



# Apache Cassandra: 8 лет в продакшене

**Владимир Дегтерёв**

degtereivy@pochtabank.ru

Telegram: @VladimirDegtereov

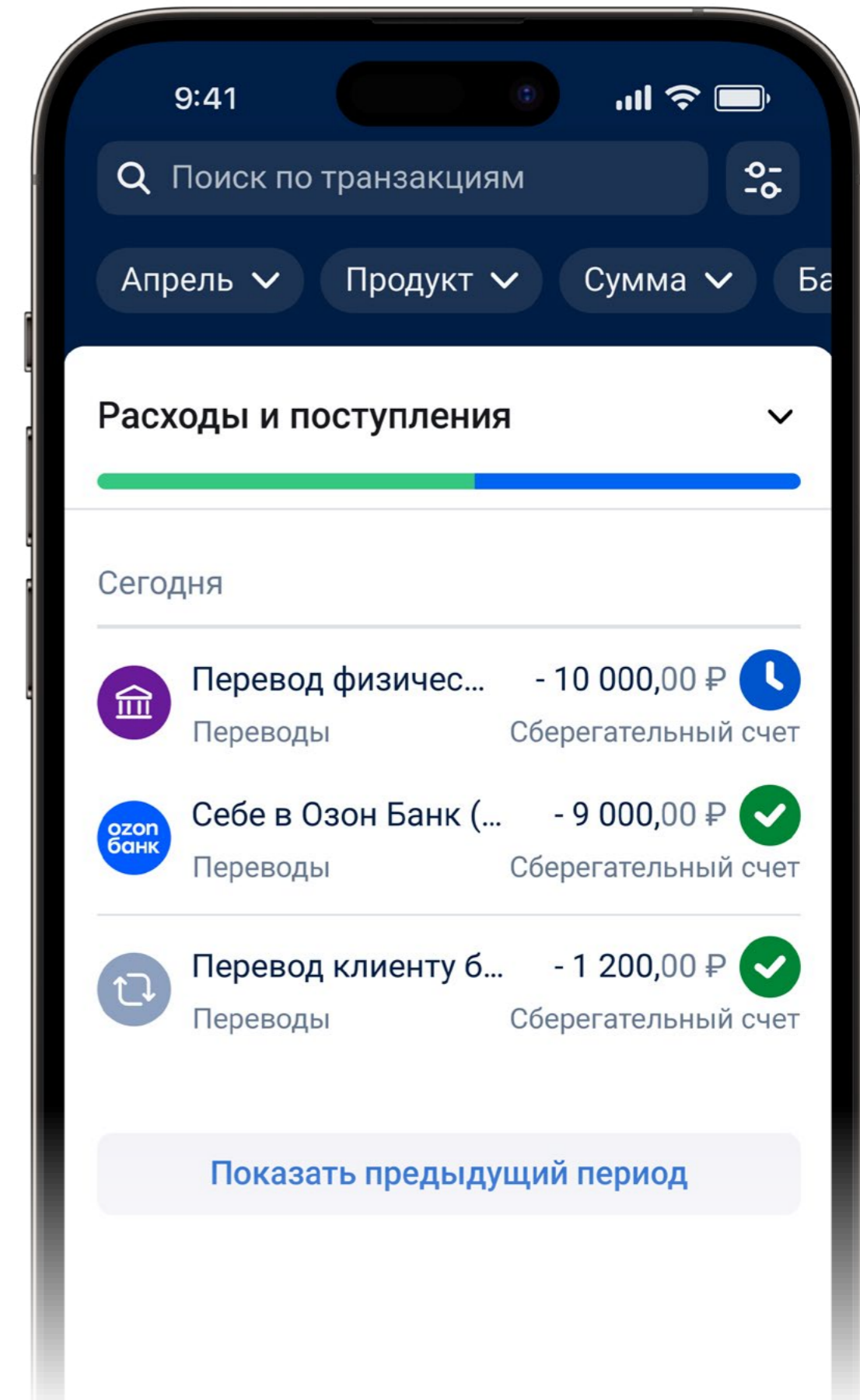
# | О чём будем говорить



- Почему выбрали Cassandra
- Как мы писали сервис ленты операций
- Как боролись с падением скорости чтения

# Постановка задачи

## Разработать интернет и мобильный банк



# Постановка задачи



## Разработать интернет и мобильный банк

Требования к БД:

# Постановка задачи



## Разработать интернет и мобильный банк

Требования к БД:

- Отказоустойчивость



**Кластер — это хорошо,  
Кластер — это надёжно**

# Постановка задачи

## Разработать интернет и мобильный банк

Требования к БД:

- Отказоустойчивость
- Масштабируемость



# Постановка задачи



## Разработать интернет и мобильный банк

Требования к БД:

- Отказоустойчивость
- Масштабируемость
- Производительность

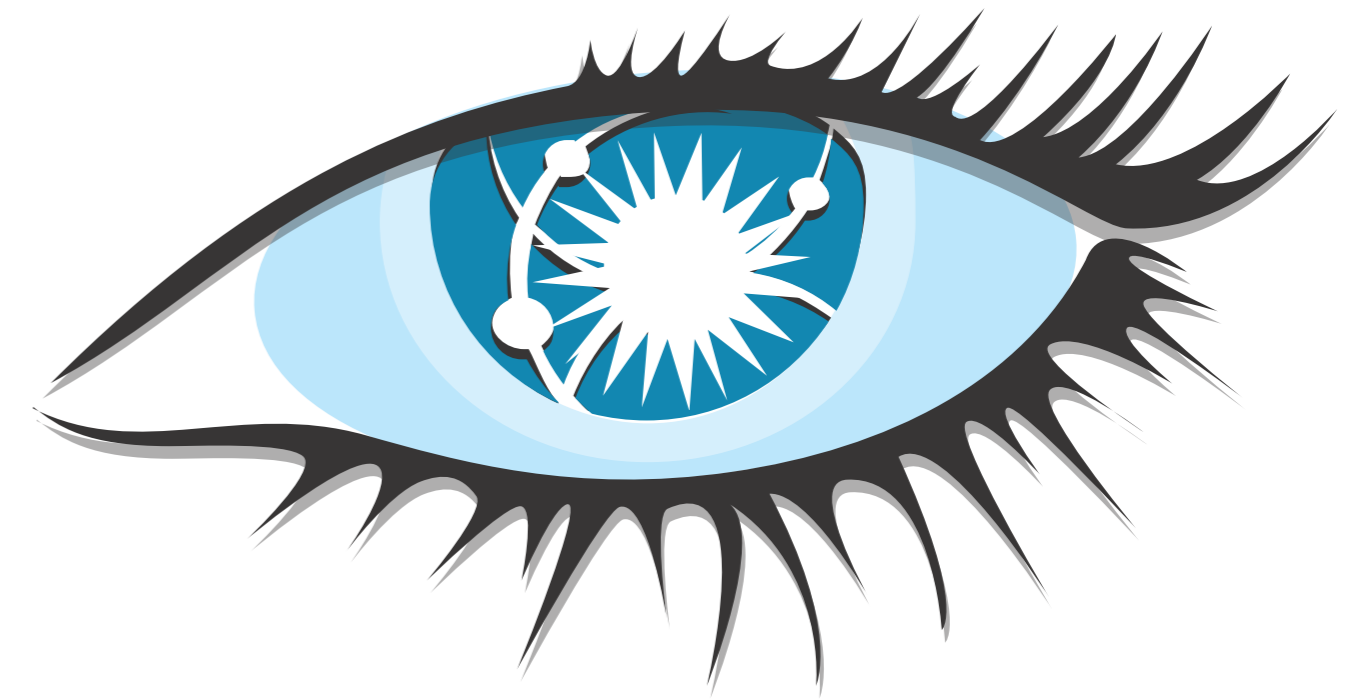


**Скорость без границ**

# Почему именно Cassandra?



- **Настраиваемая согласованность**

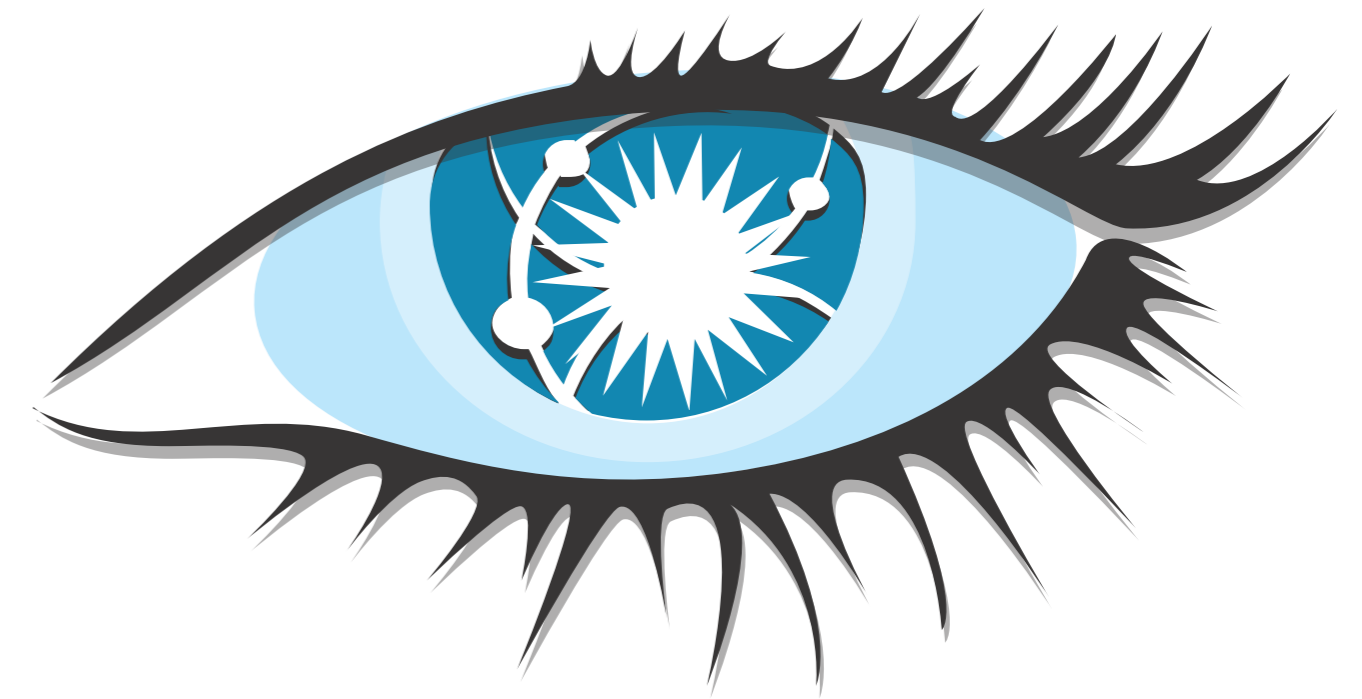




# Почему именно Cassandra?



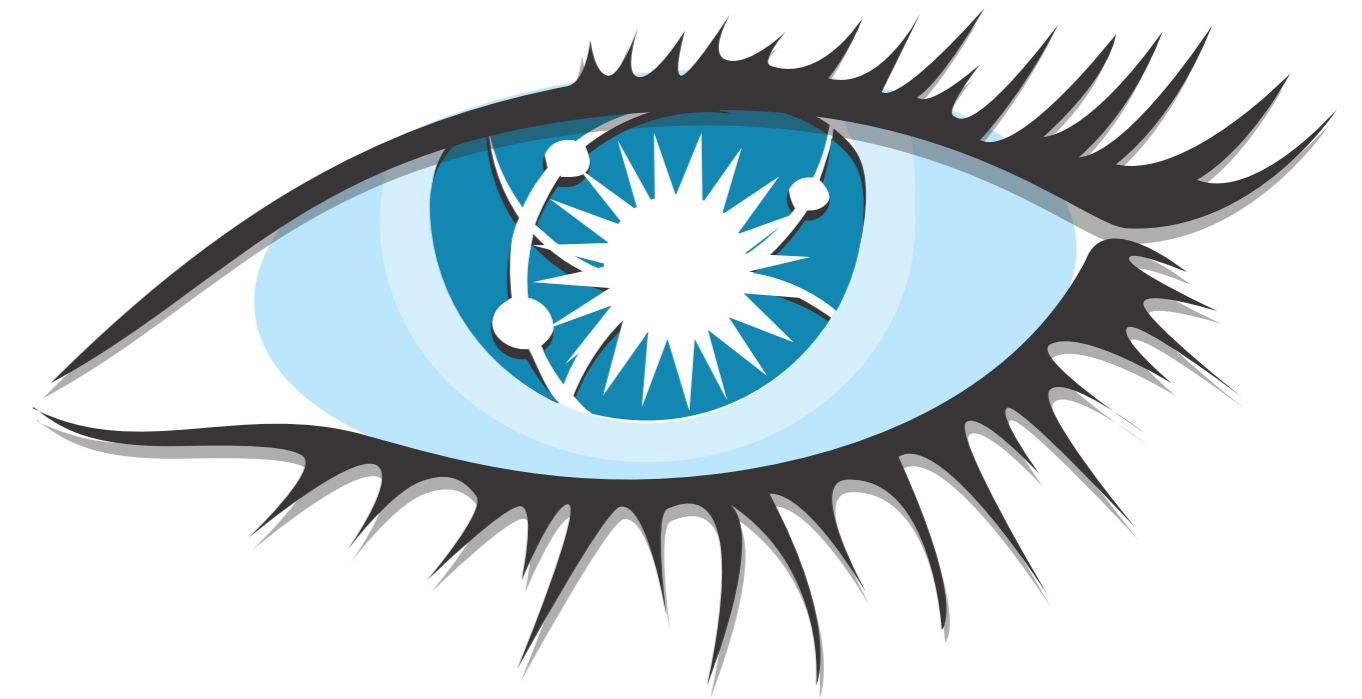
- **Настраиваемая согласованность**
- **Модель чтения/записи**



# Почему именно Cassandra?



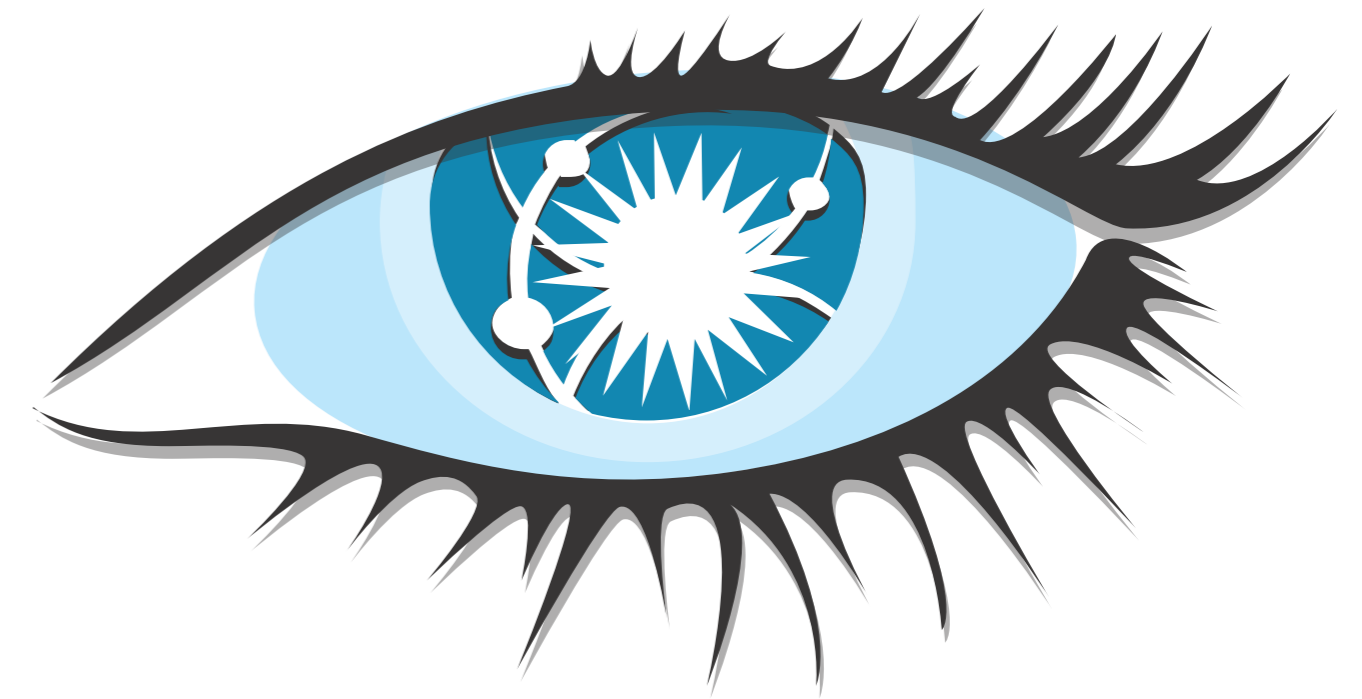
- **Настраиваемая согласованность**
- **Модель чтения/записи**
- **Partitioned wide column storage model**



# Почему именно Cassandra?



- **Настраиваемая согласованность**
- **Модель чтения/записи**
- **Partitioned wide column storage model**



