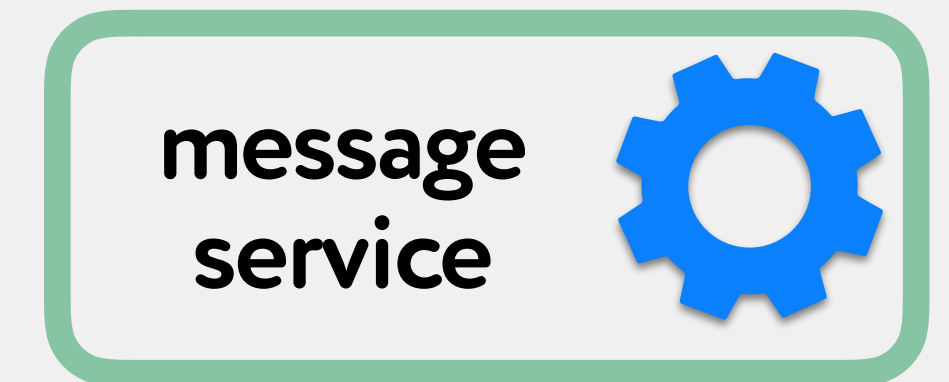


DB

Современные микросервисы

- `getMessages(viewer, chat, from, to)`

```
SELECT FROM Messages  
WHERE chatId = ? AND  
msgId BETWEEN :from AND :to
```



DB

Современные микросервисы

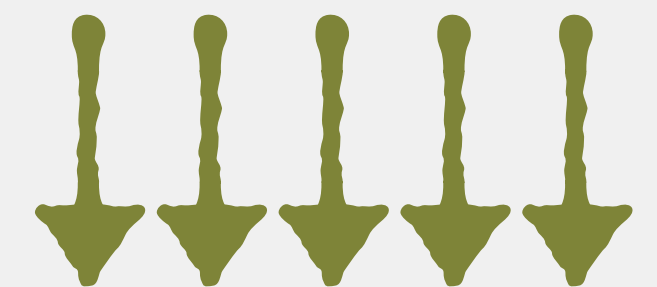
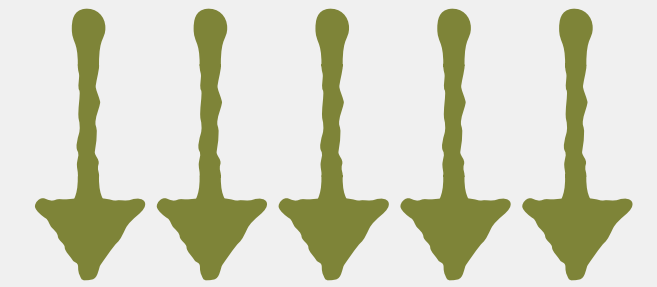
- `getMessages(viewer, chat, from, to)`
- `getLastMessages(viewer, chats)`
- `indexMessages()`

100+ k/sec

5% 95%

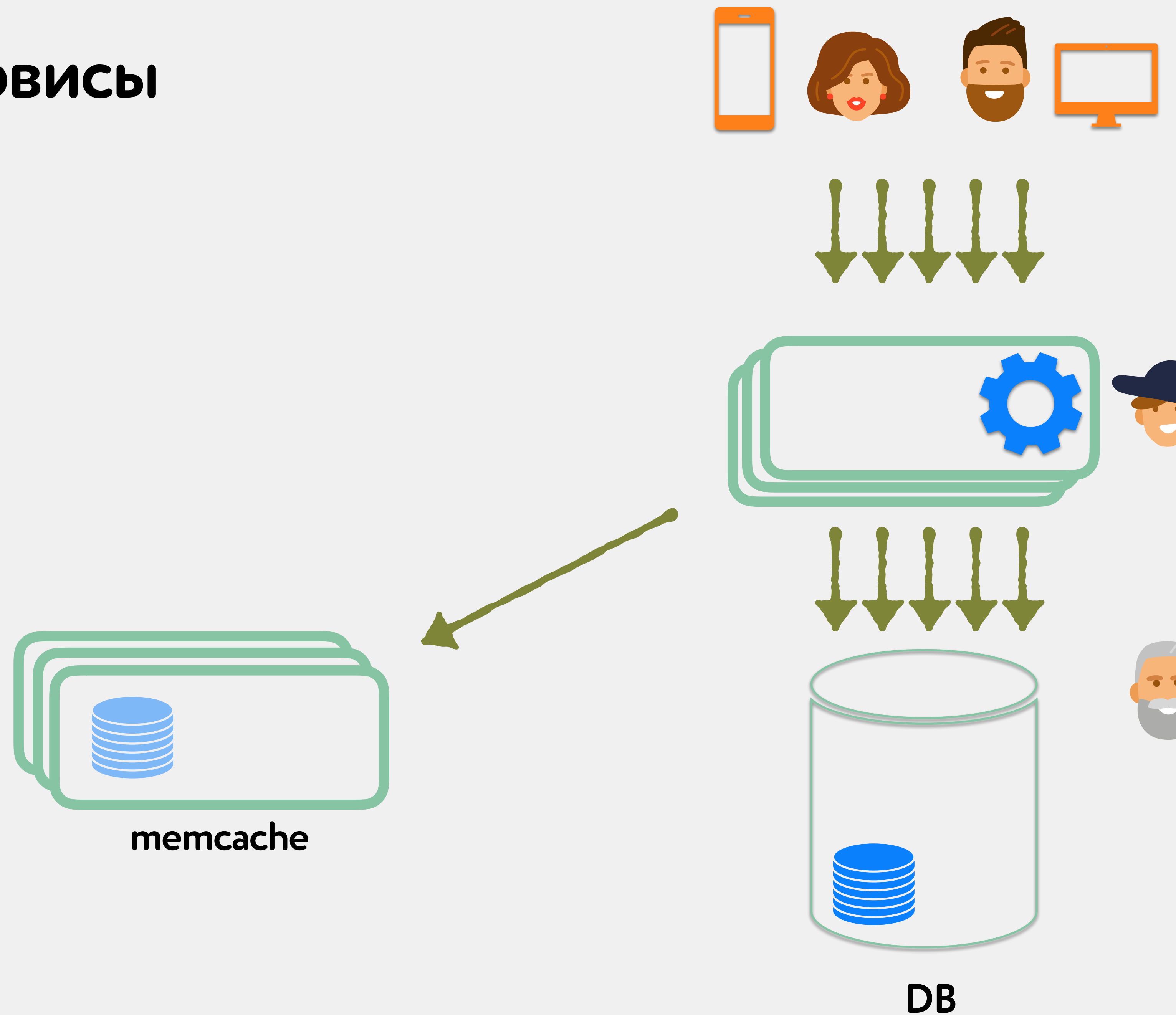
активных чатов

запросов



DB

Современные микросервисы



Микросервисы: потери

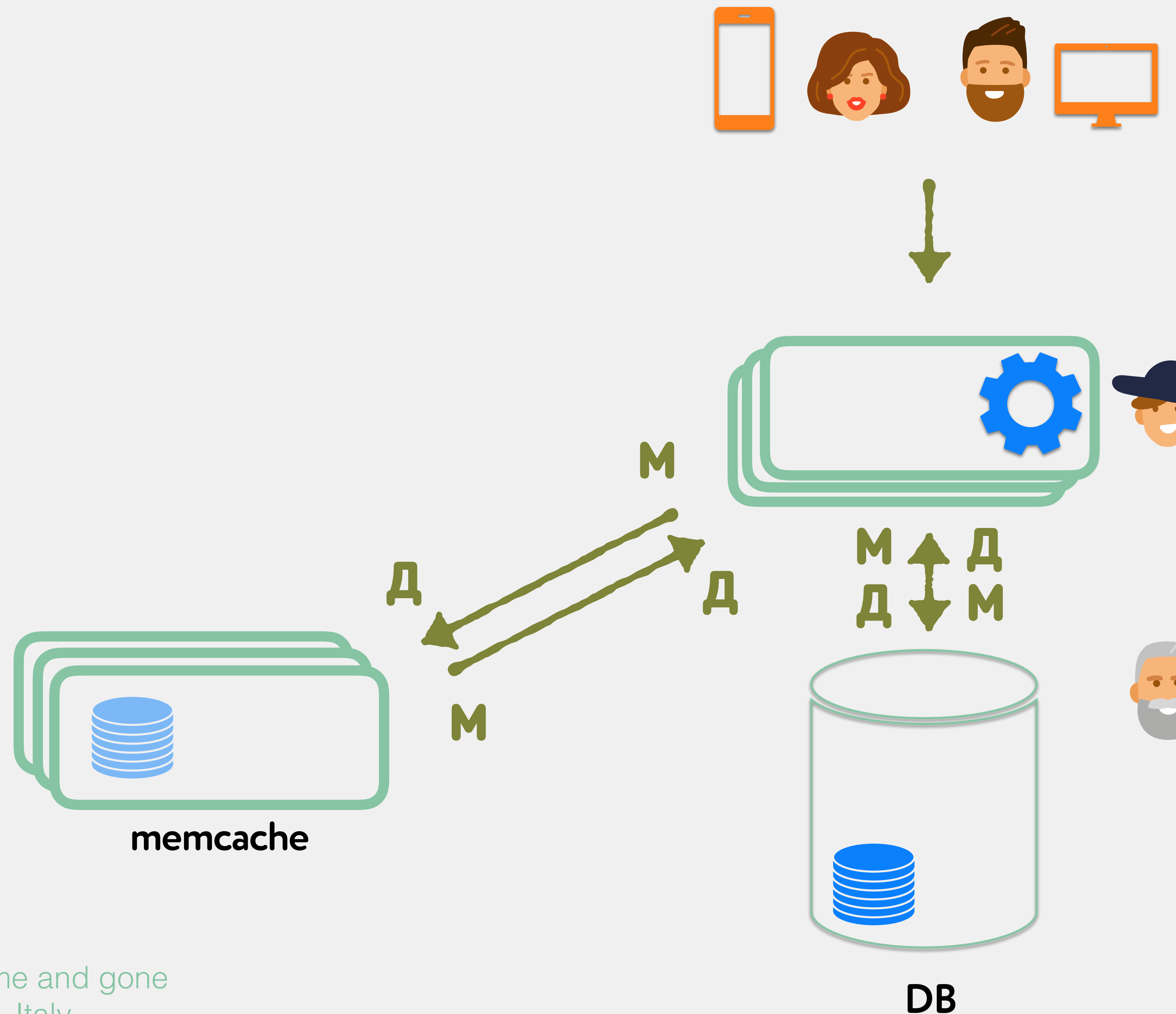
- CPU: (Де) Маршalling

+85%

к нагрузке на ЦПУ⁽¹⁾

+27%

к медианной задержке⁽¹⁾



(1) Fast key-value stores: An idea whose time has come and gone
Adya et al. HotOS '19, May 13–15, 2019, Bertinoro, Italy

Микросервисы: потери

- CPU: (Де) Маршalling
- Чрезмерные чтения/запись

если запись содержит только

10%

полезной информации, то

+46% **+86%**

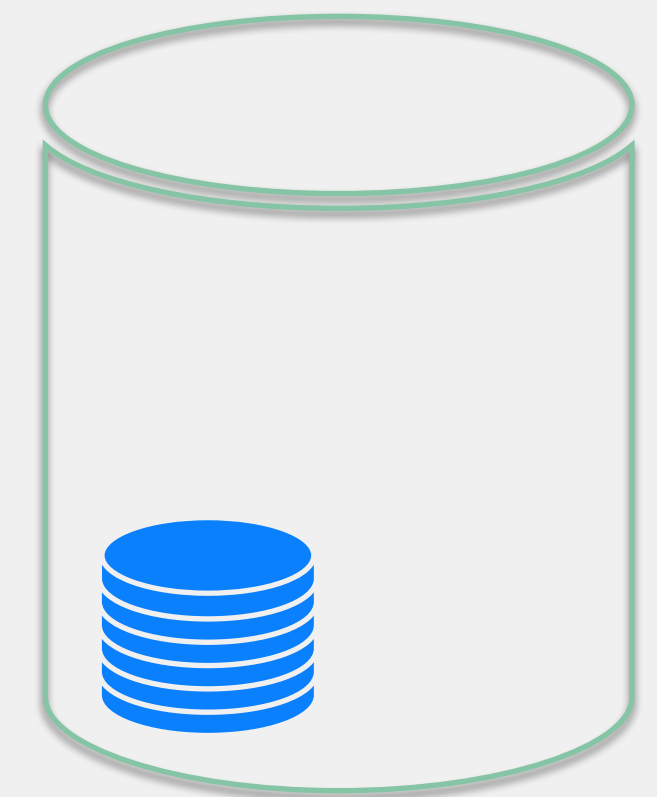
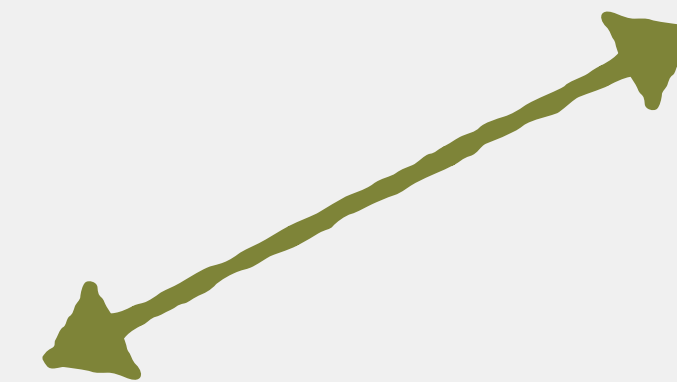
CPU

Net

потрачено зря ¹⁾



memcache



DBMS

(1) Fast key-value stores: An idea whose time has come and gone
Adya et al. HotOS '19, May 13–15, 2019, Bertinoro, Italy

Микросервисы: потери

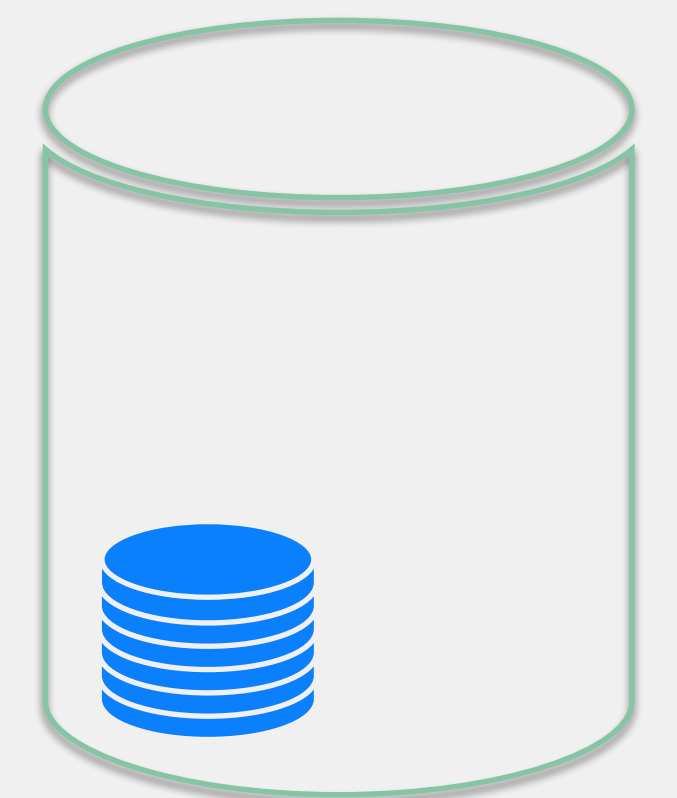
- CPU: (Де) Маршалинг
- Чрезмерные чтения/запись
- Сетевые задержки, трафик

xN

число чтений/записей
на 1 запрос



memcache



DBMS

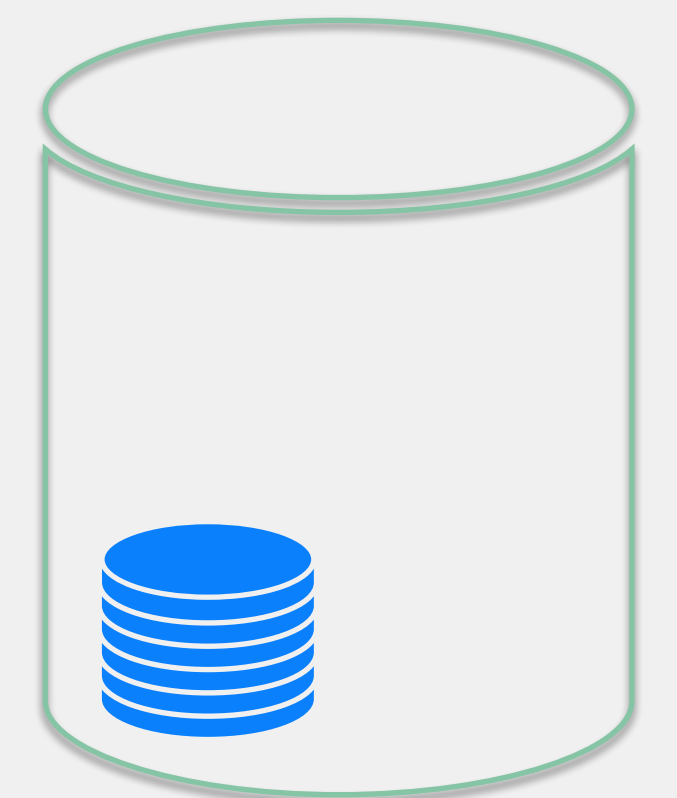
(1) Fast key-value stores: An idea whose time has come and gone
Adya et al. HotOS '19, May 13–15, 2019, Bertinoro, Italy

Микросервисы: что же делать ?