Partition key: 5

Clustering key: a

x:12

y:56

z:5

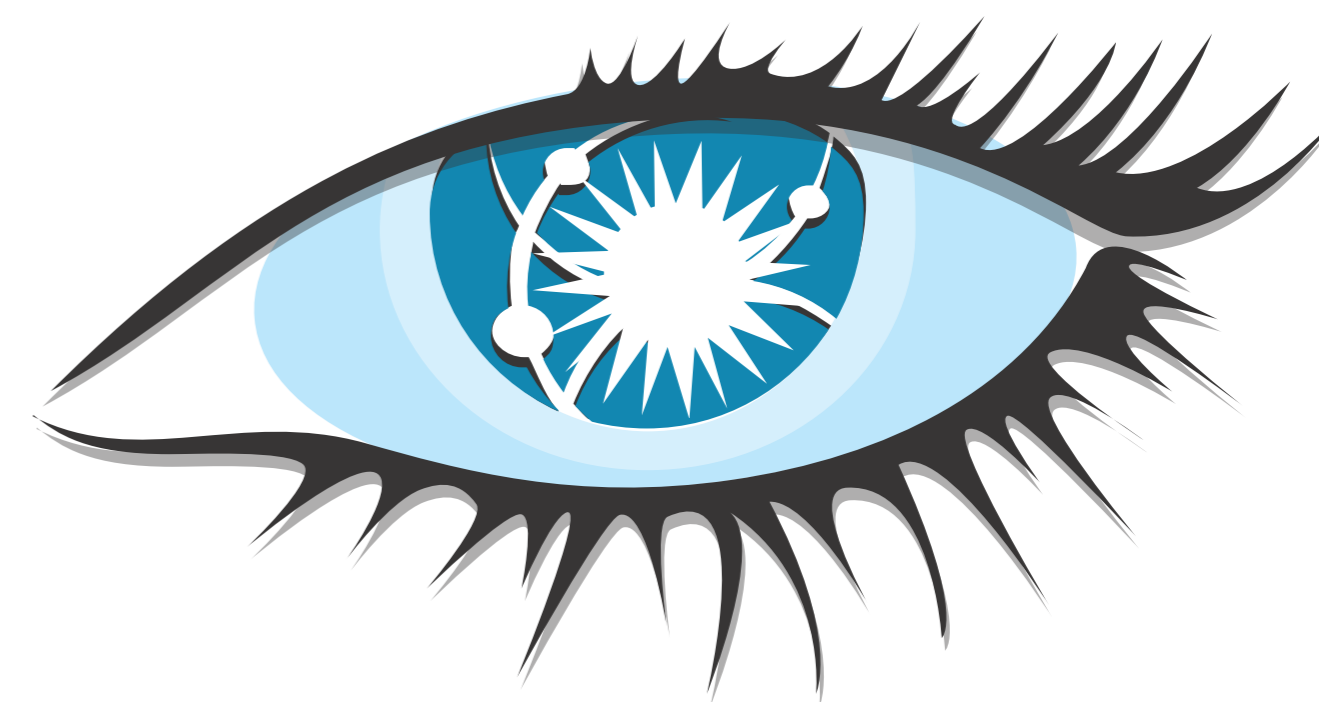
Clustering key: b

j:22

k:1

# Почему именно Cassandra?

- **Настраиваемая согласованность**
- **Модель чтения/записи**
- **Partitioned wide column storage model**



Partition key: 5

Clustering key: a

x:12

y:56

z:5

Clustering key: b

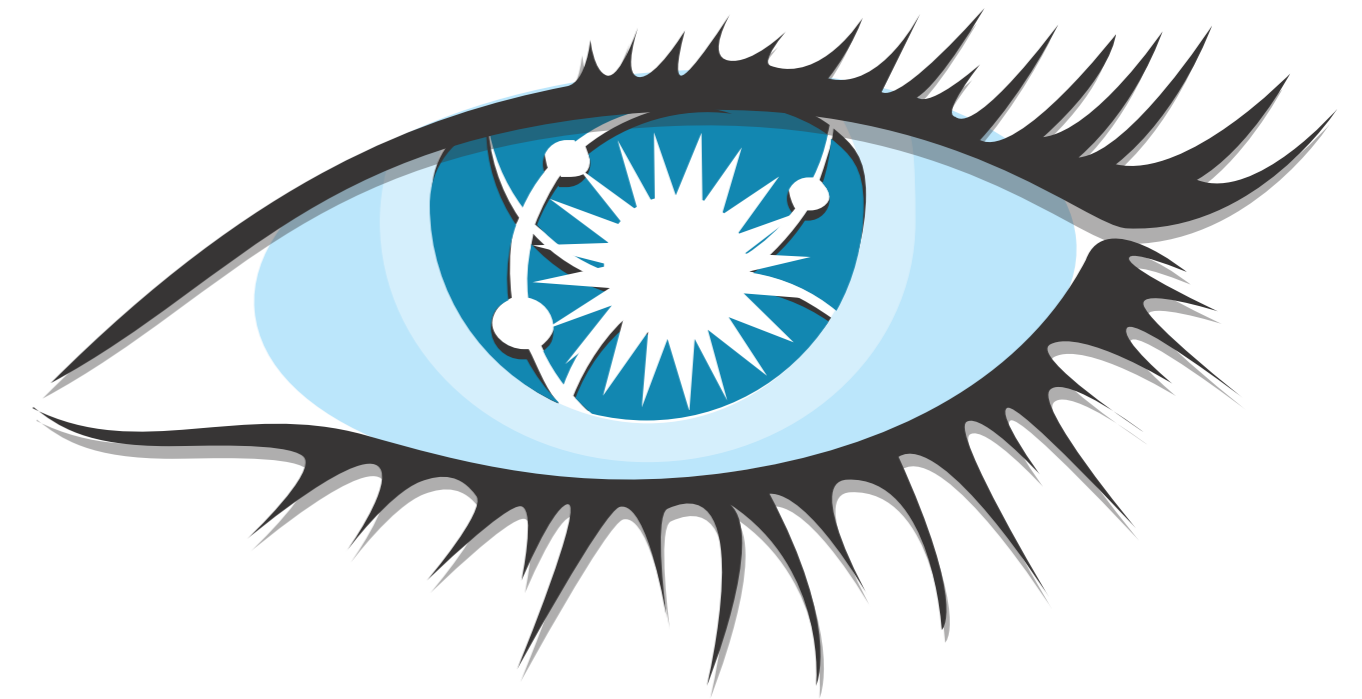
j:22

k:1

# Почему именно Cassandra?



- **Настраиваемая согласованность**
- **Модель чтения/записи**
- **Partitioned wide column storage model**