- CPU: (Де) Маршalling (5) [KV-Direct](#)
- Чрезмерные чтения/запись (2) [redis](#), (3) [tarantool](#)
- Сетевые задержки, трафик (4) [NetCache](#)



memcache



DBMS

(2) <http://redis.io>

(3) <https://tarantool.io>

(4) Netcache: Balancing key-value stores with fast in-network caching.
In X. Jin et al, Stoica. SOSP, 2017.

(5) Kv-direct: high-performance in-memory key-value store with programmable nic.
B. Li et al, In SOSP, 2017.



Микросервис с состоянием

Микросервис с состоянием

- ✓ CPU: (De)Маршalling
- ✓ Чрезмерные чтения/запись
- ✓ Сетевые задержки, трафик



Прикладная логика



Встроенный прикладной кеш
custom in-memory store



Встроенная распределенная БД
embedded distributed store

Микросервис с состоянием

всегда быстрее !

а что с отказоустойчивостью ?



Прикладная логика

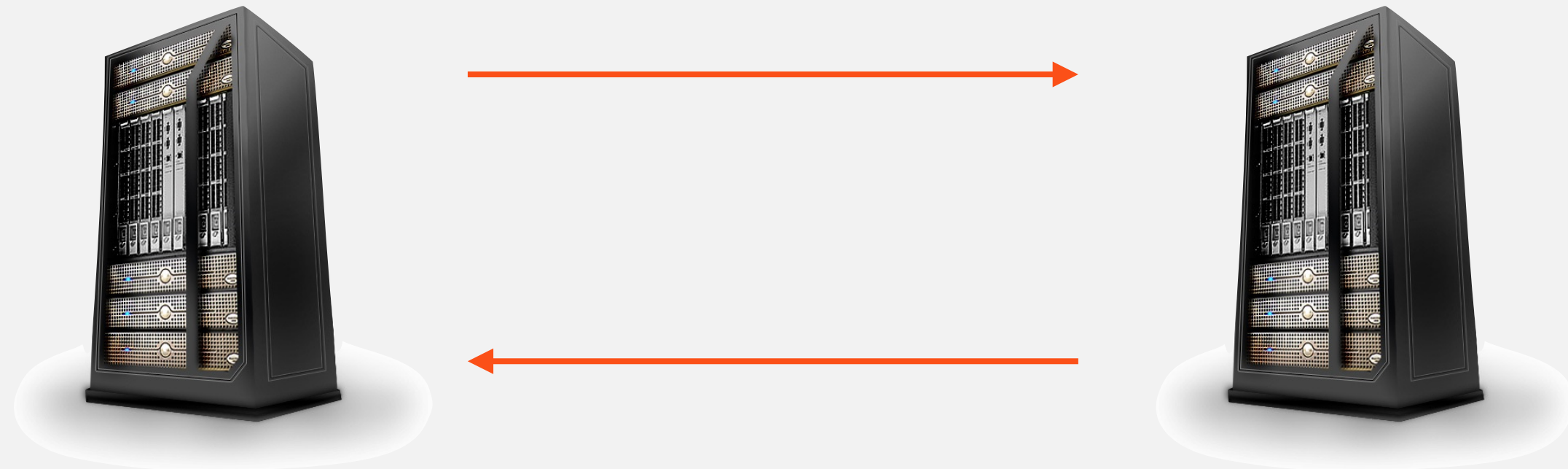


Встроенный прикладной кеш
custom in-memory store



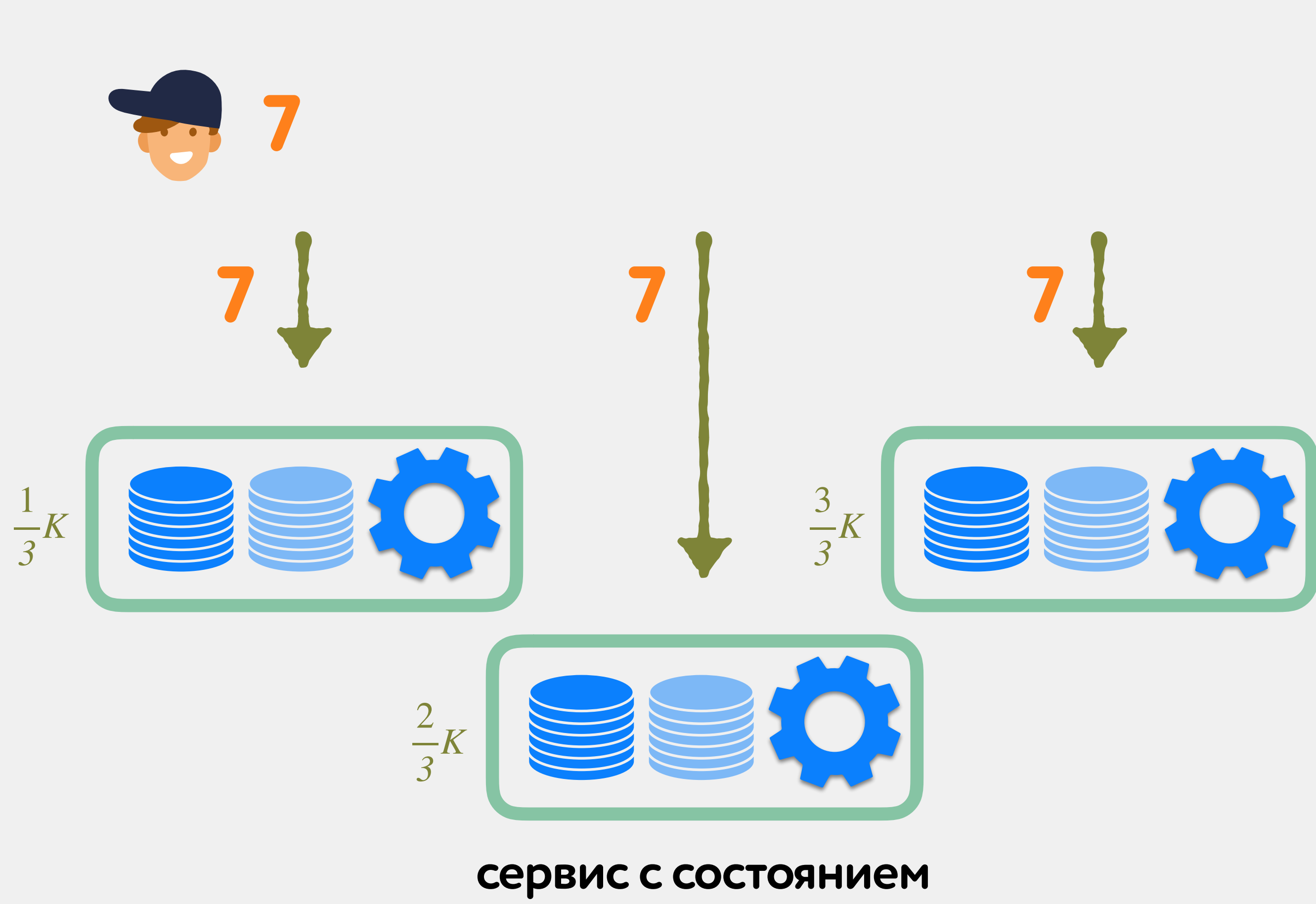
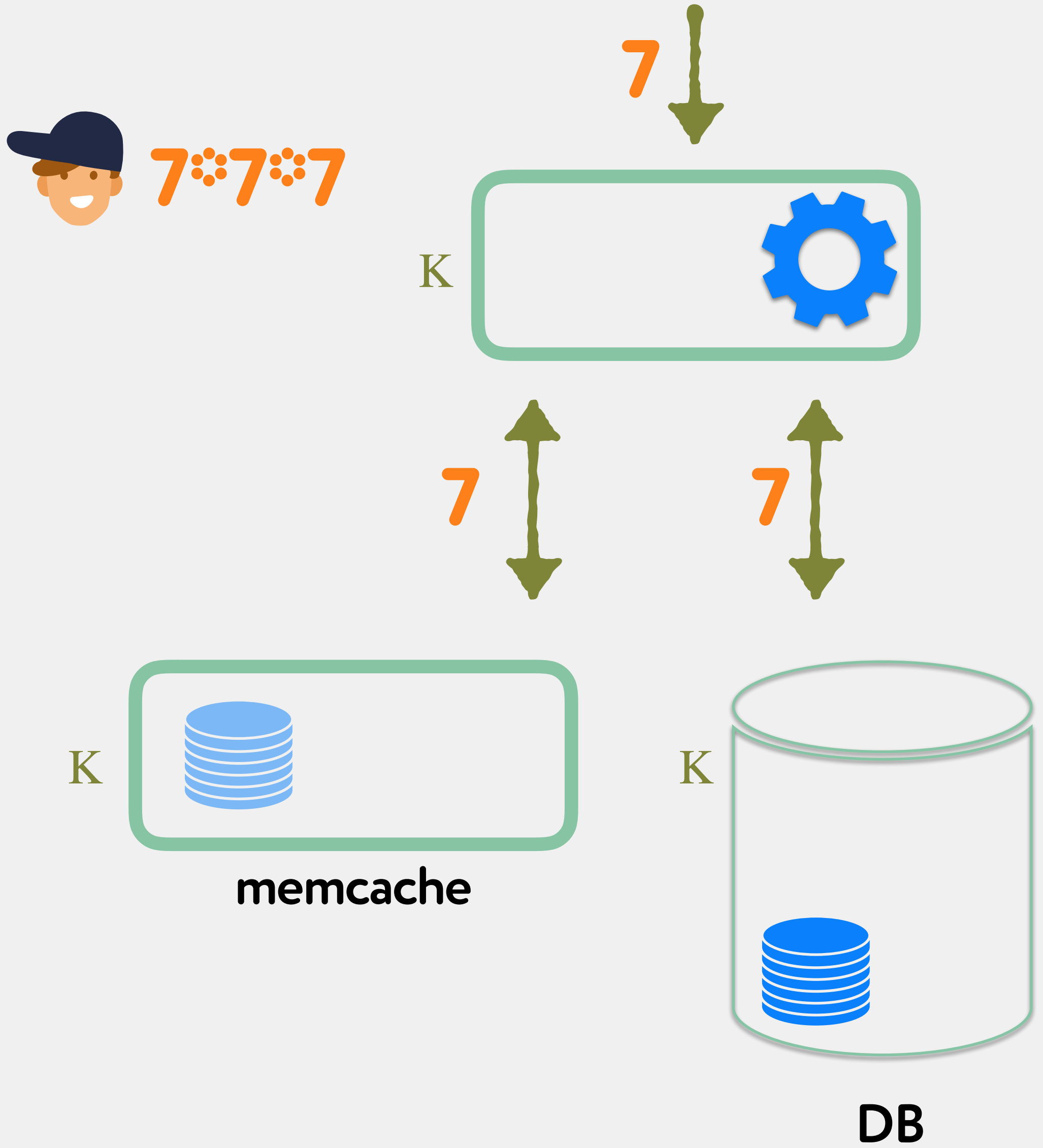
Встроенная распределенная БД
embedded distributed store

Что может пойти не так ?



1. Пропажа клиента
2. Пропажа сервера
3. Потеря исходящего сообщения
4. Потеря входящего сообщения
5. Таймаут сервера
6. Неправильный ответ
7. Произвольный отказ

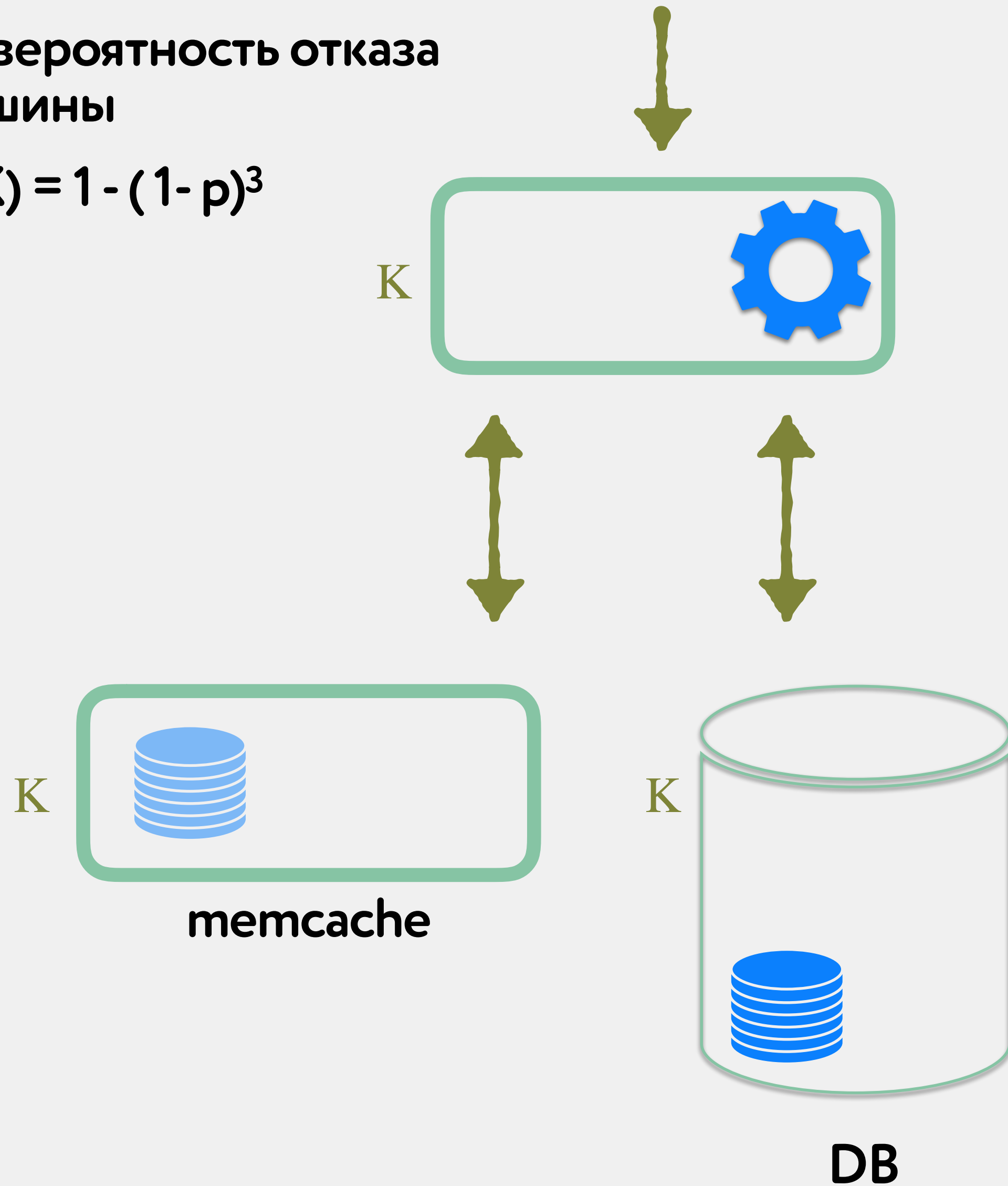
Сценарии отказов



Вероятности отказов

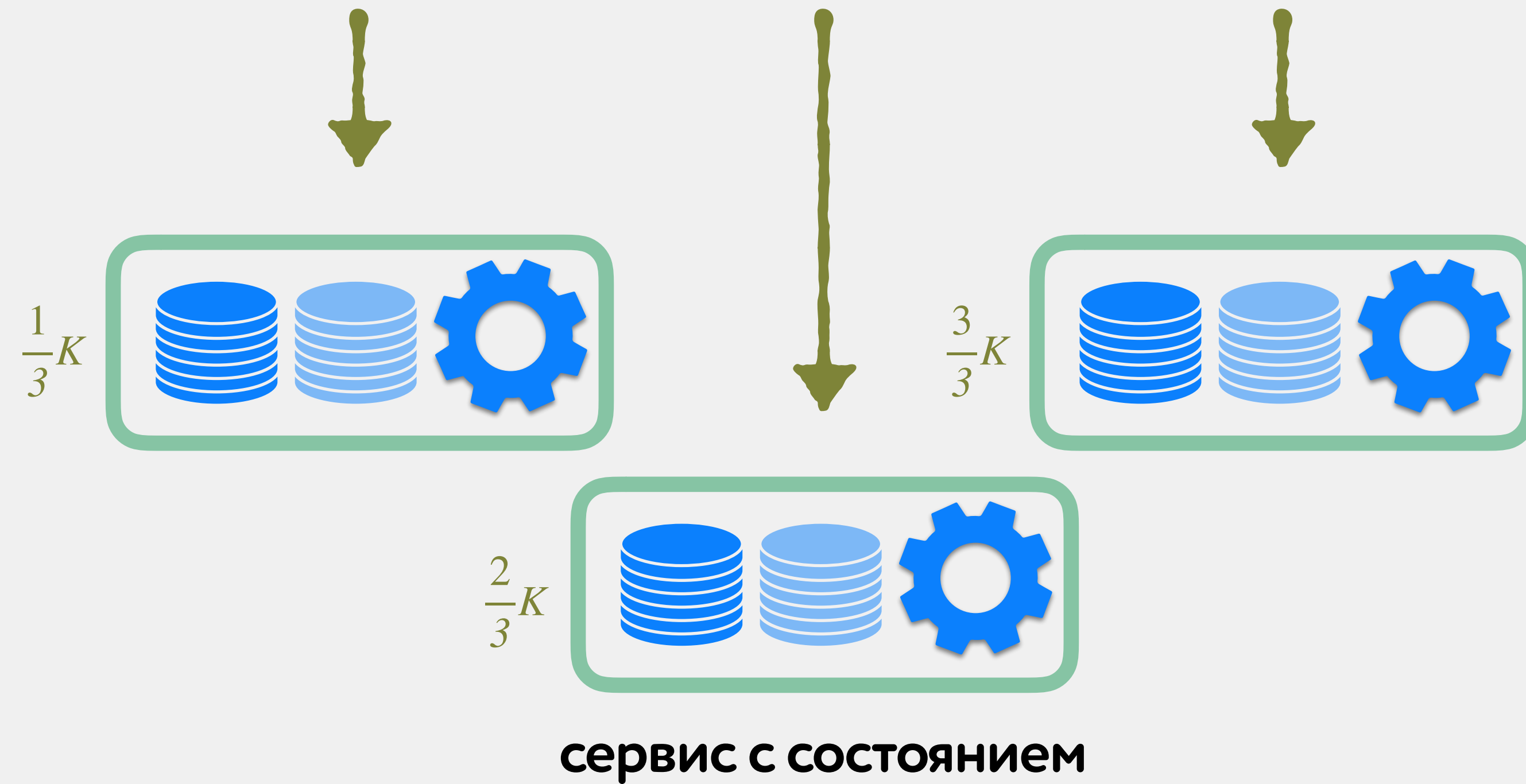
p - вероятность отказа
машины

$$P(K) = 1 - (1 - p)^3$$



$$P(K) = p^3$$

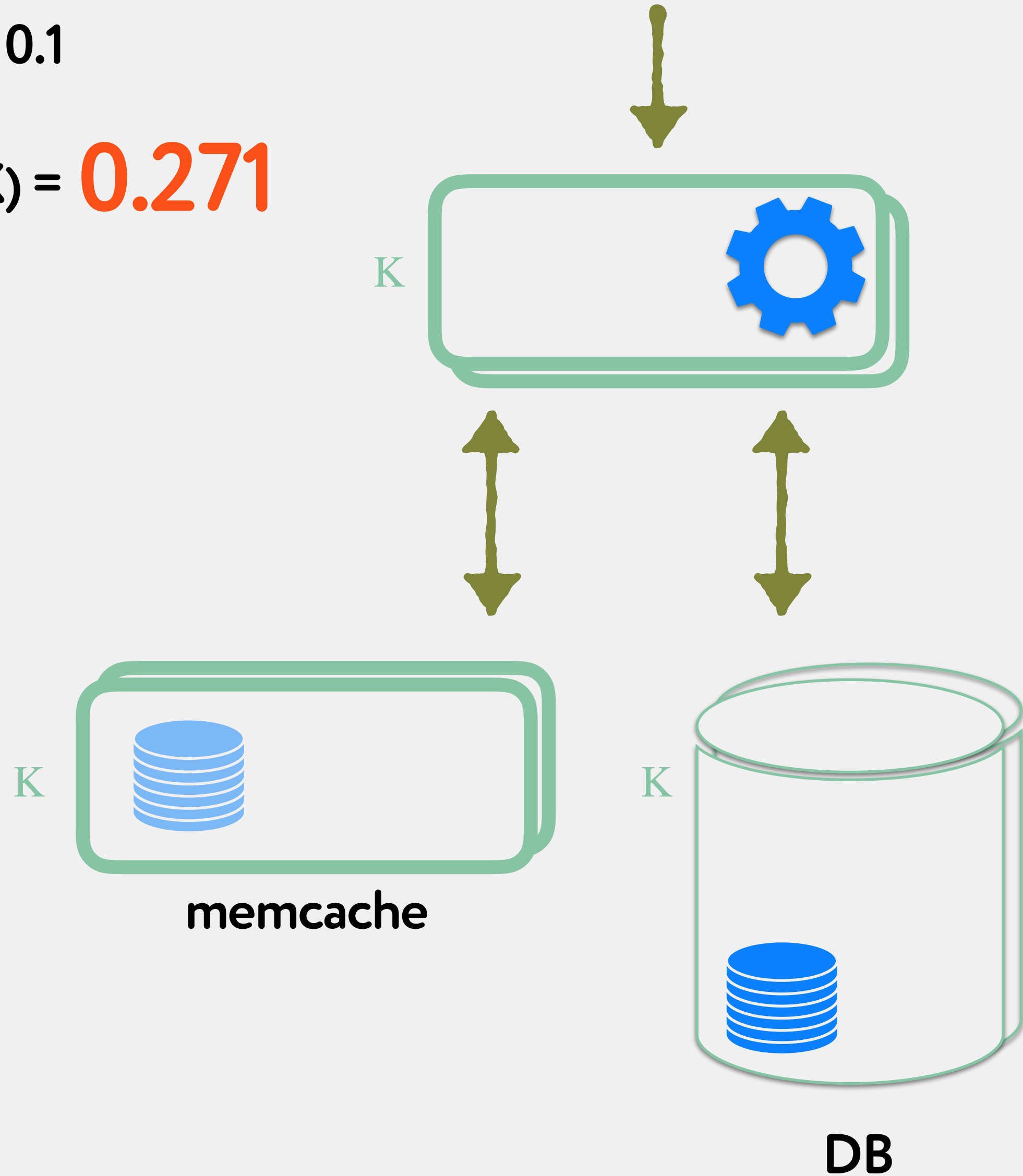
$$P(\frac{1}{3} K) = 1 - (1 - p)^3$$



Вероятности отказов

$p = 0.1$

$P(K) = 0.271$



$P(K) = 0.001$

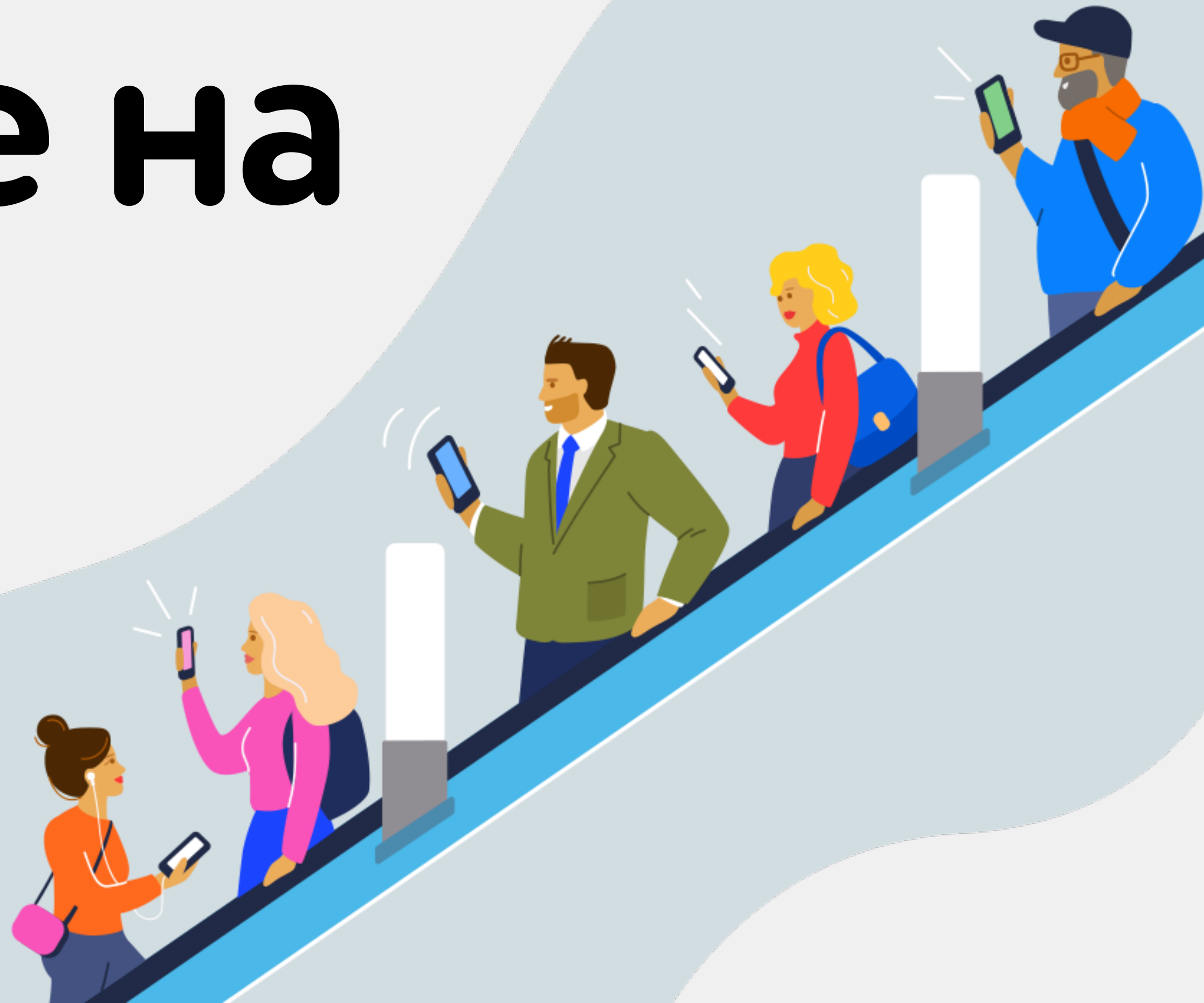
$P(\frac{1}{3}K) = 0.271$

Всегда надежнее!





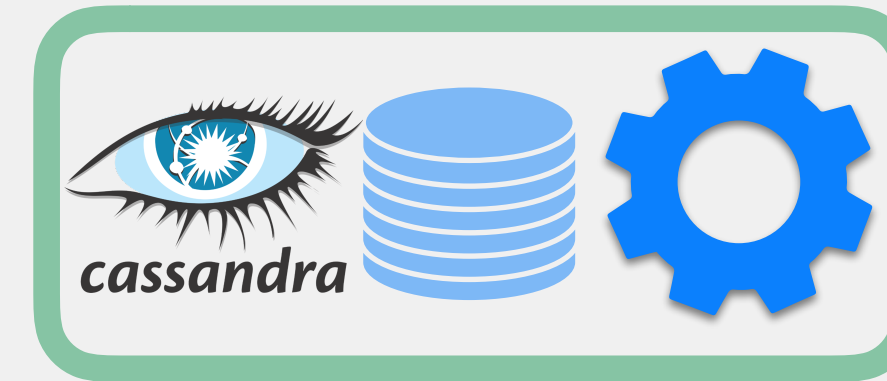
Хорошо все на бумаге



Встраиваем БД в сервис

должно быть:

- **Высокодоступным**
Репликация, консистентность
- **Масштабируемым**
Решардинг
- **На языке приложения**
Минимизация (де)маршалинга,
Интеграция с приложением
- **Open Source**
для доработок



Встраиваем БД в сервис

1. `-cp cassandra/lib/*.jar`

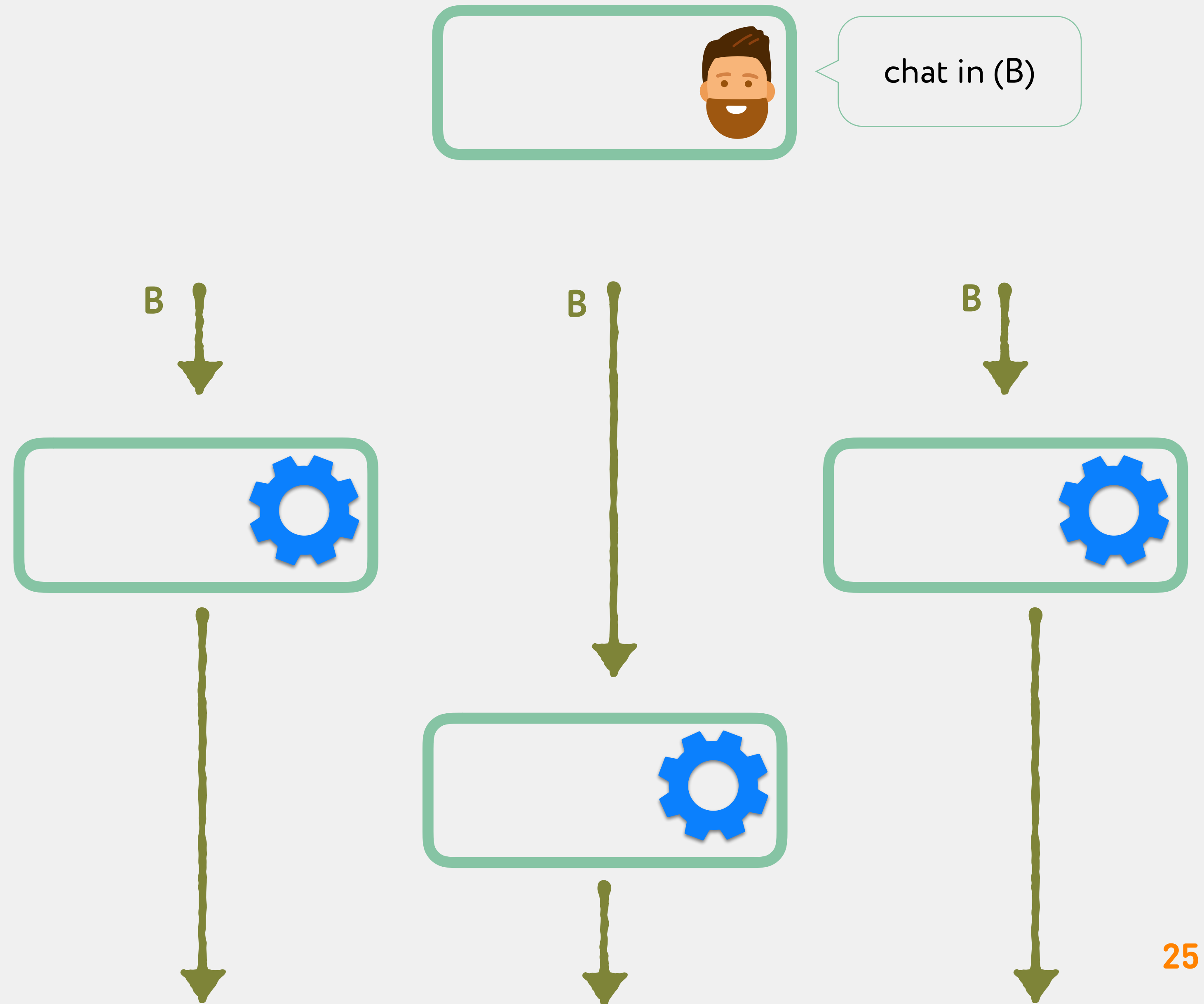
```
package org.apache.cassandra.service;  
  
public class CassandraDaemon  
{
```

2. `System.setProperty("cassandra.config", "file:///whatever/cassandra.yaml");`
`CassandraDaemon.instance.activate();`

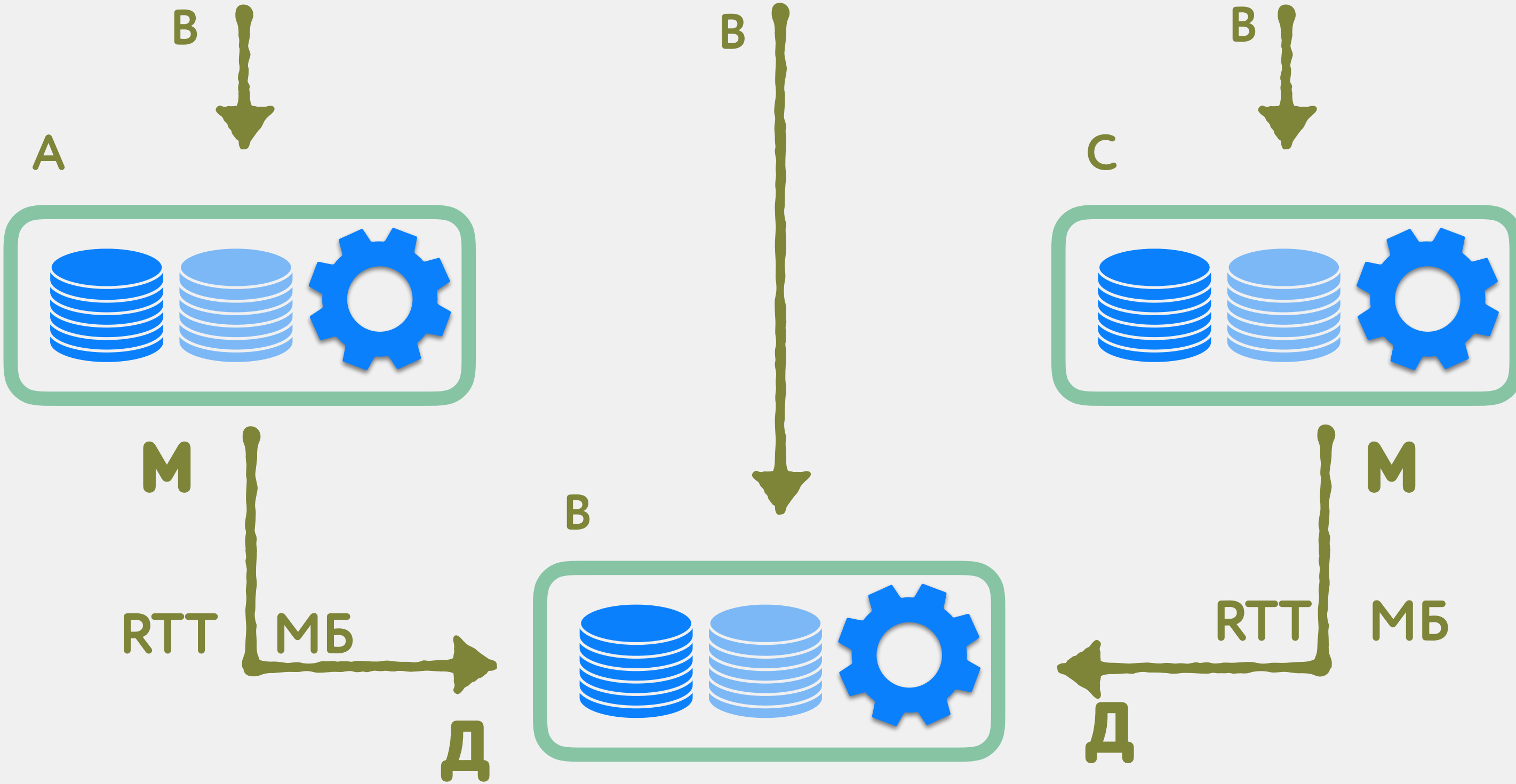
3.

```
public interface ConfigurationLoader  
{  
    Config loadConfig() ;  
}
```