Partition key: 5

Clustering key: a

x:12

y:56

z:5

Clustering key: b

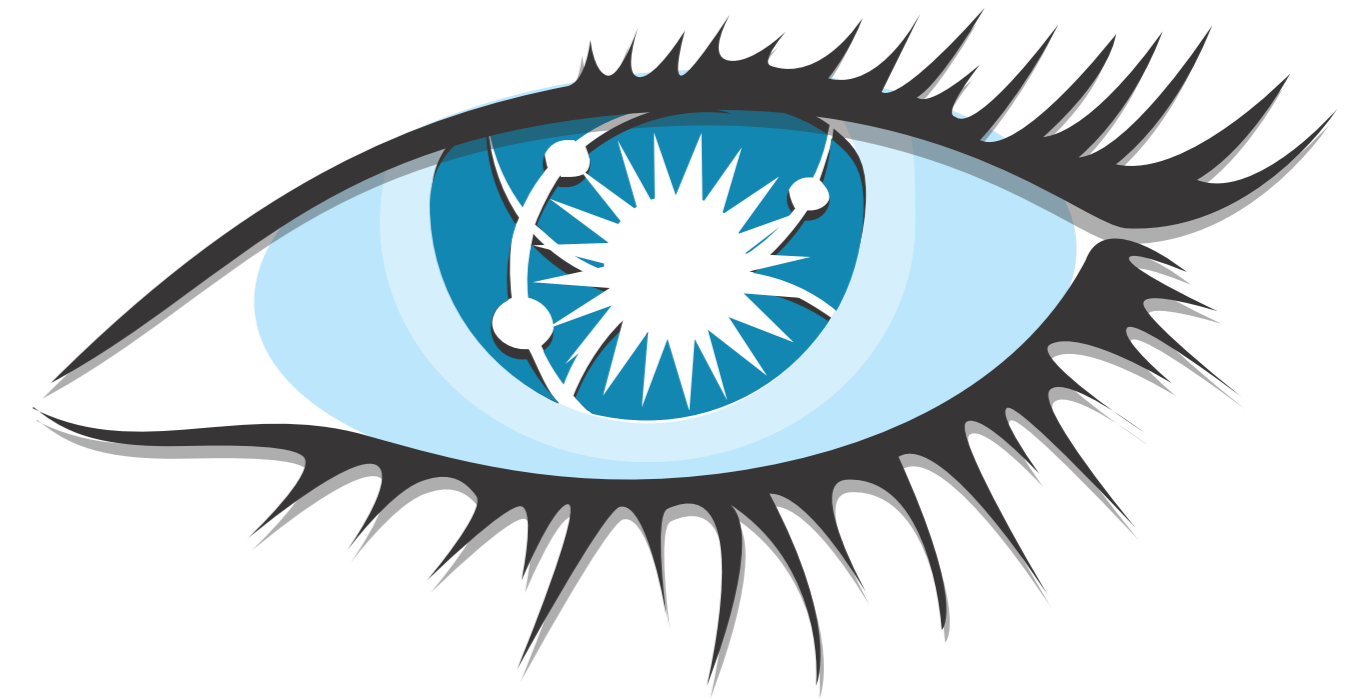
j:22

k:1

# Почему именно Cassandra?



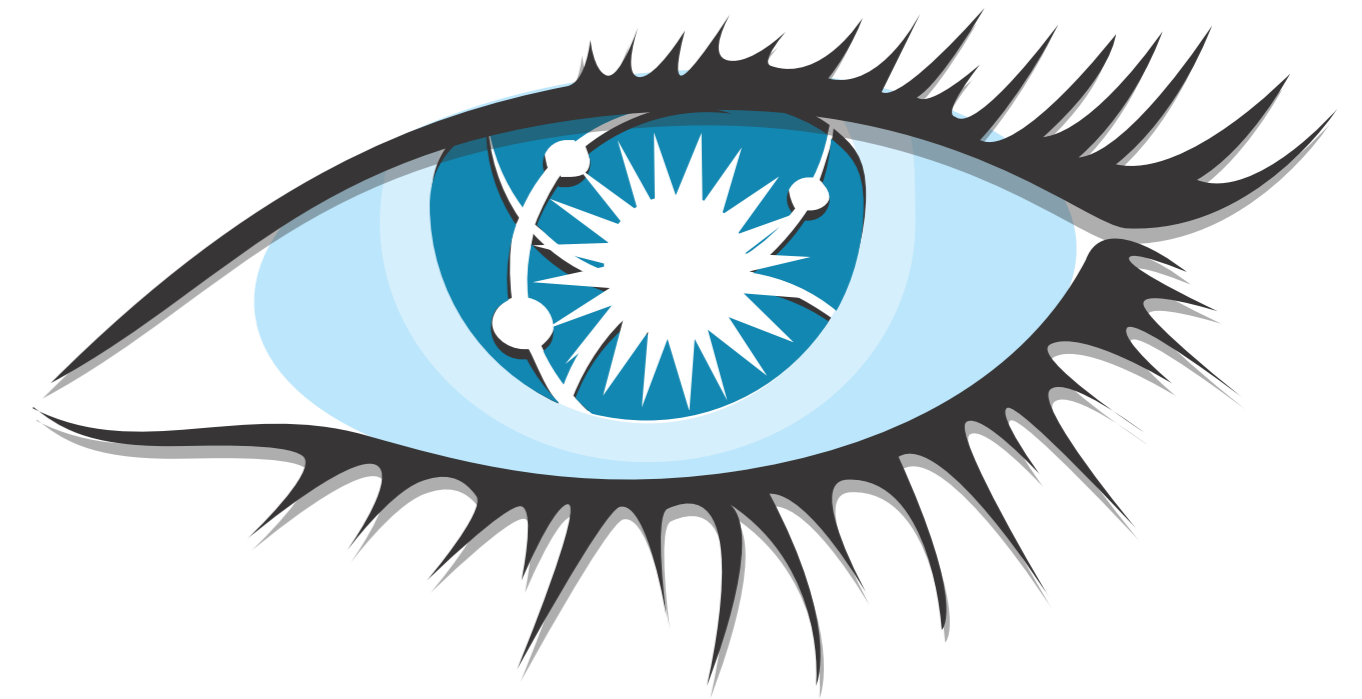
- **Настраиваемая согласованность**
- **Модель чтения/записи**
- **Partitioned wide column storage model**



# Почему именно Cassandra?



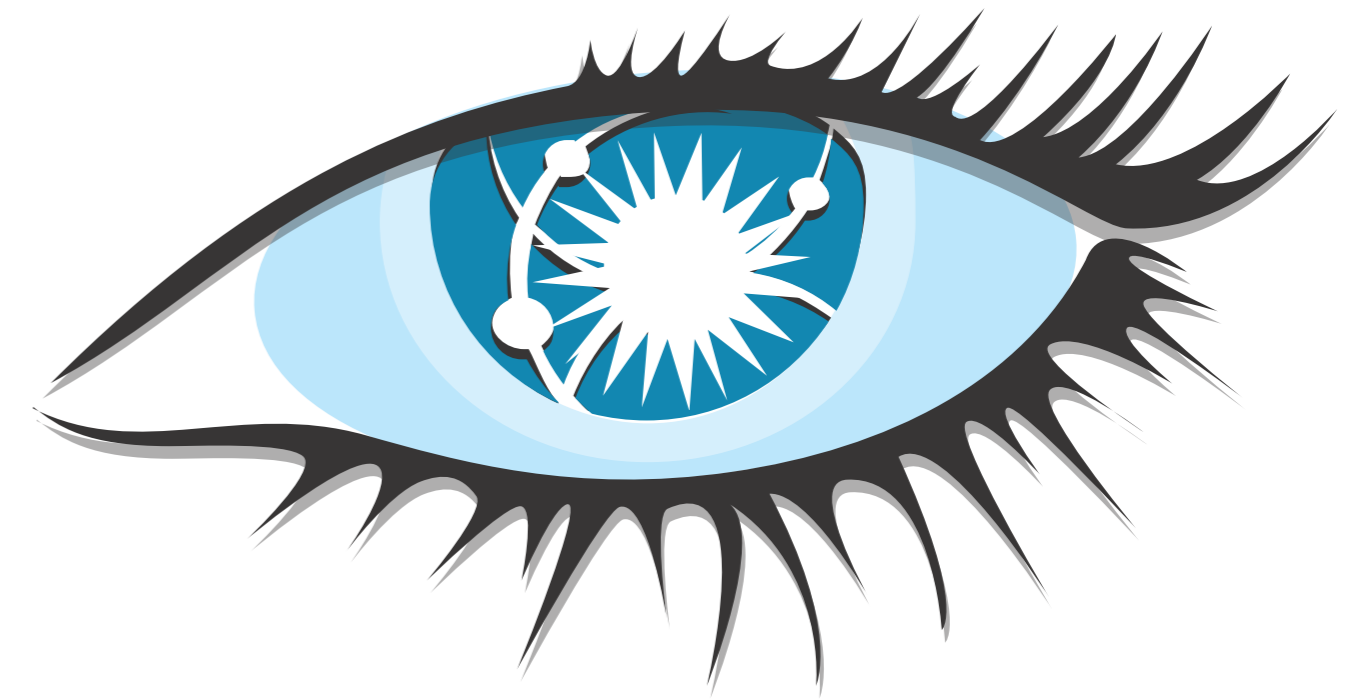
- **Настраиваемая согласованность**
- **Модель чтения/записи**
- **Partitioned wide column storage model**



# Почему именно Cassandra?



- **Настраиваемая согласованность**
- **Модель чтения/записи**
- **Partitioned wide column storage model**
- **Open source**

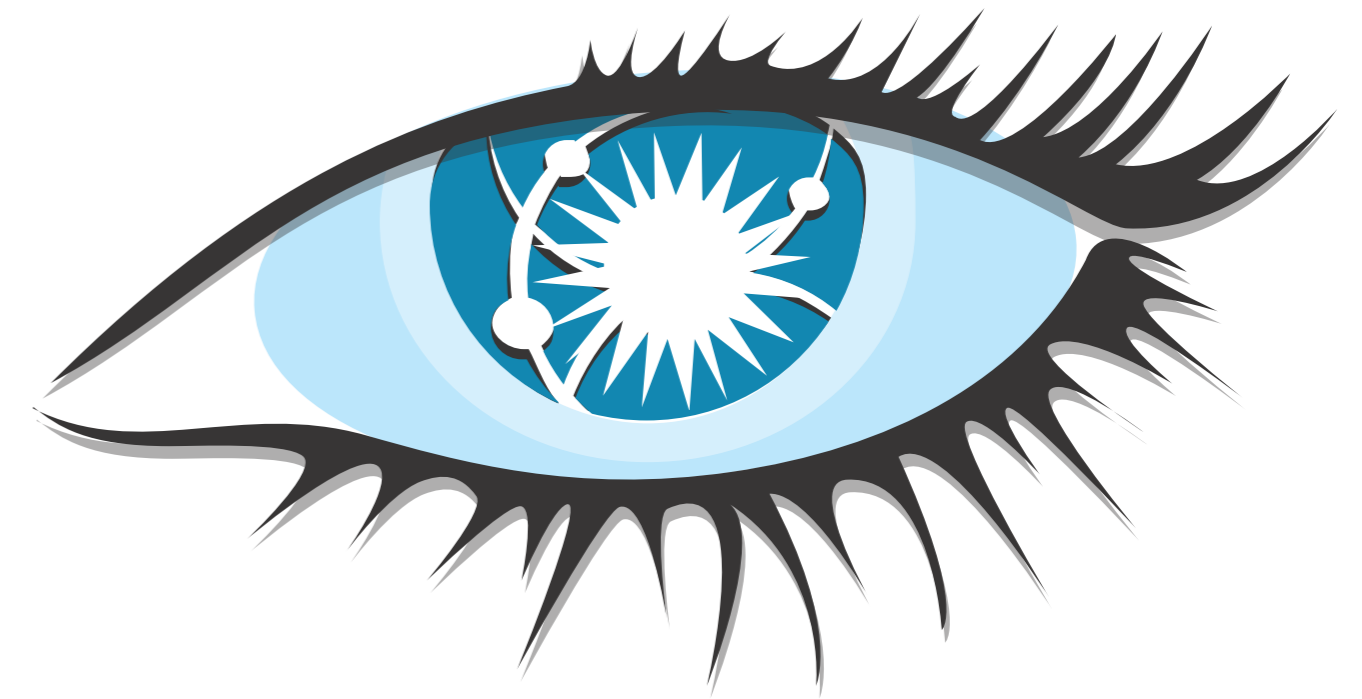




# Почему именно Cassandra?



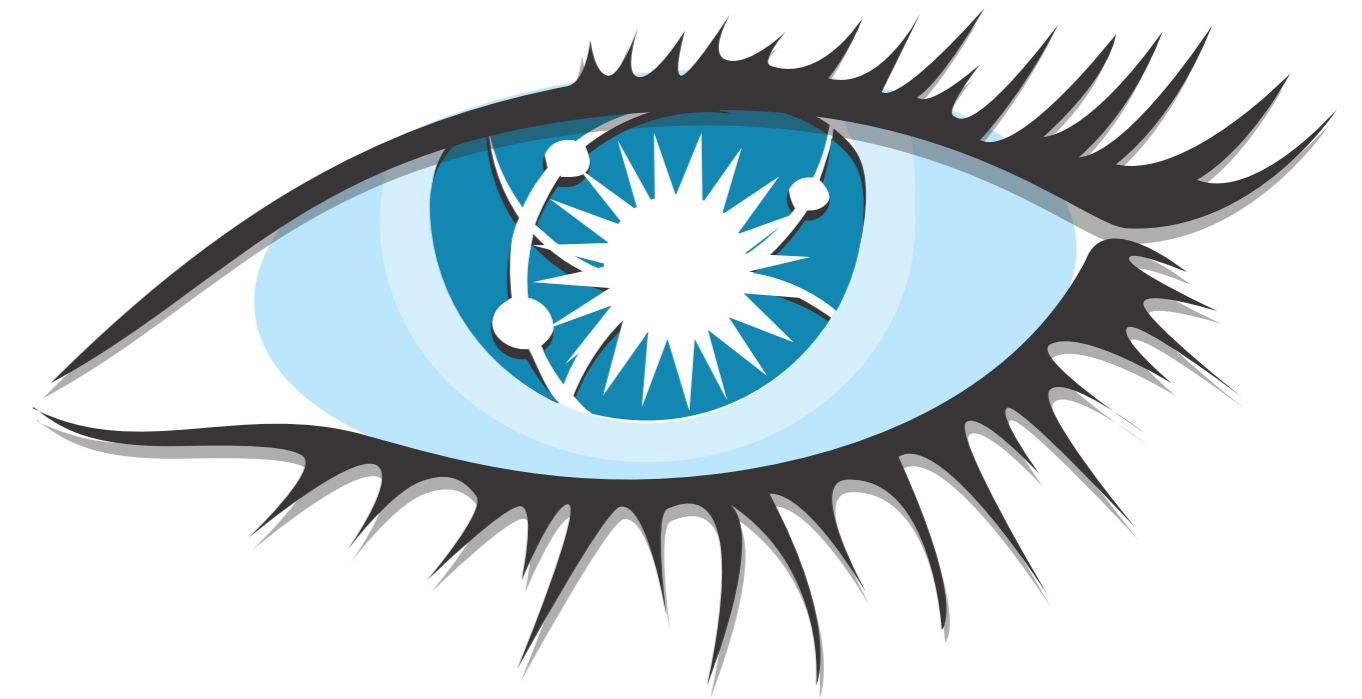
- **Настраиваемая согласованность**
- **Модель чтения/записи**
- **Partitioned wide column storage model**
- **Open source**
- **Java**



# Почему именно Cassandra?



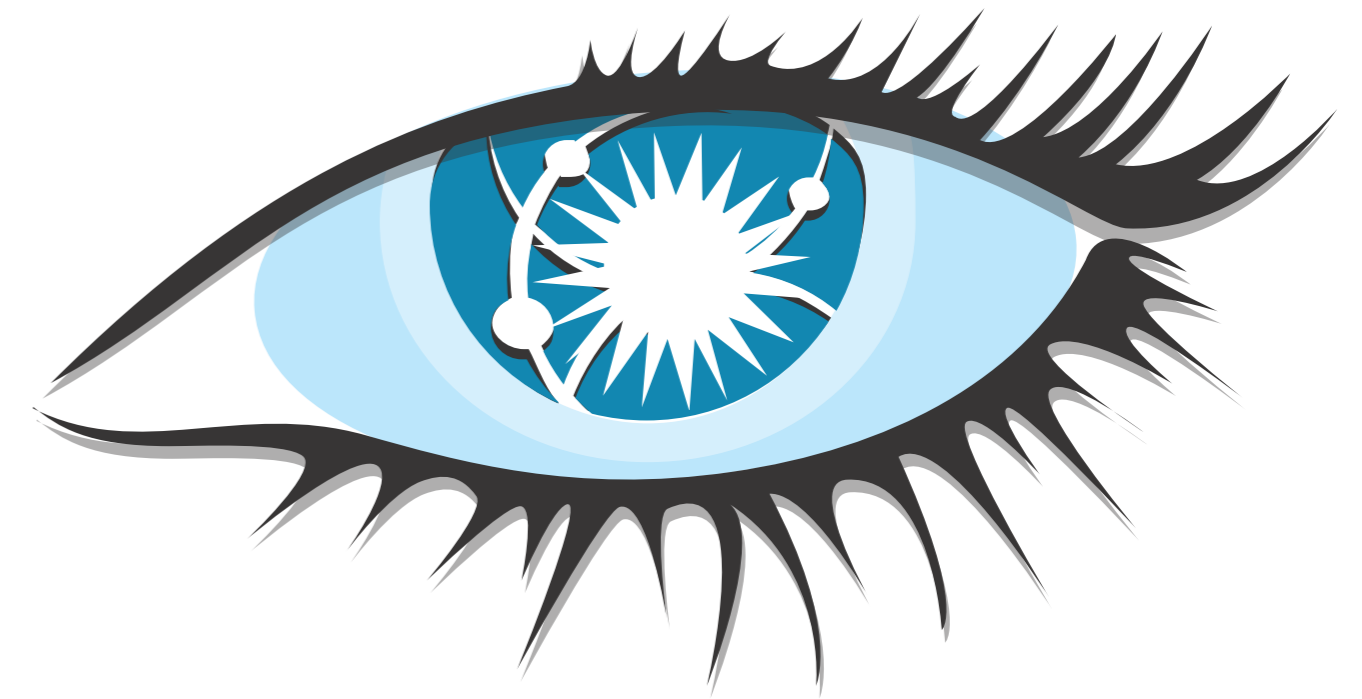
- **Настраиваемая согласованность**
- **Модель чтения/записи**
- **Partitioned wide column storage model**
- **Open source**
- **Java**
- **Зрелость проекта, сообщество**



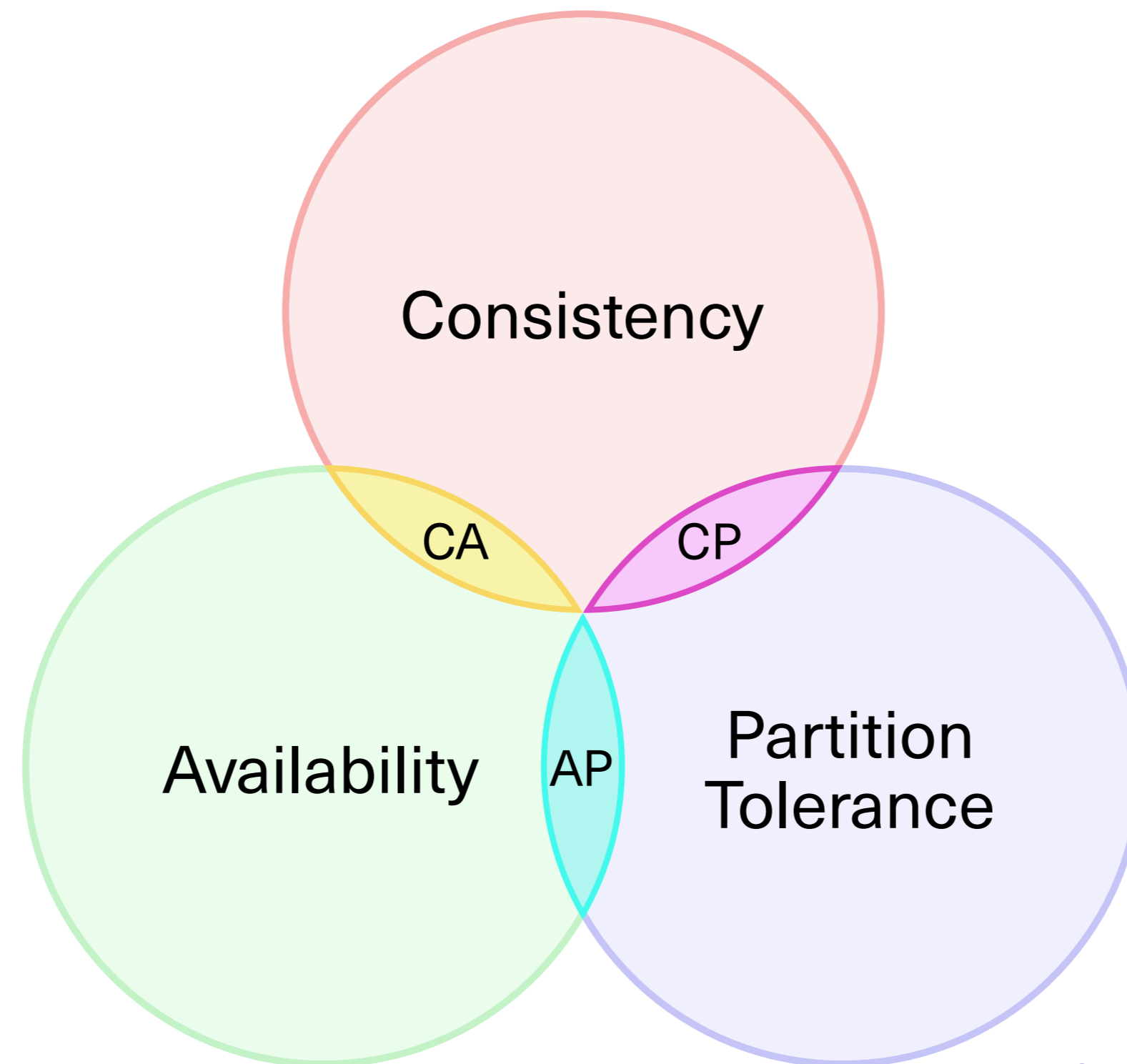
# Почему именно Cassandra?