Маршрутизация запросов



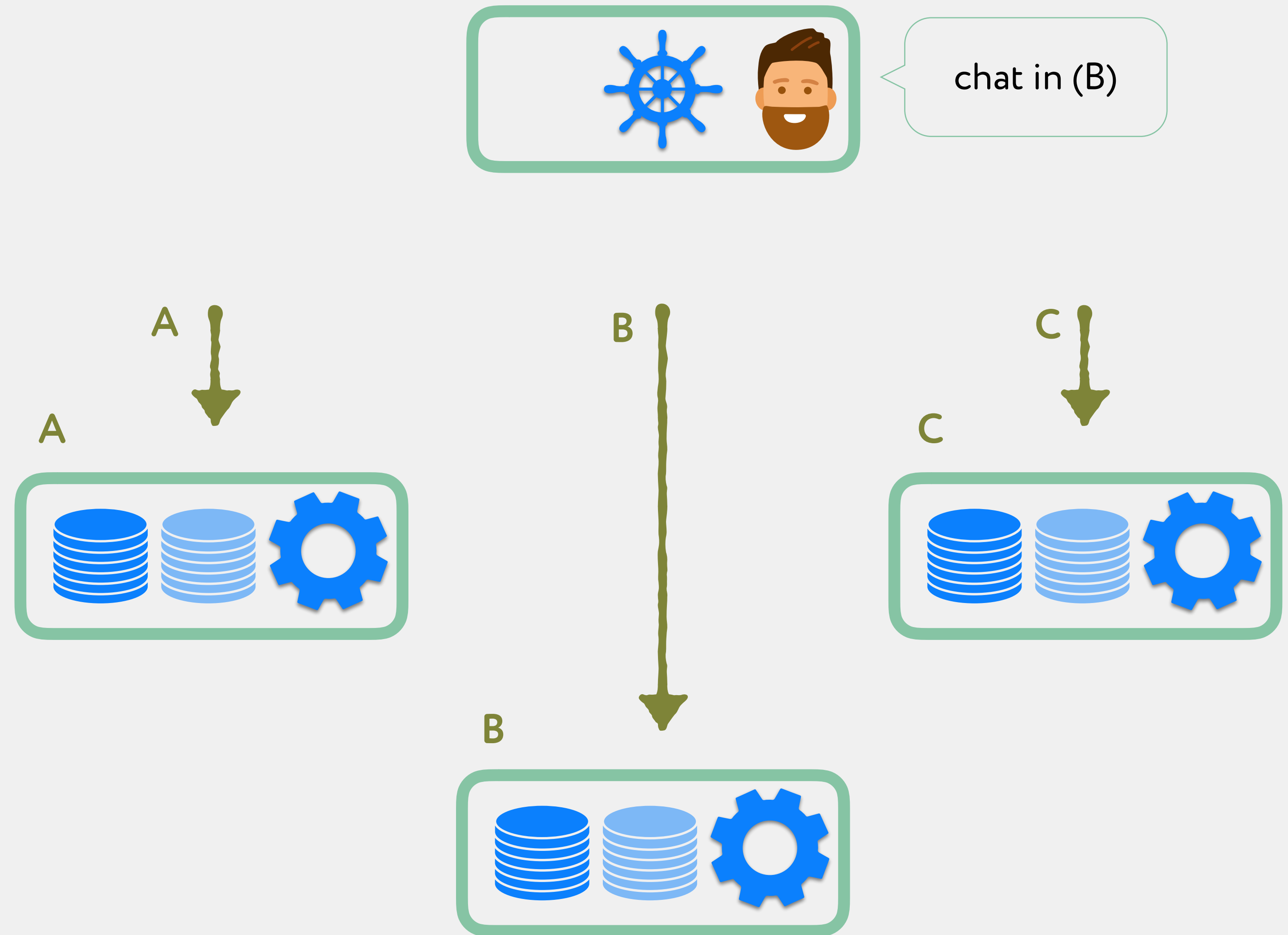
Маршрутизация запросов



Маршрутизация запросов

- **Partition-aware client routing library**

Обеспечивает вызов реплики владеющей данными по ключу на основании информации о кластере



Распределение данных

- **Partition Key (chatId)**

Определяет положение записи на ноде

- **Clustering Key (msgId)**

Сортирует записи внутри партиции

```
CREATE TABLE Messages (  
  chatId, msgId  
  
  user, type, text, attachments[], terminal, deletedBy[], replyTo  
  
  PRIMARY KEY ( chatId, msgId )  
)
```

Распределение данных

- **Partitioner**

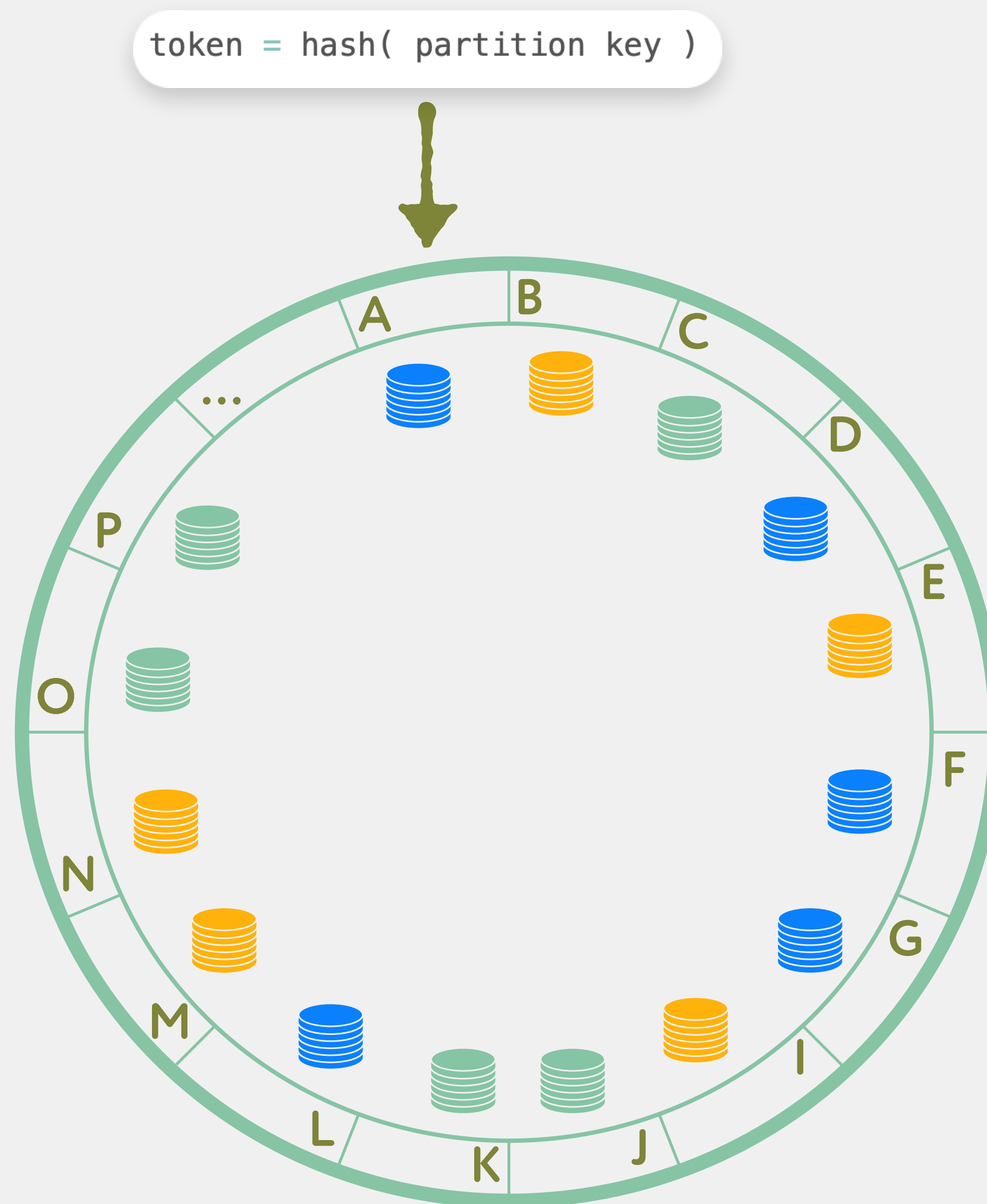
Вычисляет токен, положение на кольце

- **TokenMetadata**

Сопоставляет первичный интервал
нодам кластера

- **Replication Strategy**

Определяет распределение реплик
данных



Распределение данных

- **Partitioner**

Вычисляет токен, положение в кольце

- **TokenMetadata**

Сопоставляет первичный интервал
нодам кластера

- **Replication Strategy**

Определяет распределение реплик
данных

```
SortedMap<Token, List<InetAddress>> endpointMap = ...
```

```
AbstractReplicationStrategy replication = ...
```

```
for ( Token token : tokenMetadata.sortedTokens() ) {  
    endpointMap.put( token, replication.getNaturalEndpoints( token ) );  
}
```

+ Изменение топологии

Обновление,
устаревшая информация

Мессенджер: работаем с БД

- `getMessages(viewer, chat, from, to)`
- `add(chat, message)`

```
CREATE TABLE Messages (  
    chatId, msgId  
  
    user, type, text, attachments[], terminal, delete  
  
    PRIMARY KEY ( chatId, msgId )  
)
```

```
package org.apache.cassandra.cql3;  
  
import java.nio.ByteBuffer;  
  
public class QueryProcessor  
{  
    public static UntypedResultSet execute(String query,  
                                            ConsistencyLevel cl, Object... values)  
        throws RequestExecutionException
```

```
UntypedResultSet rs = QueryProcessor.execute(  
    "SELECT * FROM Messages "  
    + "WHERE chatId = ? AND msgId < ? AND msgId > ?",  
    ConsistencyLevel.QUORUM, chatId, from, to );  
  
rs.forEach( row -> {} );
```

Кеш сообщений

600 млрд
сообщений

5 млрд
чатов

100 ТБ

5%
активных чатов

80%
последние 13 сообщений



кеш сообщений в памяти

3+ млрд
сообщений

250 млн
чатов

500 ГБ

Кеш сообщений

- `getMessages(viewer, chat, from, to)`
- `getLastMessages(viewer, chats)`
- `add(chat, message)`



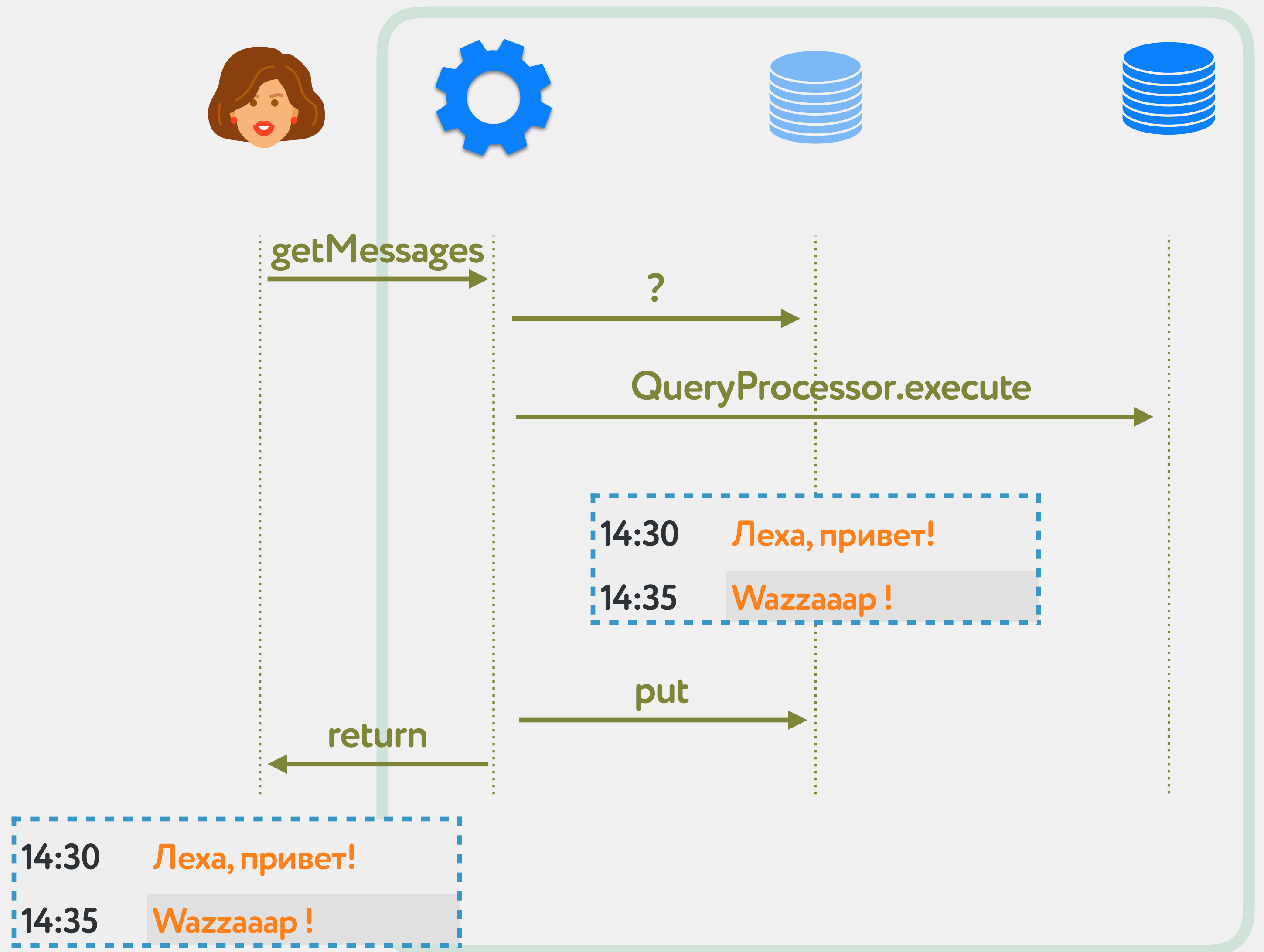
кеш сообщений в памяти

3+ млрд
сообщений

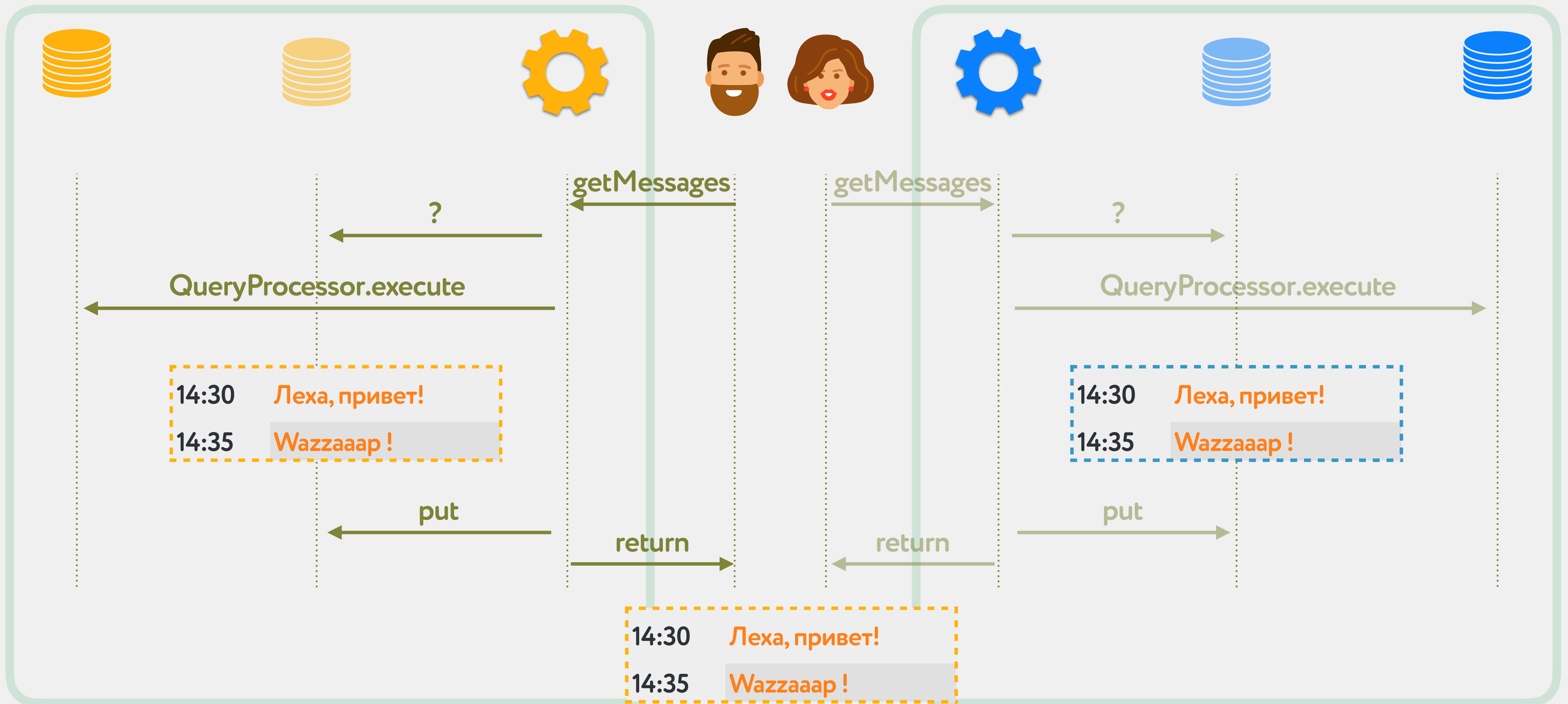
250 млн
чатОВ

500 ГБ

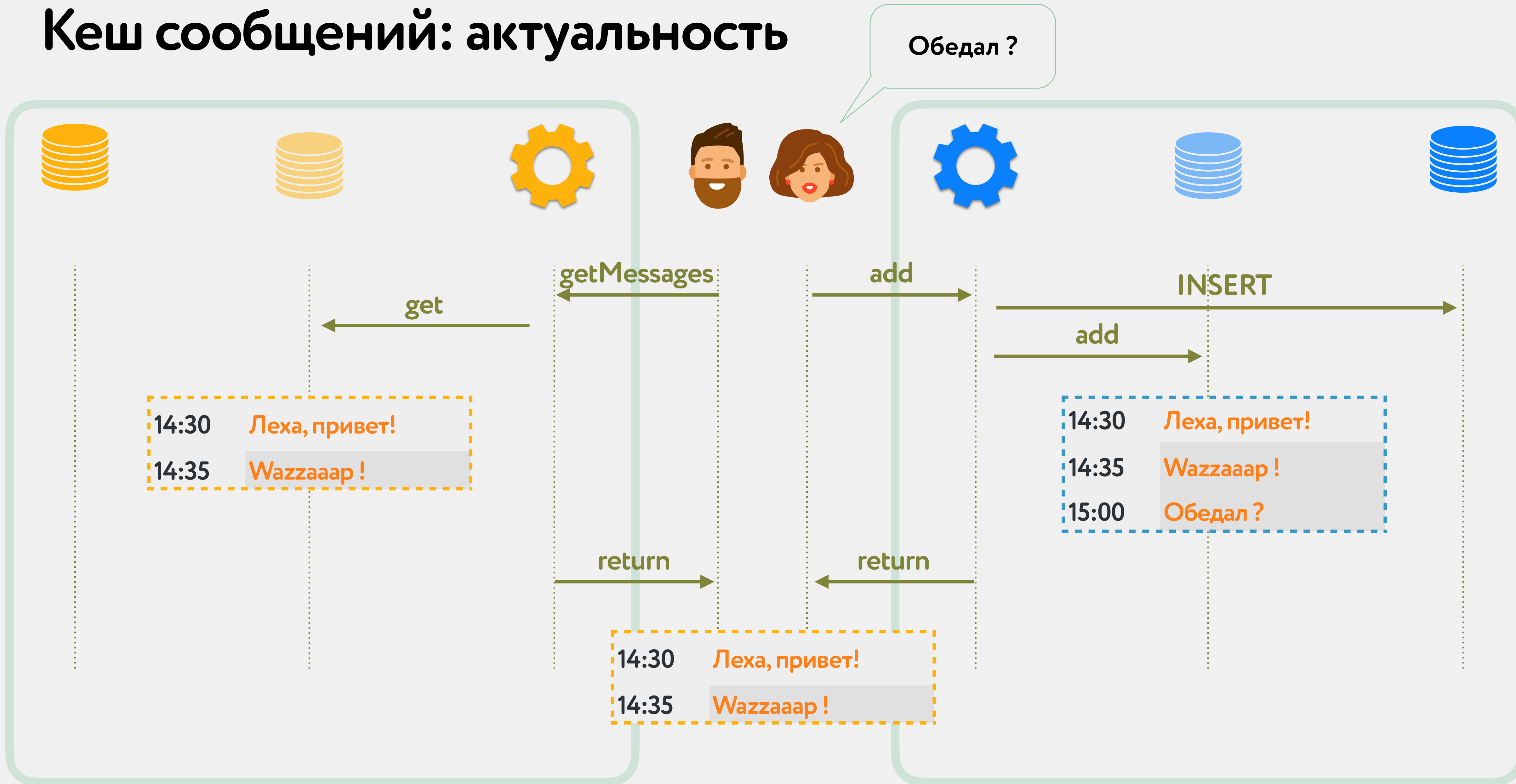
Кеш сообщений: getMessages



Кеш сообщений

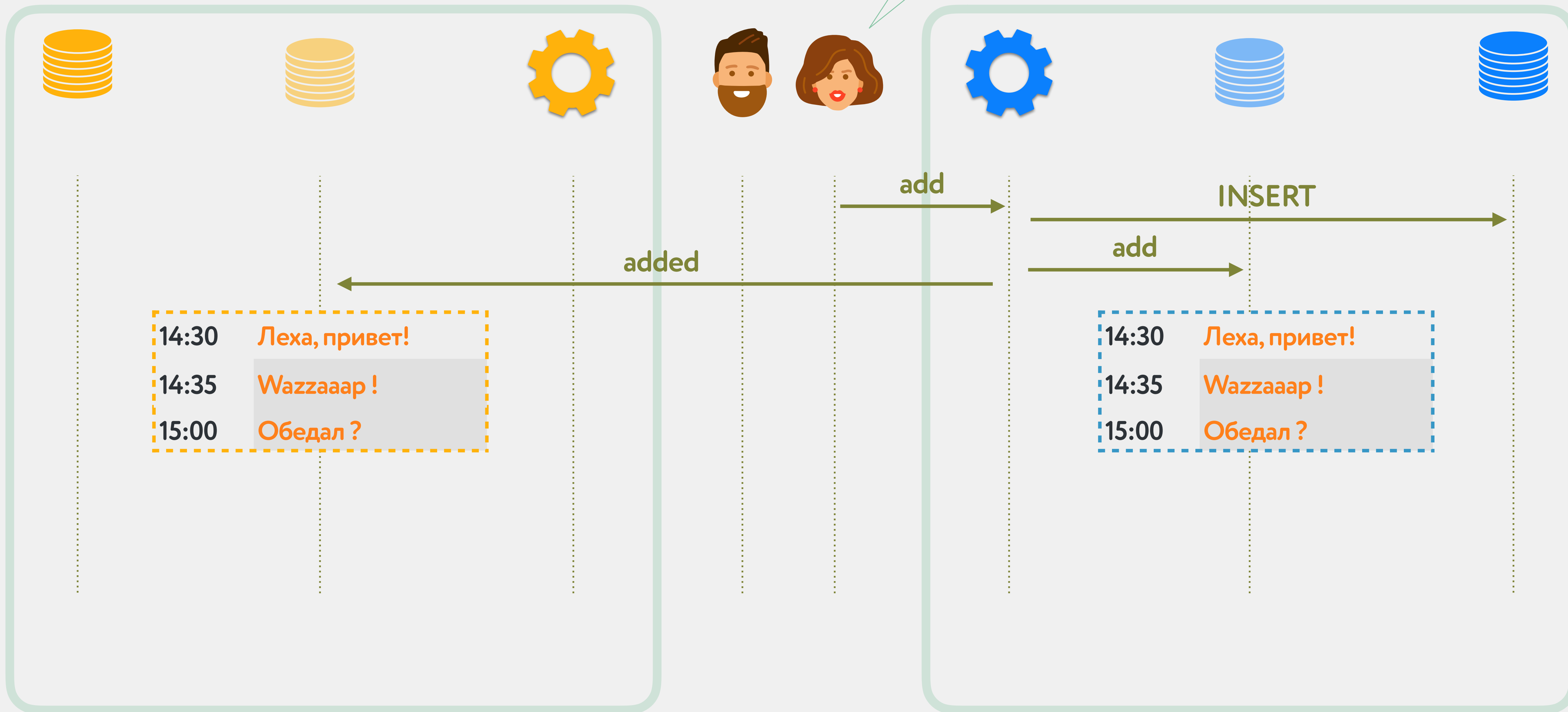


Кеш сообщений: актуальность



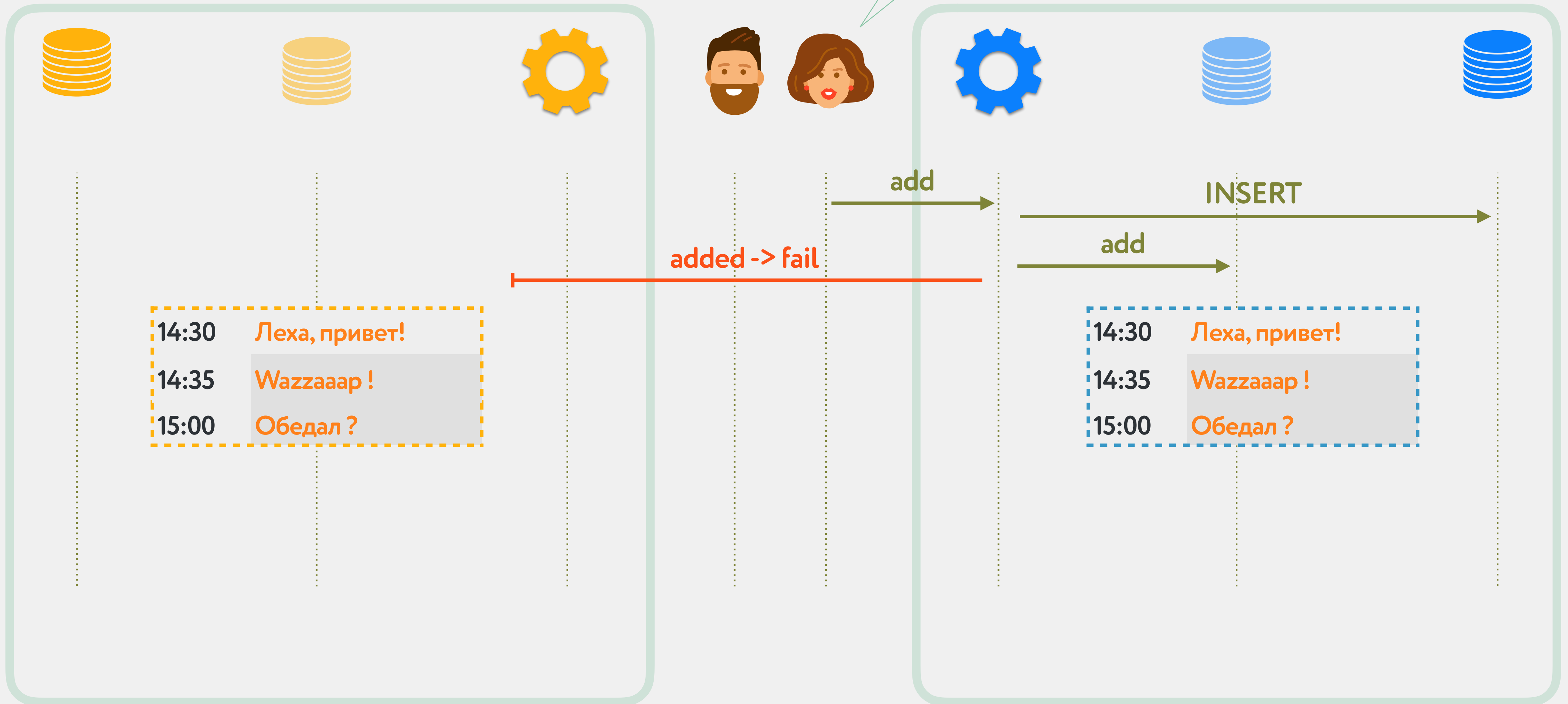
Кеш сообщений: актуальность

Обедал ?



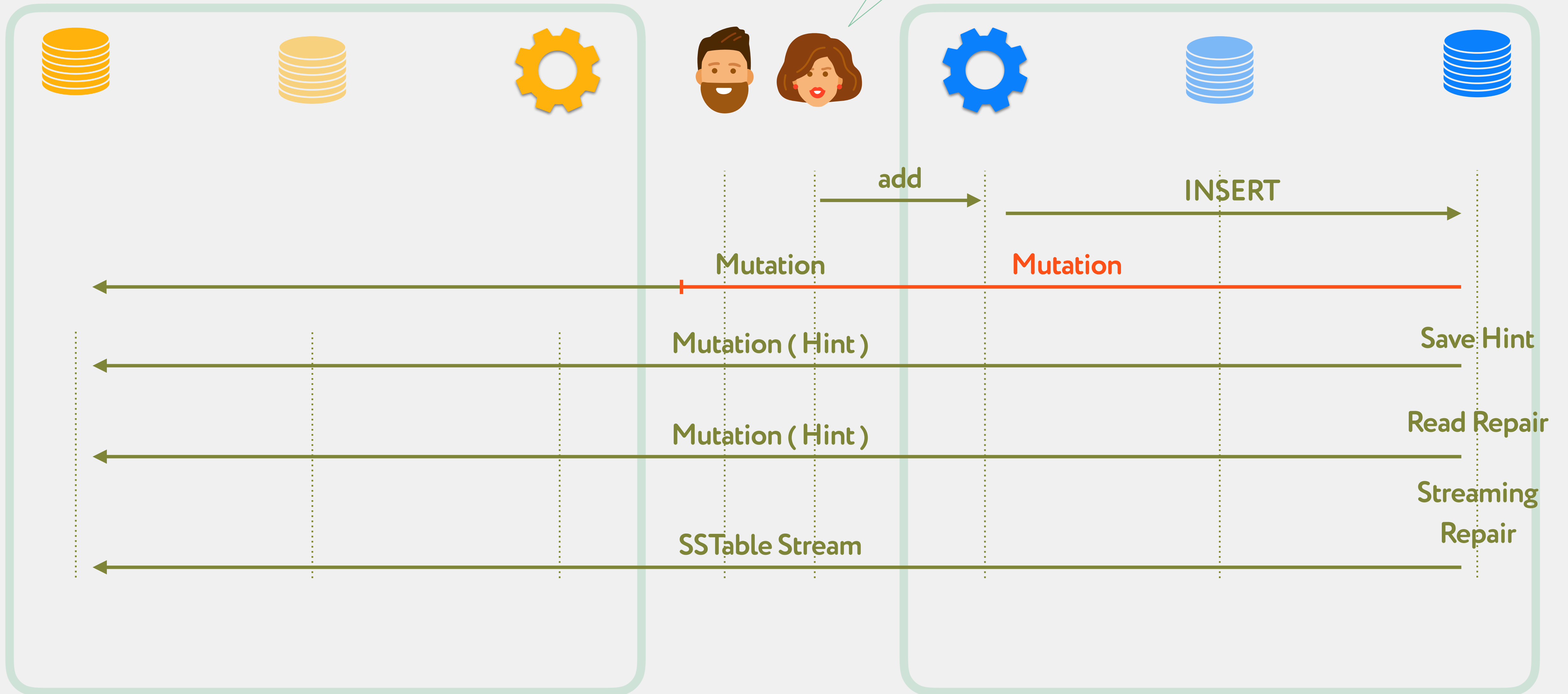
Кеш сообщений: актуальность

Обедал ?