- **Настраиваемая согласованность**
- **Модель чтения/записи**
- **Partitioned wide column storage model**
- **Open source**
- **Java**
- **Зрелость проекта, сообщество**
- **ScyllaDB уже вышла, но побоялись**



## CAP теорема



[https://en.wikipedia.org/wiki/CAP\\_theorem](https://en.wikipedia.org/wiki/CAP_theorem)

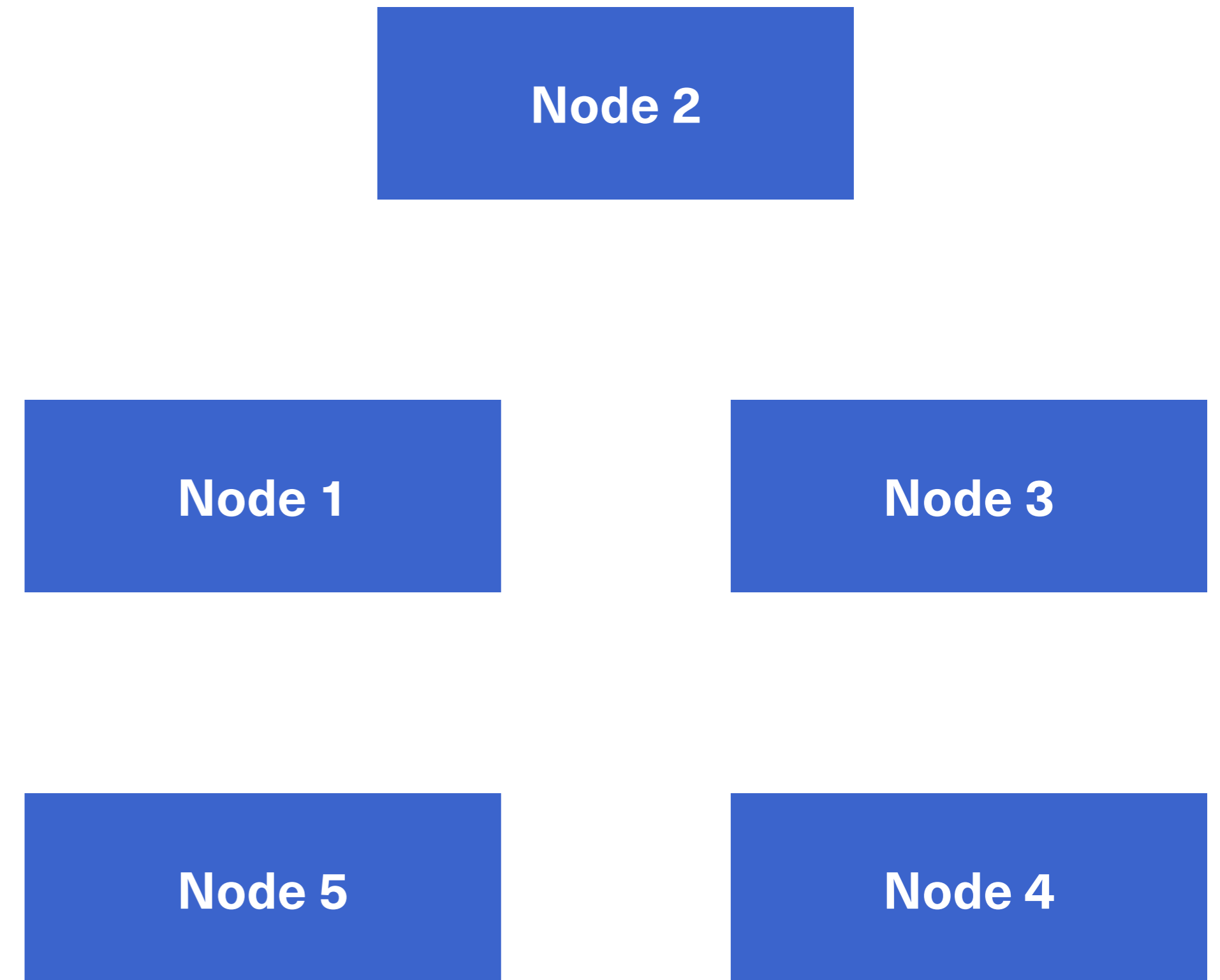
[https://en.wikipedia.org/wiki/PACELC\\_theorem](https://en.wikipedia.org/wiki/PACELC_theorem)

# Настраиваемая согласованность



**Replication factor (RF)** —  
сколько узлов хранят  
копии одной партиции

**Consistency level (CL)** —  
сколько из этих узлов  
должны подтвердить  
операцию, чтобы вернуть  
клиенту «успех»



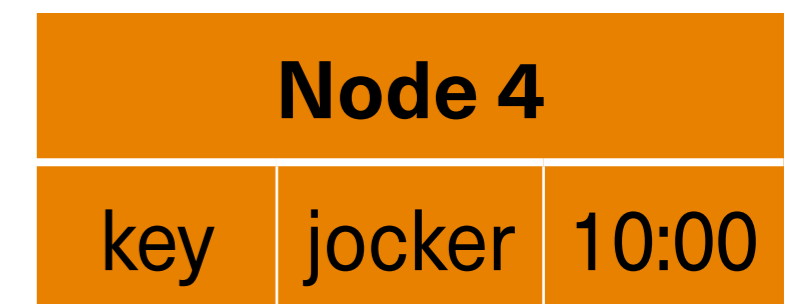
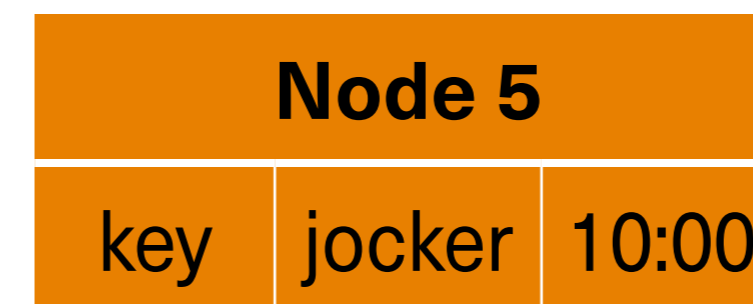
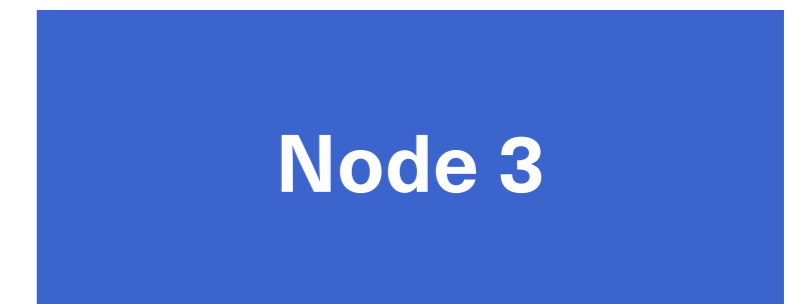
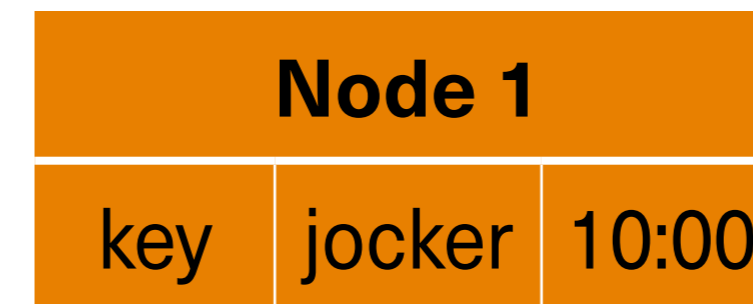
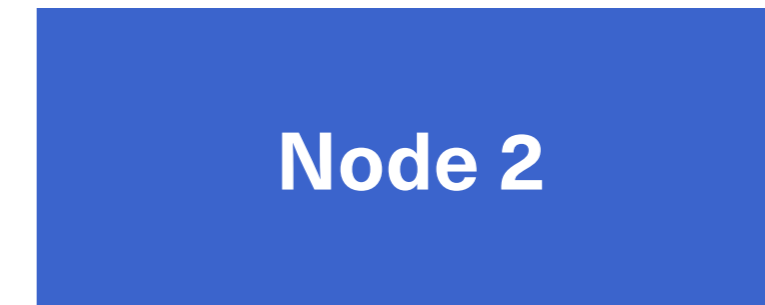
# Настраиваемая согласованность