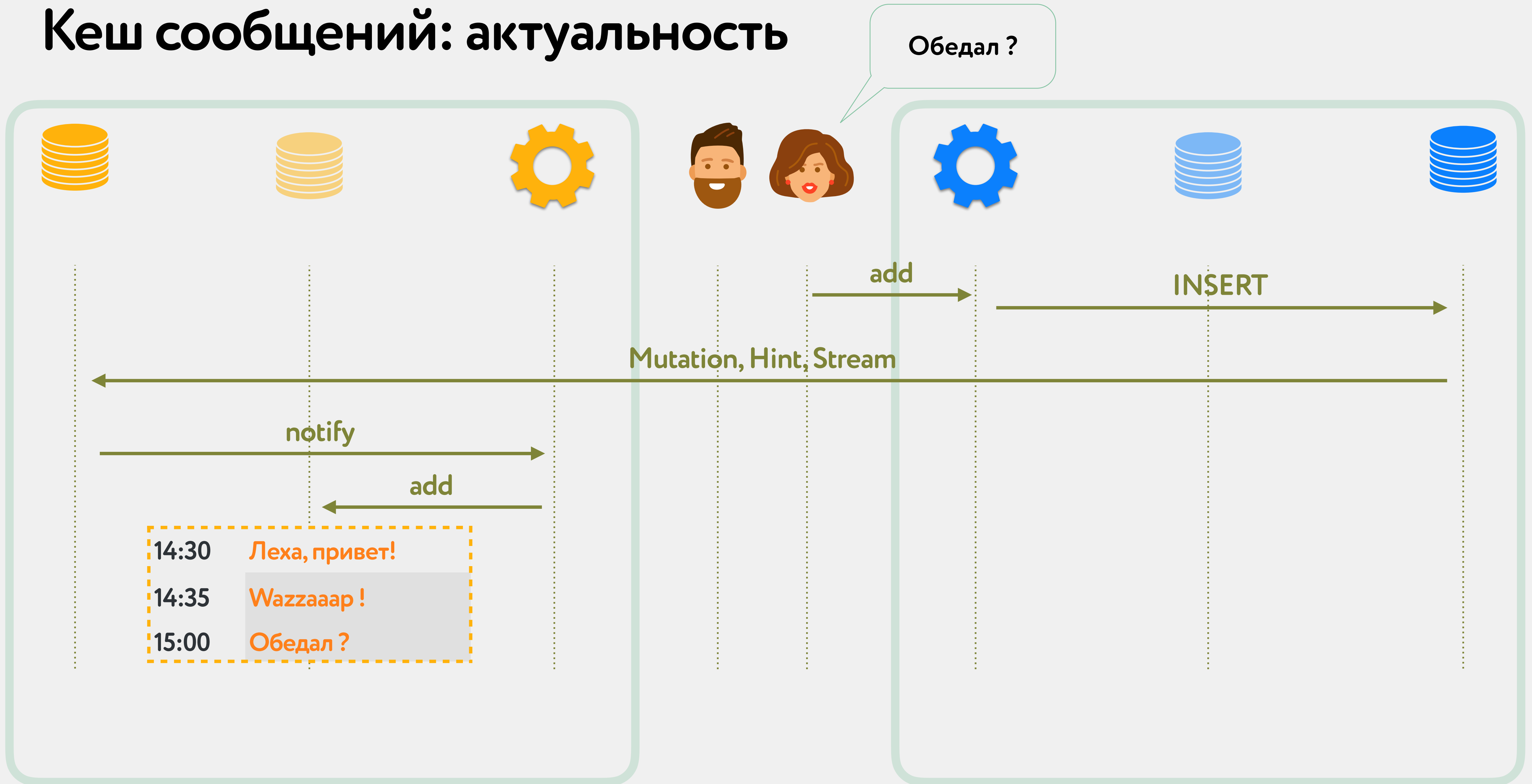


Кеш сообщений: актуальность

Обедал ?



Кеш сообщений: актуальность



Mutation - записи:

```
interface ApplyMutationListener
{
    void onApply(ByteBuffer key,
                 DeletionTime deletion,
                 Iterator<Unfiltered> atoms);
}
```

```
package org.apache.cassandra.db;

public class Keyspace
{
    public void apply( Mutation mutation,
                     boolean writeCommitLog,
                     boolean updateIndexes,
                     boolean isDroppable )
```

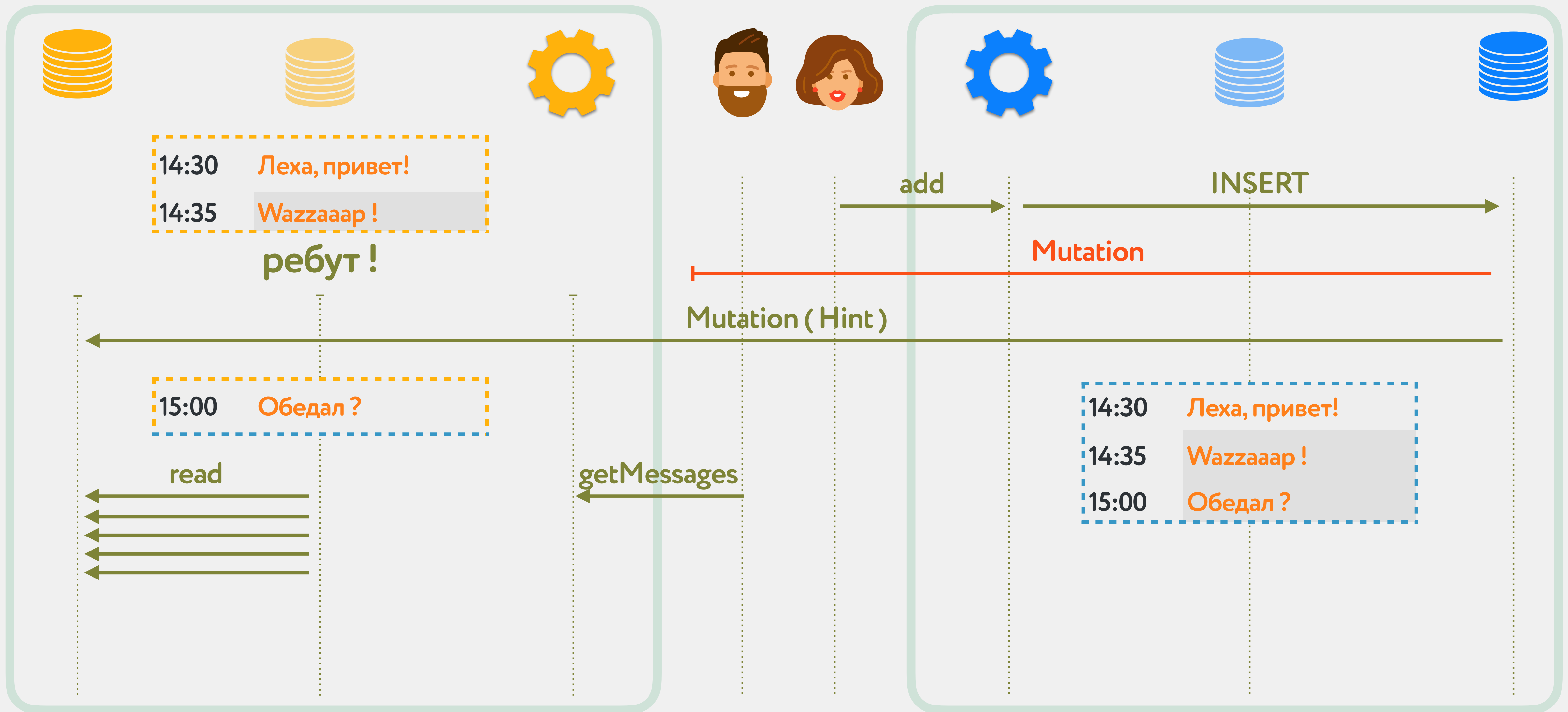
Streaming - таблицы:

```
package org.apache.cassandra.streaming;

public class StreamReceiveTask extends StreamTask
{
    // holds references to SSTables received
    protected Collection<SSTableReader> sstables;

    private static class OnCompletionRunnable implements Runnable {
```

Кеш сообщений: потеря состояния



Кеш сообщений: потеря состояния



```
CREATE KEYSPACE Caches  
  WITH REPLICATION = {  
    'class': 'LocalStrategy'  
  }
```

```
CREATE TABLE Caches.MessagesSnapshot (  
  rowkey blob,  
  value blob,  
  PRIMARY KEY ( rowkey )  
)
```

```
SELECT * FROM MessagesSnapshot
```

Кеш: оптимизируем рестарты

- **Разделяемая память**

Shared Memory

- **/dev/shm/messages-cache.mem**

но это необязательно

- **tmpfs**

- **hugetlbfs**

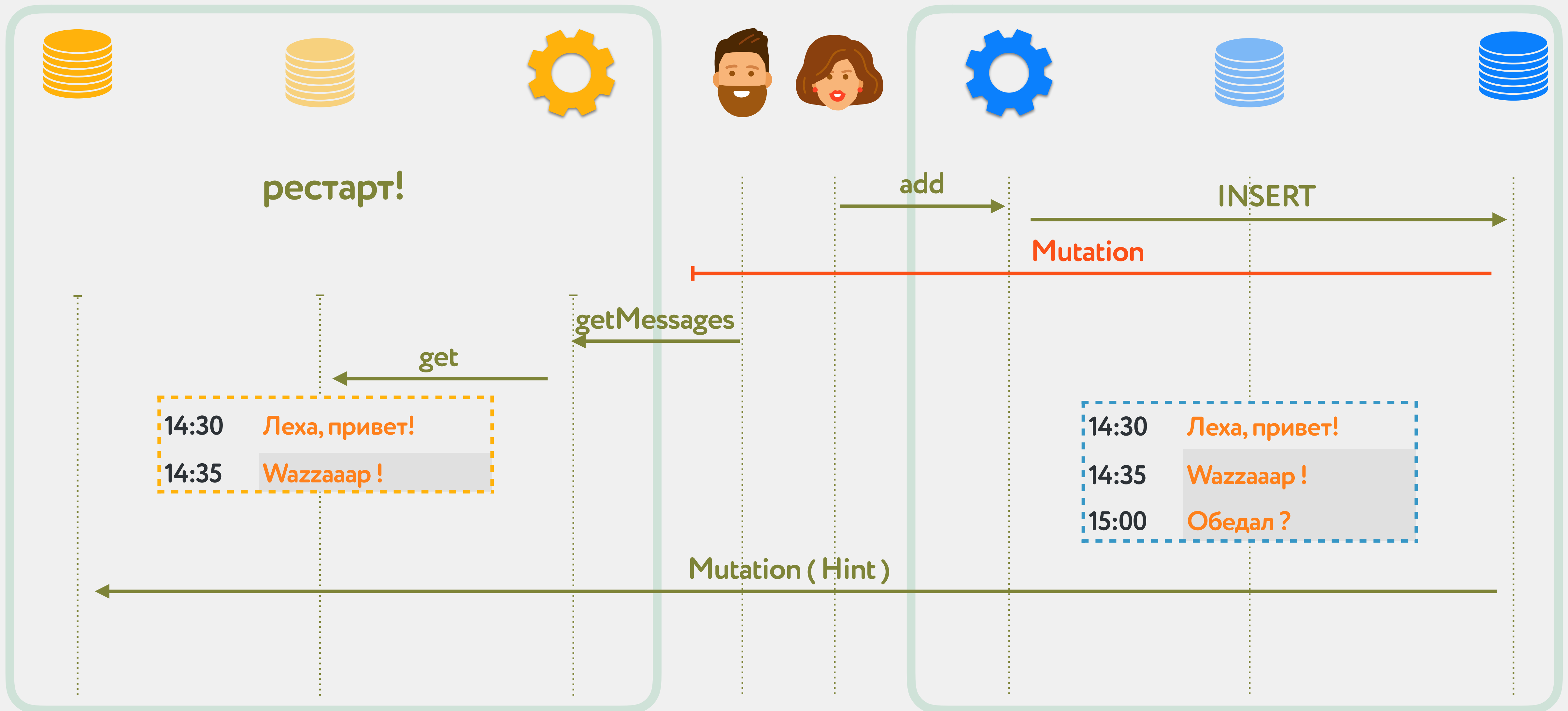
4К страницы -> 2М, 1Г

[**https://github.com/odnoklassniki/one-nio**](https://github.com/odnoklassniki/one-nio)

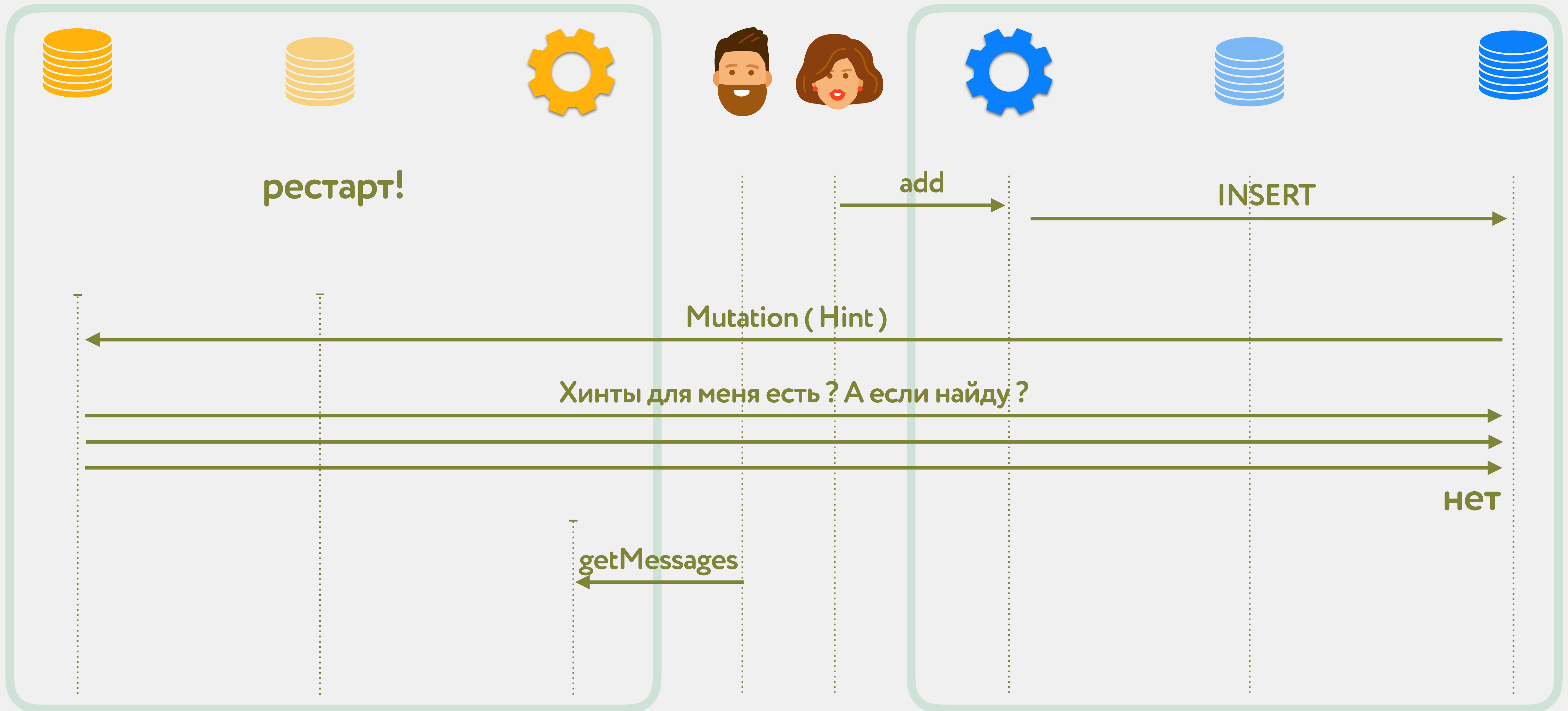
one.nio.mem

SharedMemoryMap

Кеш сообщений: ожидание согласованности



Кеш сообщений: ожидание согласованности





Мы хотим большего!



getLastMessages(chats[])

- **Множество чатов**

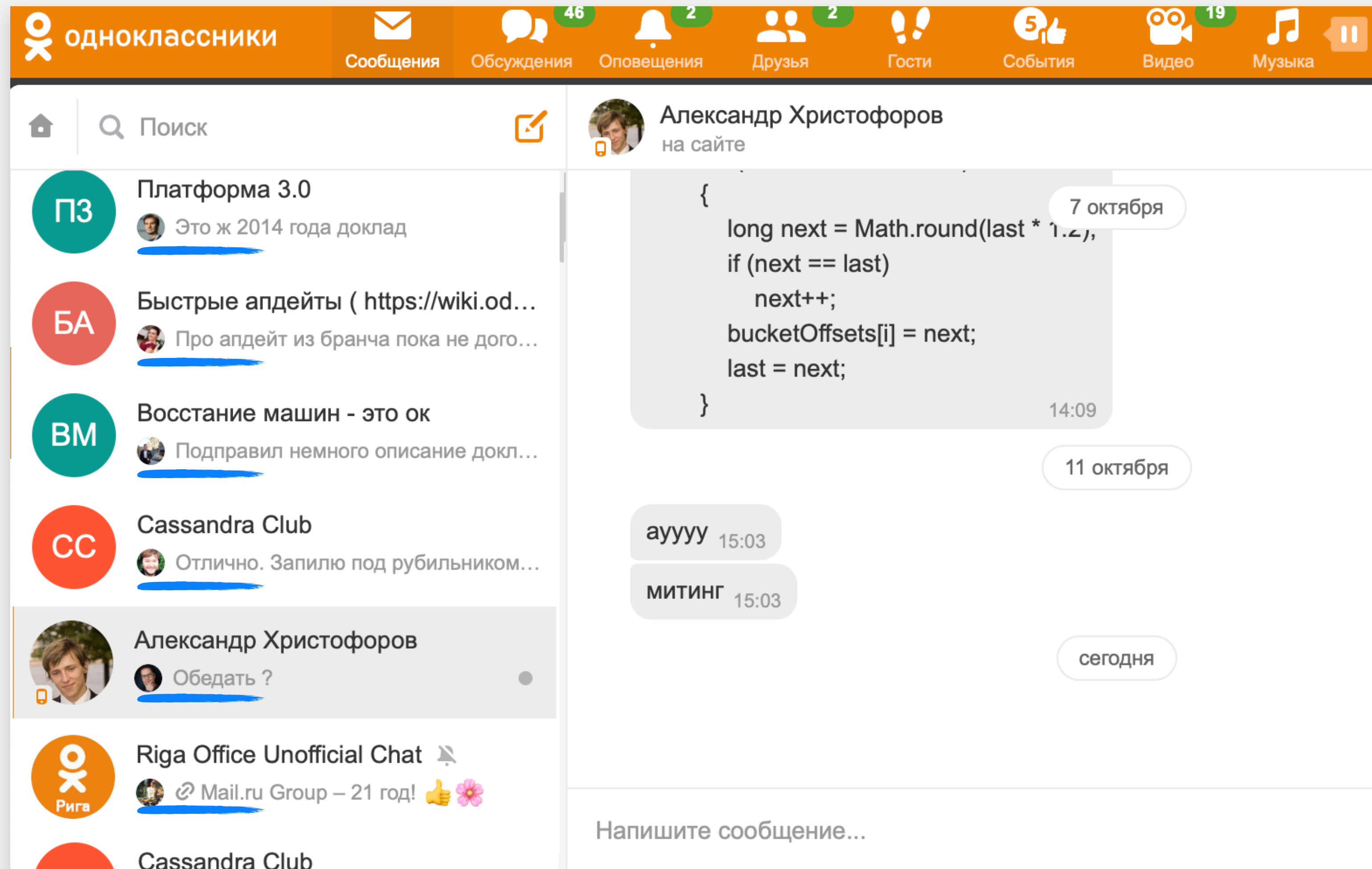
Нет 1 ноды владеющей необходимыми данными

- **Часть в кеше часть нет**

Грузить в кеш старые чаты смысла мало

- **Старые чаты не в кеше**

Они согласованы



getLastMessages(chats[])

- **Множество чатов**

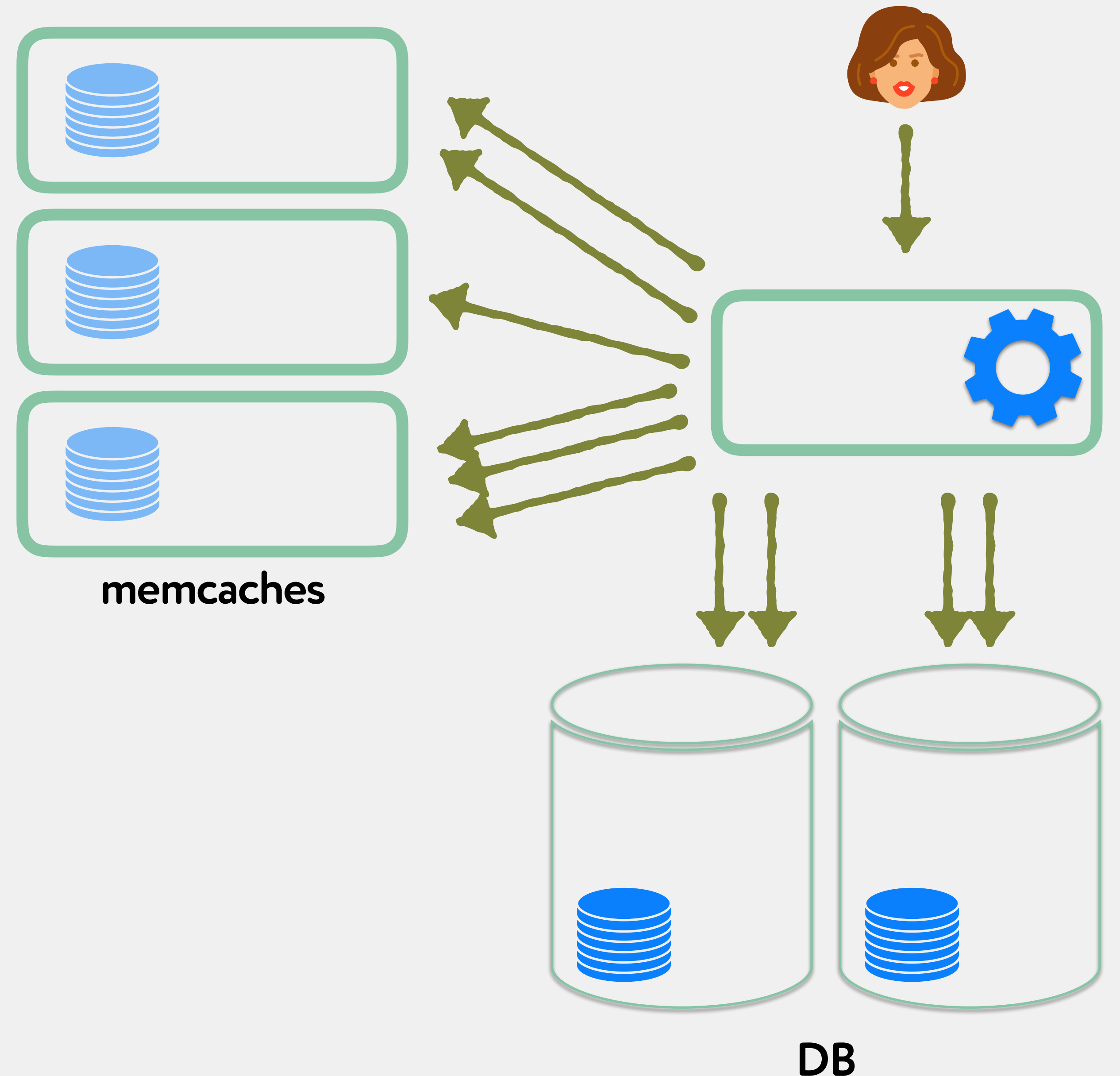
Нет 1 ноды владеющей необходимыми данными

- **Часть в кеше часть нет**

Грузить в кеш старые чаты смысла мало

- **Старые чаты не в кеше**

Они согласованы



split & merge

```
Map<Long, Message> getLastMessages( Long[] chatIds )
```

- **Множество чатов**

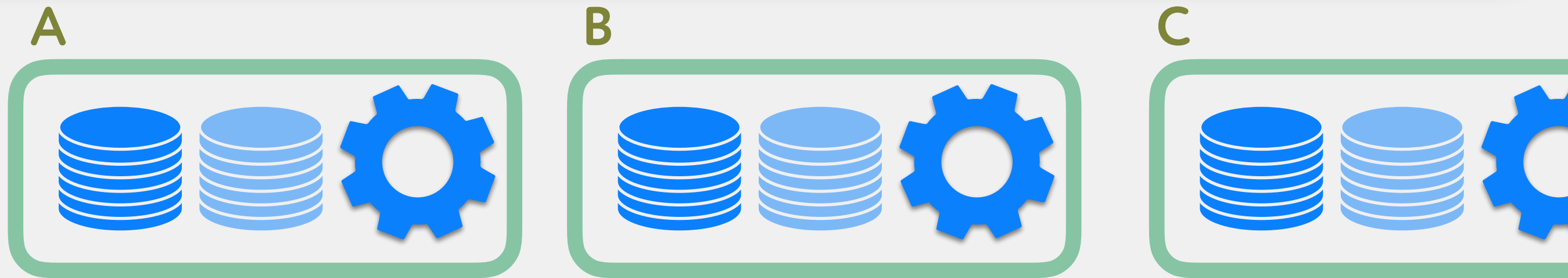
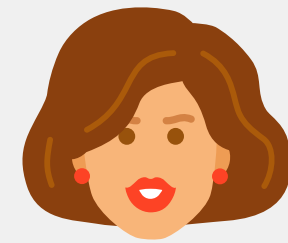
Нет 1 ноды владеющей необходимыми данными

- **Часть в кеше часть нет**

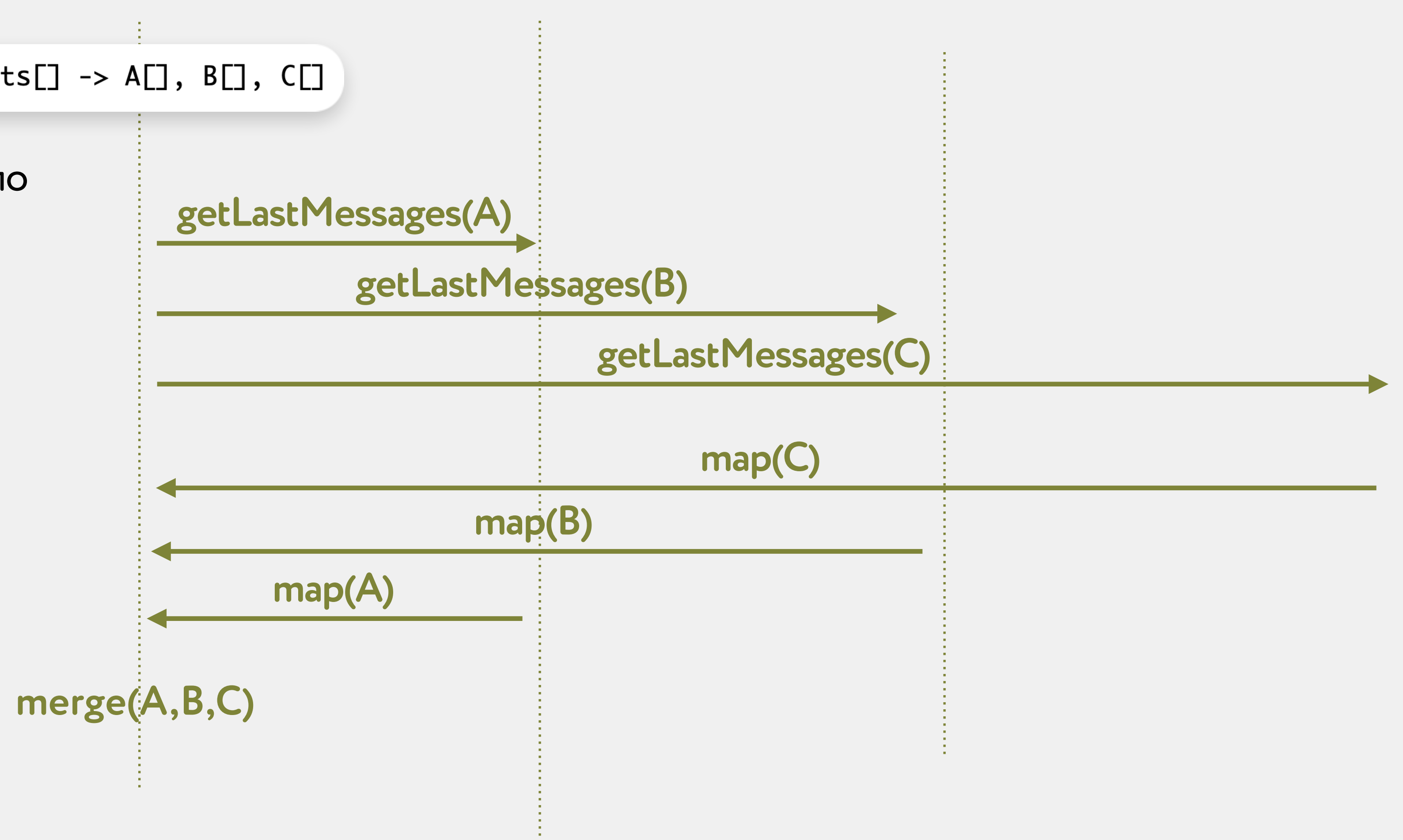
Грузить в кеш старые чаты смысла мало

- **Старых чатов больше**

Согласованность не так важна



chats[] -> A[], B[], C[]



getLastMessages: локальное чтение

- **Множество чатов**

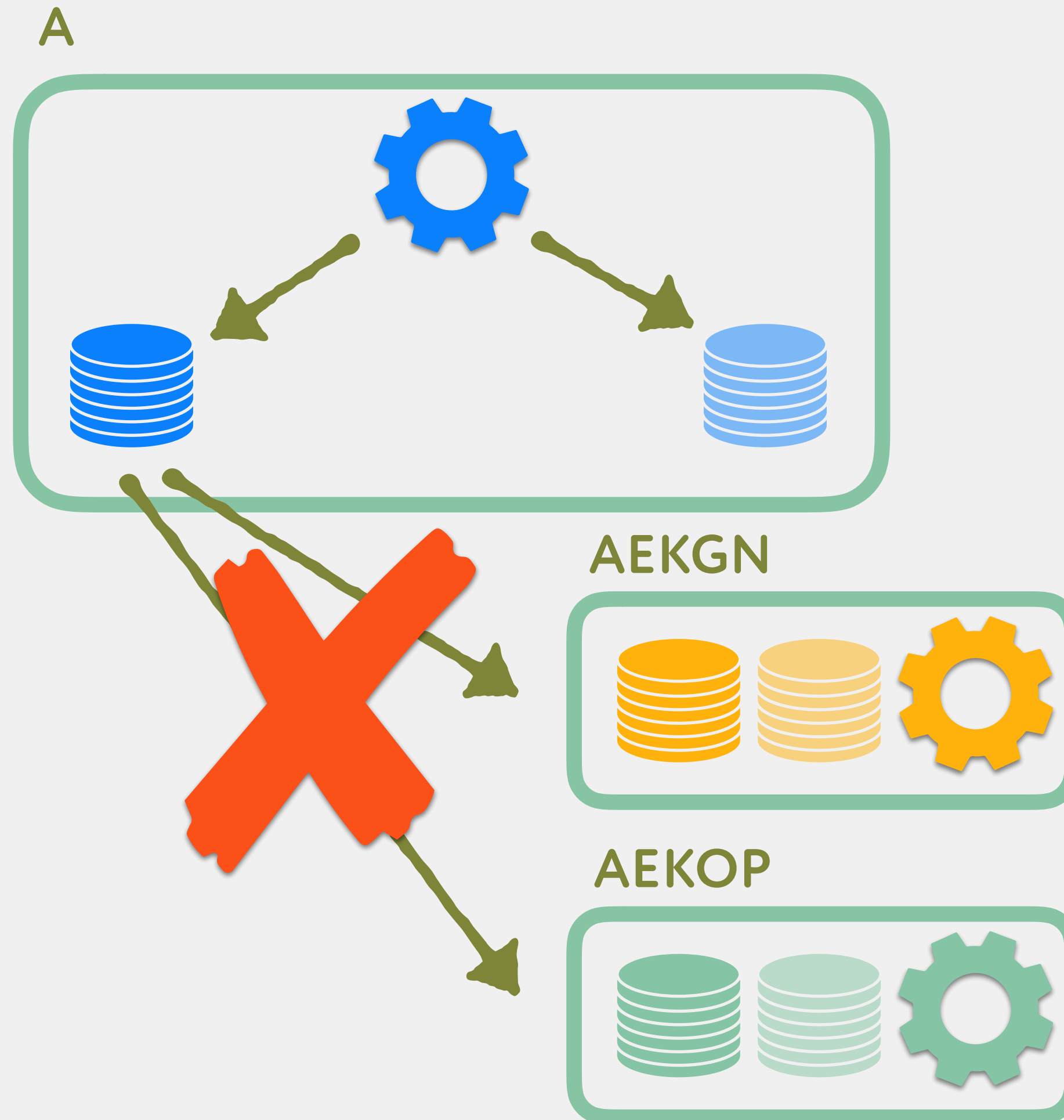
Нет 1 ноды владеющей необходимыми данными