**Replication factor (RF)** — сколько узлов хранят копии одной партиции

**Consistency level (CL)** — сколько из этих узлов должны подтвердить операцию, чтобы вернуть клиенту «успех»

**RF=3**



# Настраиваемая согласованность

**Replication factor (RF)** —  
сколько узлов хранят  
копии одной партиции

**Consistency level (CL)** —  
сколько из этих узлов  
должны подтвердить  
операцию, чтобы вернуть  
клиенту «успех»

**RF=3**

**Node 2**

**UPDATE**  
CL = QUORUM  

key	jpoint
-----	--------

**Node 1**

key	jocker	10:00
-----	--------	-------

**Node 3**

**Node 5**

key	jocker	10:00
-----	--------	-------

**Node 4**

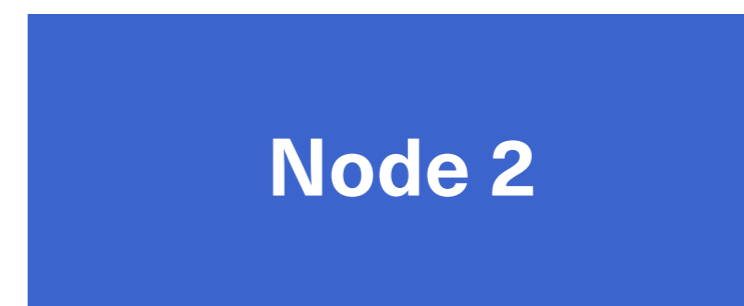
key	jocker	10:00
-----	--------	-------

# Настраиваемая согласованность

**Replication factor (RF)** —  
сколько узлов хранят  
копии одной партиции

**Consistency level (CL)** —  
сколько из этих узлов  
должны подтвердить  
операцию, чтобы вернуть  
клиенту «успех»

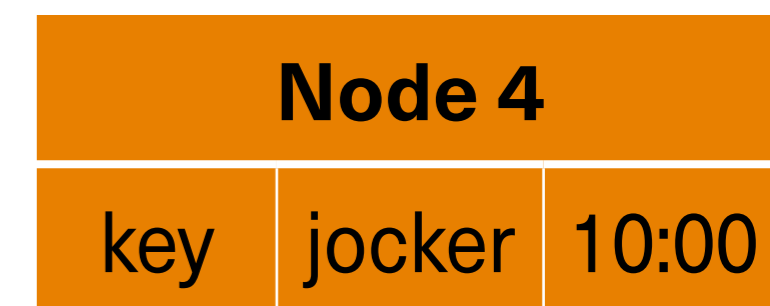
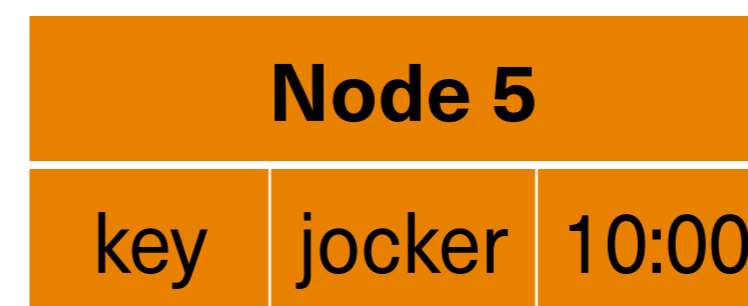
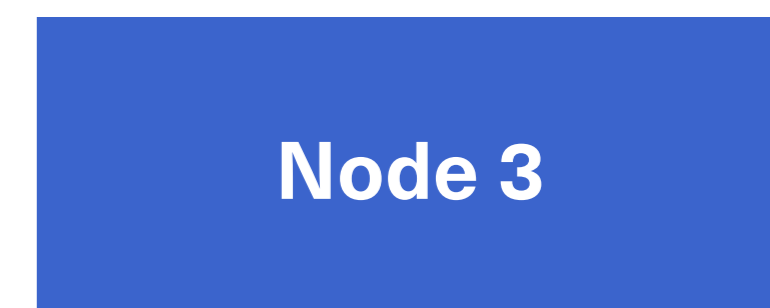
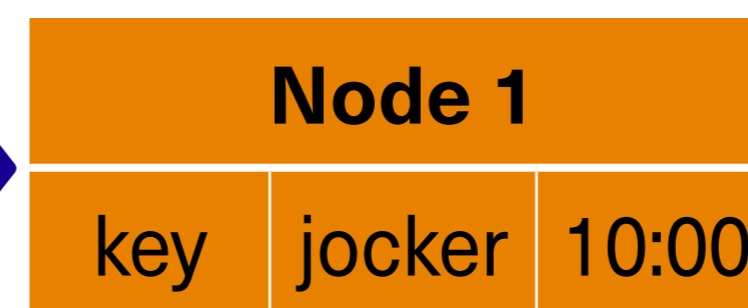
**RF=3**