- **Часть в кеше часть нет**

Грузить в кеш старые чаты смысла мало

- **Старых чатов больше**

Согласованность не так важна



Локальное чтение

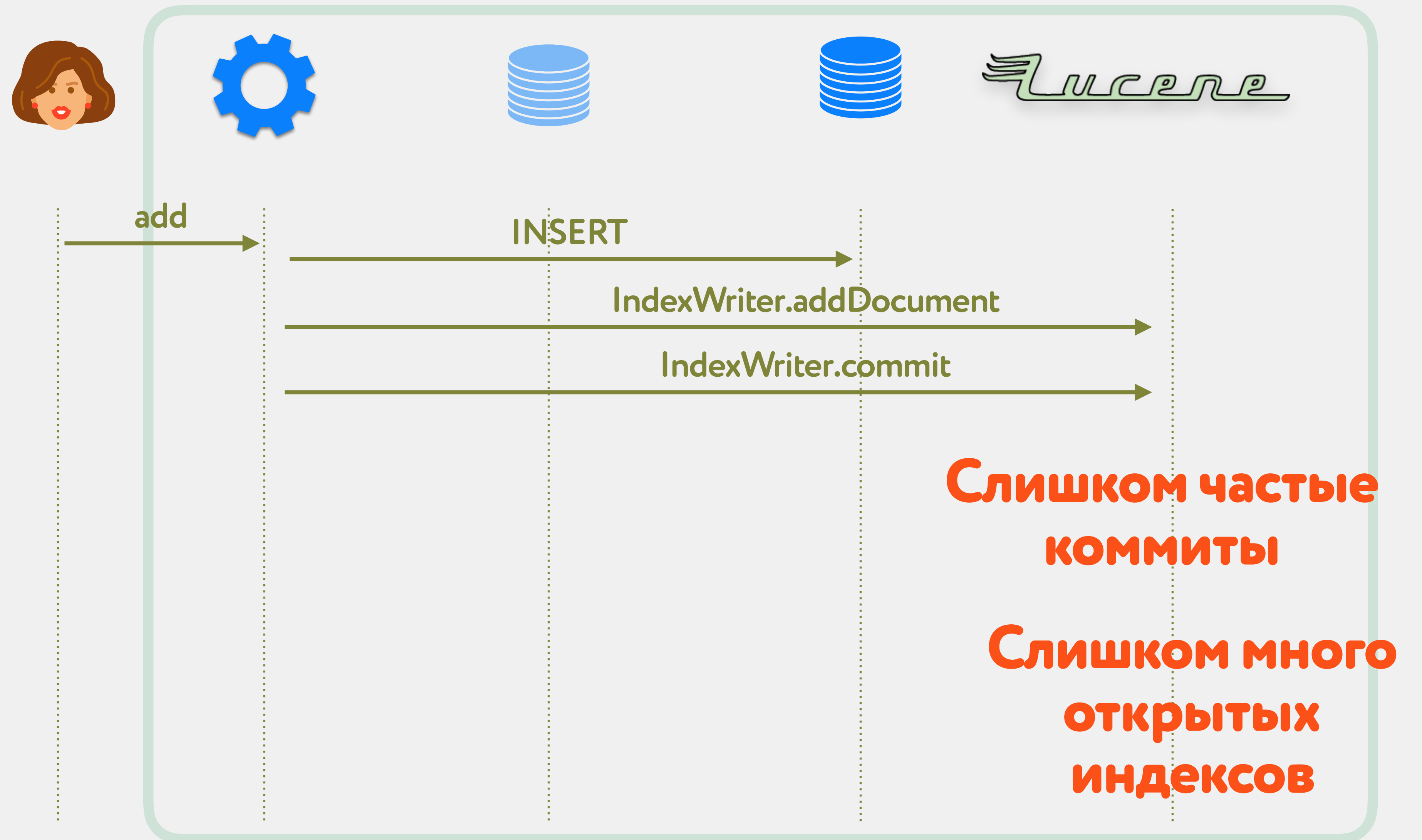
```
public Message getLastMessage(Long chatId) {  
  
    Prepared prepared  
        = QueryProcessor.prepareInternal( "SELECT ... WHERE chatId=?" );  
    SelectStatement select  
        = (SelectStatement) prepared.statement;  
  
    QueryOptions opt  
        = QueryProcessor.makeInternalOptions( prepared, chatId );  
  
    ReadQuery query = select.getQuery( opt, FBUtilities.nowInSeconds() );  
  
    UnfilteredPartitionIterator partitions  
        = query.executeLocally( query.executionController() );  
  
    UnfilteredRowIterator partition = partitions.next();  
    partition.forEachRemaining( atom -> { /* unmarshalling code */ } );  
}
```

Полнотекстовый П

- **Строим индекс**
lucene.apache.org
- **Отдельный для каждого чата**
Один большой индекс не сработает
- **Для больших чатов**
Для небольших индекс можно строить непосредственно перед поиском

The screenshot shows the VK mobile app interface. At the top, there's a navigation bar with icons for Messages (46), Discussions, Notifications (2), Friends (2), Guests, Events (5), Video (19), and Music. Below the navigation bar, a search bar contains the text 'joker'. The search results list several chat groups, including 'ХабраЧат', 'Александр Анисимов', 'Восстание машин - это ок', 'ODKL:all IT сюда публикуется инт...', 'Cassandra unconf/meetup', and 'Cassandra Club'. The right side of the screen shows a chat conversation with 'Александр Христофоров на сайте'. The chat history includes a code block with Java code, a date separator for '7 октября', a message 'ауууу' at 15:03, a 'МИТИНГ' at 15:03, a date separator for '11 октября', and a 'сегодня' separator. At the bottom, there's a text input field with the placeholder 'Напишите сообщение...'.

Полнотекстовый П



Compaction

- **Слияние поколений данных**
В файлах на диске.
- **В порядке следования ключей**
PartitioningKey, Clustering Key

```
package org.apache.cassandra.db.compaction;  
  
public class CompactionManager implements CompactionManager  
{
```

**Бесплатная массовая
обработка данных**



Эксплуатация

4+



Более долгий старт

КТО ВИНОВАТ:

- **Инициализация БД**

Зависит от скорости дисков с данными
Скорости проигрывания WAL

- **Загрузка кеша из CF**

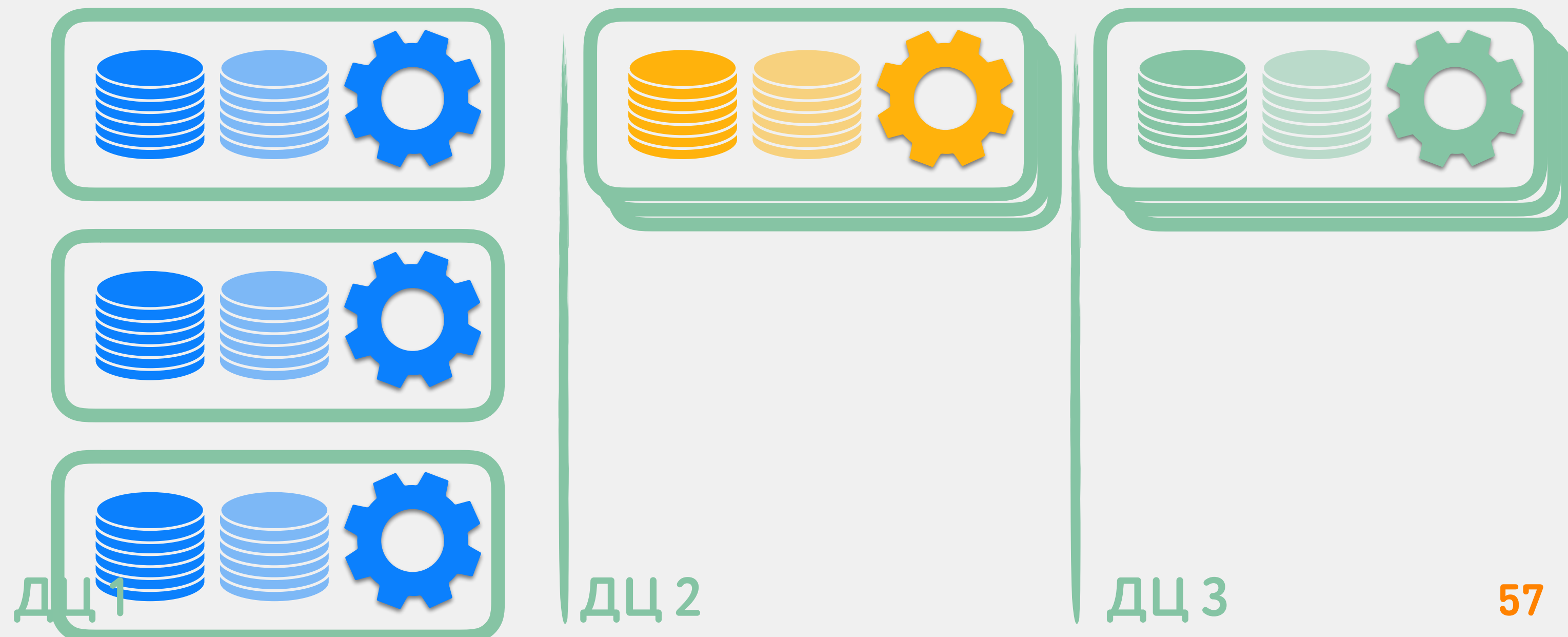
Размер кеша, contention, CPU

- **Ожидание согласованности**

Количество пропущенных мутаций

что делать:

- **Параллелизировать выкладку по зонам доступности**



Более частые рестарты БД

КТО ВИНОВАТ:

- БД теперь вместе с приложением

ЧТО ДЕЛАТЬ:

- Ничего, это хорошо

Позволяет отладить отказы зон доступности в контролируемой среде



Не трогай, это же
БАЗА ДАННЫХ, ее
нельзя рестартать - мы
же все п*ем¹ !!11

(1) Потеряем все данные



ДЦ 1



ДЦ 2



ДЦ 3

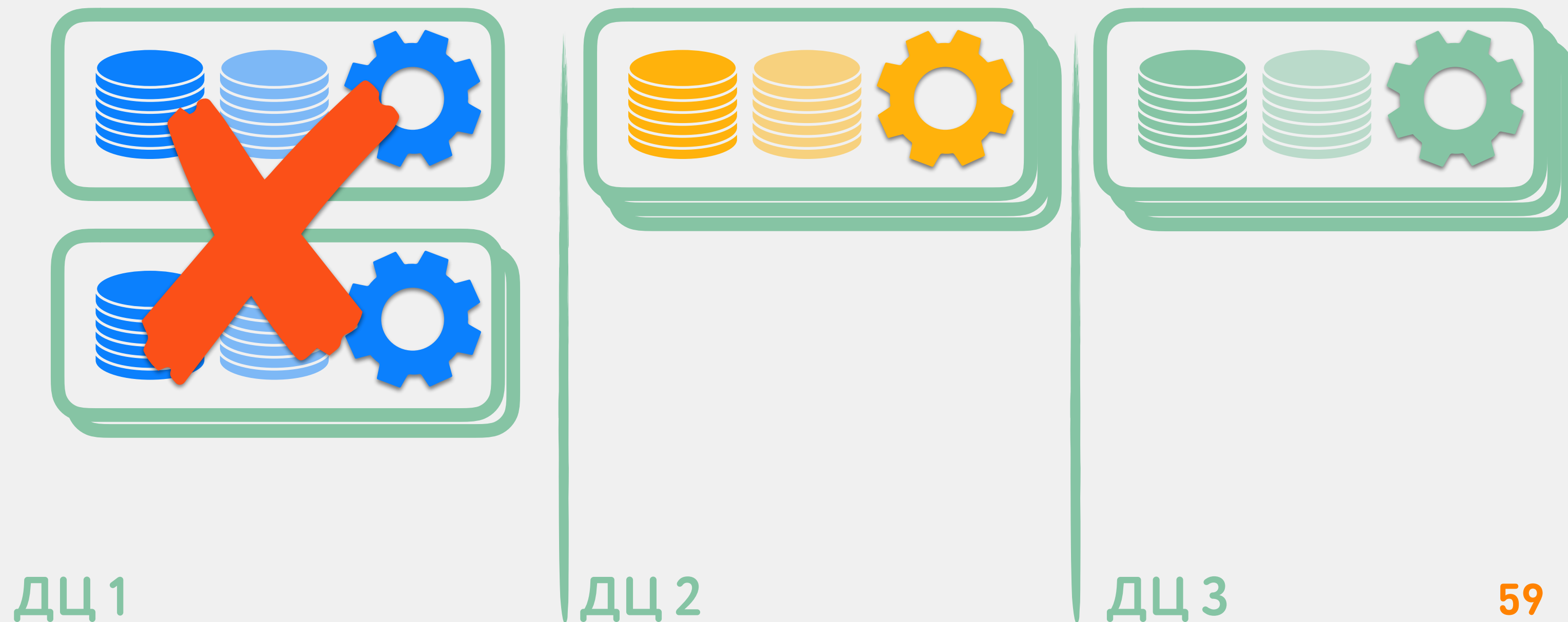
Взаимное влияние прикладухи на БД

КТО ВИНОВАТ:

- **БД теперь вместе с приложением**
Прикладные баги в потреблении CPU, памяти могут уложить все реплики сервиса

что делать:

- **Продленная тестовая эксплуатация новой версии**
В одной зоне доступности (остальные остаются на старой)
- **Рубильники на опасный функционал**



Апгрейд на новую версию встроенной БД

КТО ВИНОВАТ:

- **Внутреннее API БД**

Которое может измениться от версии к версии БД



ОНИ все поменяют и мы застрянем на старой версии БД - и потом мы же все п*ем¹ !!11

что делать:

- **Архитектура БД меняется медленно**

9 лет наблюдения за Cassandra выявило только страсть к переименованию.

- **Современные БД поддерживают rolling upgrades**

Делаем версию сервиса с новой библиотекой БД,

Выкатываем на 1 ДЦ,

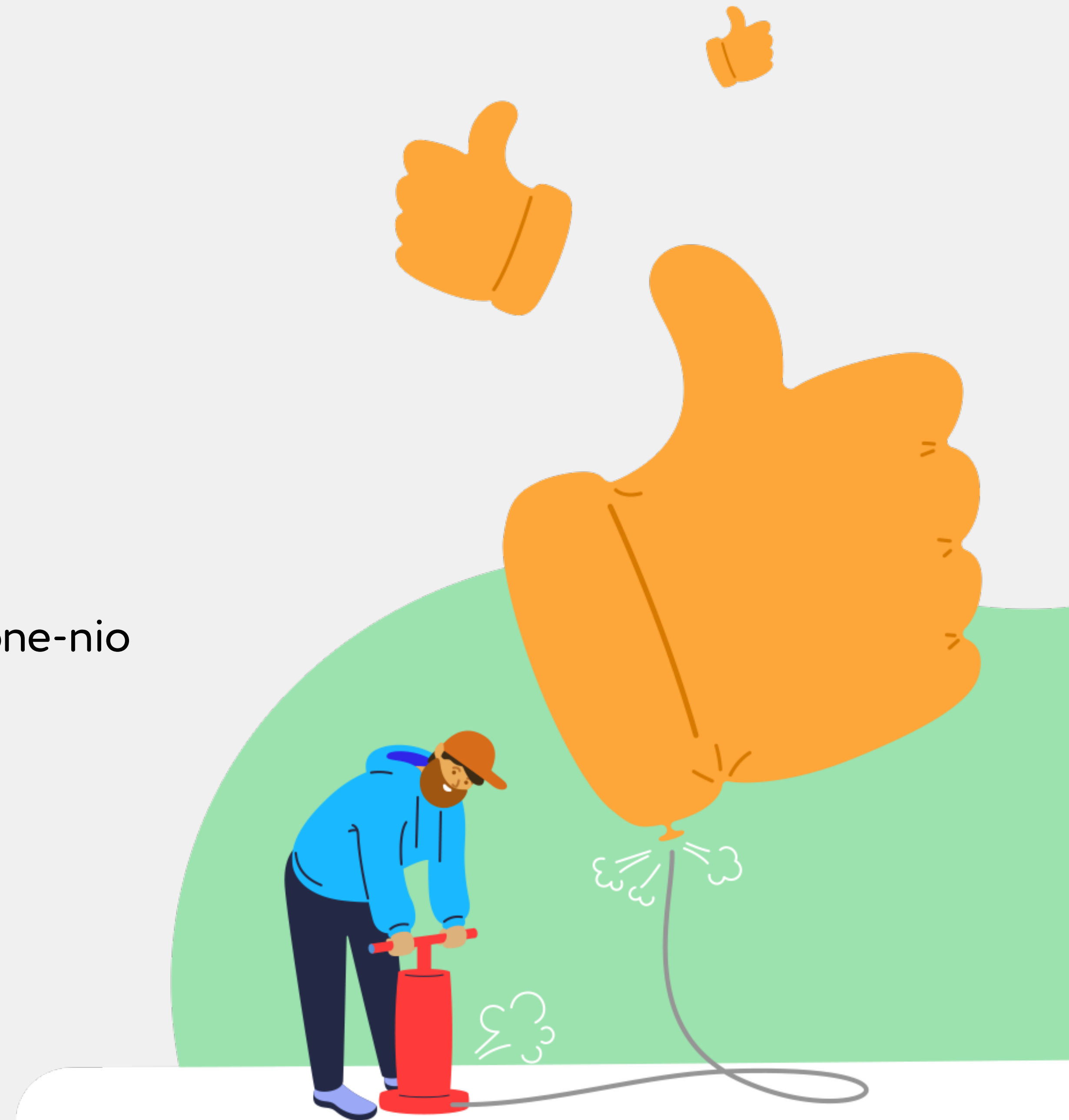
Ждем что сломается,

В тяжелых случаях - откат и восстанавливаем данные из реплик

(1) Потеряем все данные, клиентов, работу и не сможем их потом восстановить

Итог

- **Эффективнее и надежнее**
за счет отсутствия потерь на маршалинг и сеть
- **Гарантии консистентности в кешах**
За счет совмещения кеша и БД в 1 процессе
- **Относительно просто реализуется**
Все самое сложное уже реализовано в Cassandra и one-nio
- **Мы широко и давно используем**
лента, обсуждения, посты, нотификации, ...



Тут можно узнать больше:



Распределенные системы в Одноклассниках

Олег Анастасьев

ok.ru/video/97105087088



Новый граф Одноклассников

Антон Иванов

highload.ru/moscow/2019/abstracts/5793



Восстание машин - это ОК

Леонид Талалаев

highload.ru/moscow/2019/abstracts/5784





ОДНОКЛАССНИКИ