**UPDATE**

CL = QUORUM →

key	jpoint
-----	--------

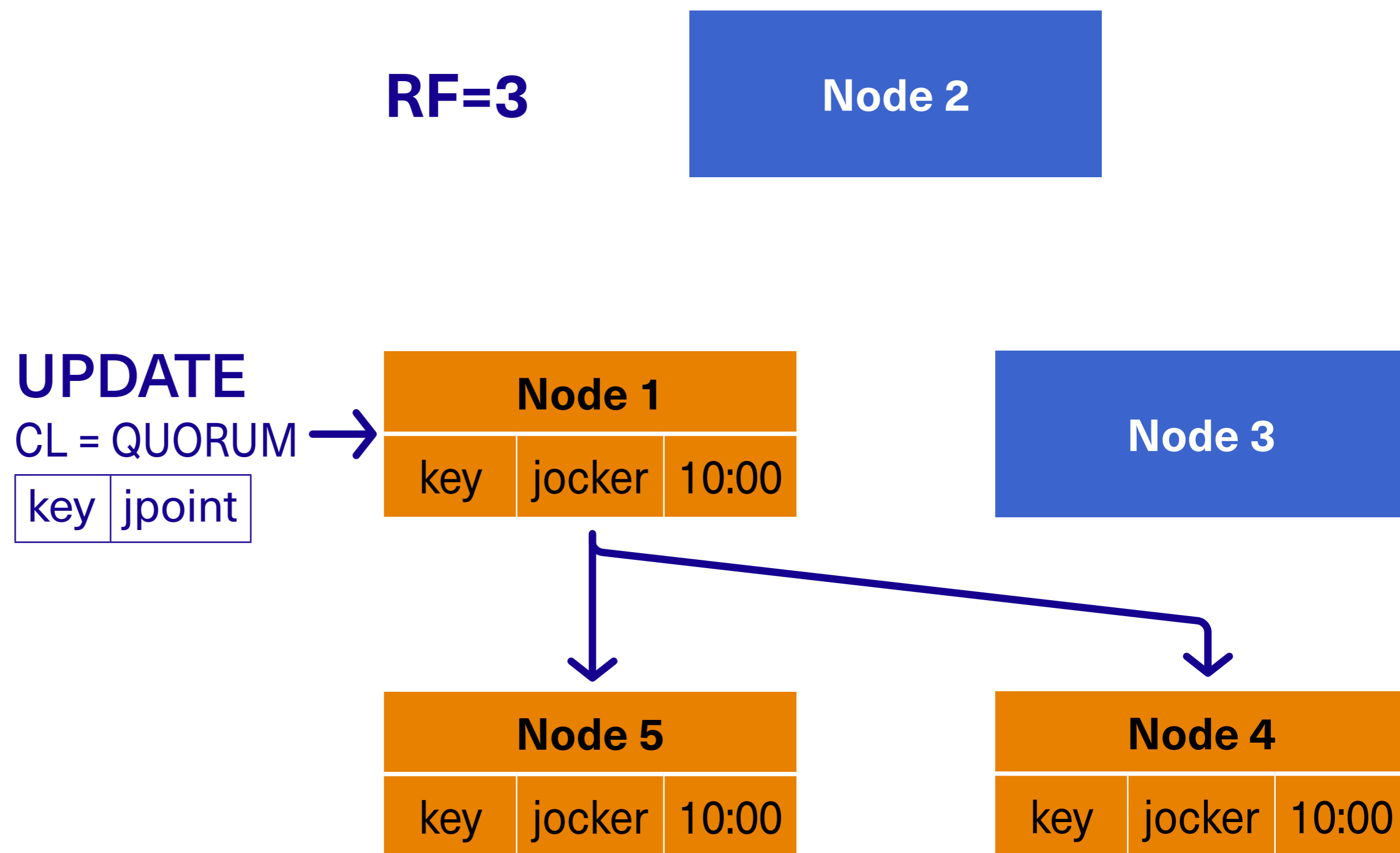




# Настраиваемая согласованность

**Replication factor (RF)** — сколько узлов хранят копии одной партии

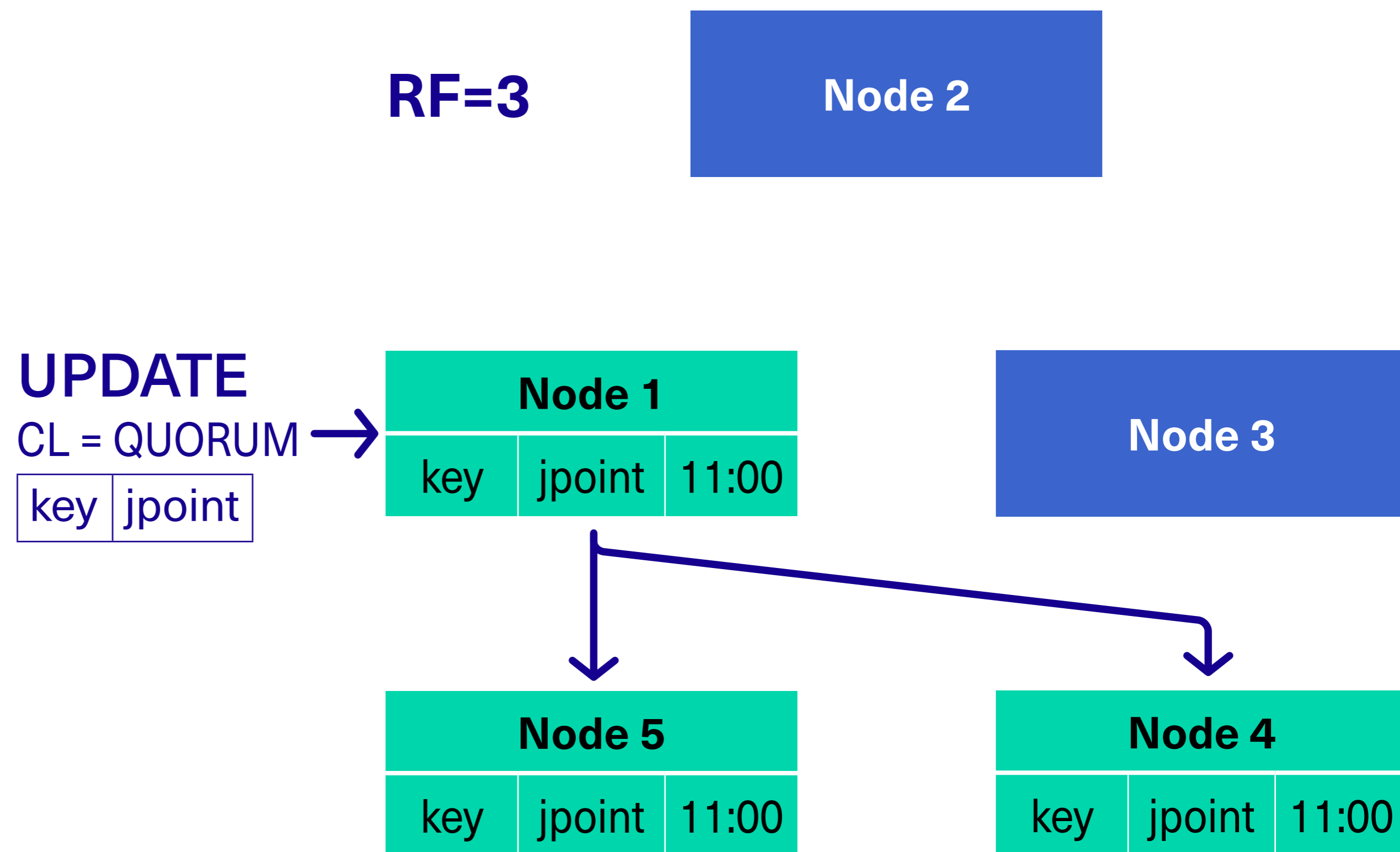
**Consistency level (CL)** — сколько из этих узлов должны подтвердить операцию, чтобы вернуть клиенту «успех»



# Настраиваемая согласованность

**Replication factor (RF)** — сколько узлов хранят копии одной партии

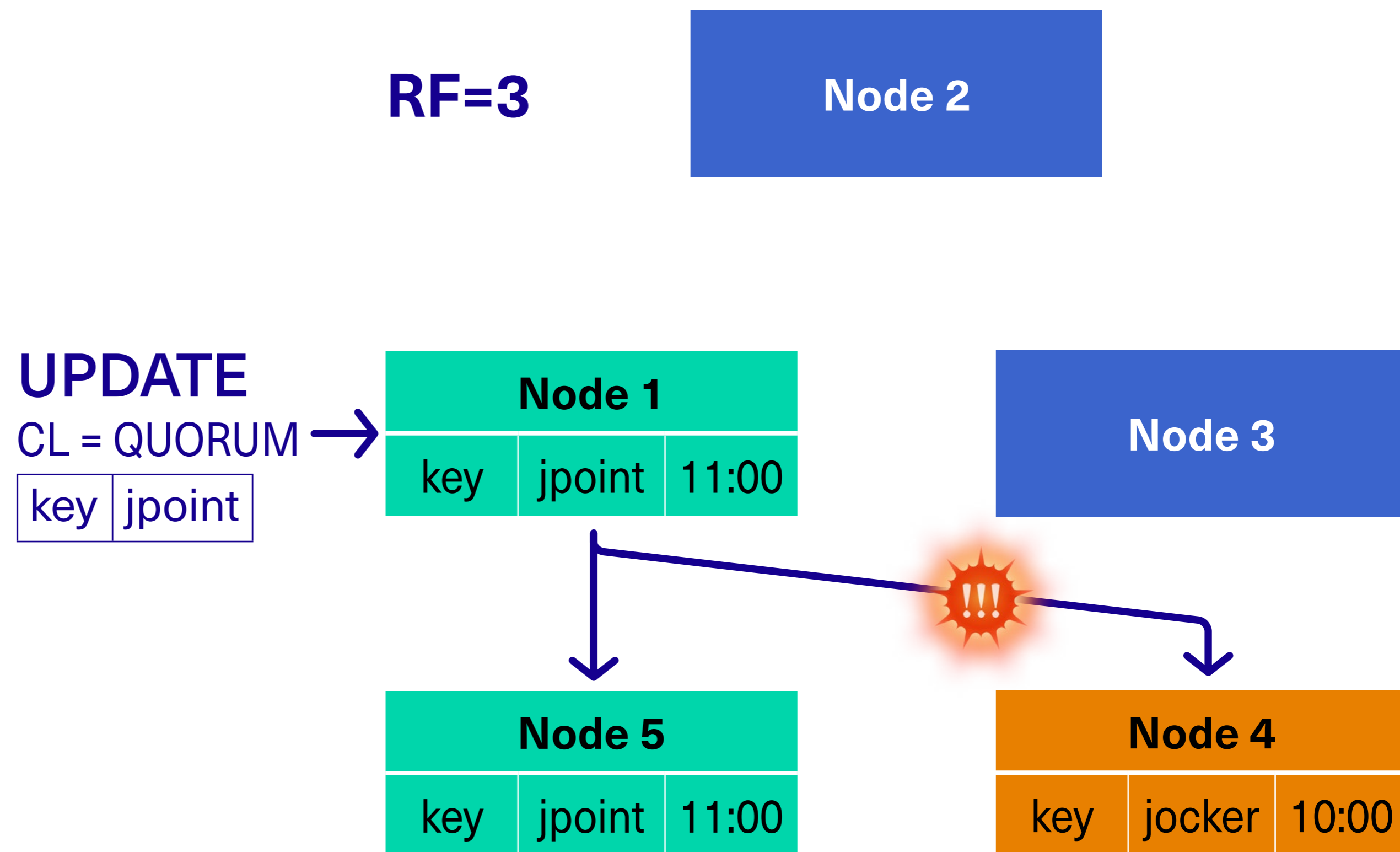
**Consistency level (CL)** — сколько из этих узлов должны подтвердить операцию, чтобы вернуть клиенту «успех»



# Настраиваемая согласованность

**Replication factor (RF)** — сколько узлов хранят копии одной партии

**Consistency level (CL)** — сколько из этих узлов должны подтвердить операцию, чтобы вернуть клиенту «успех»

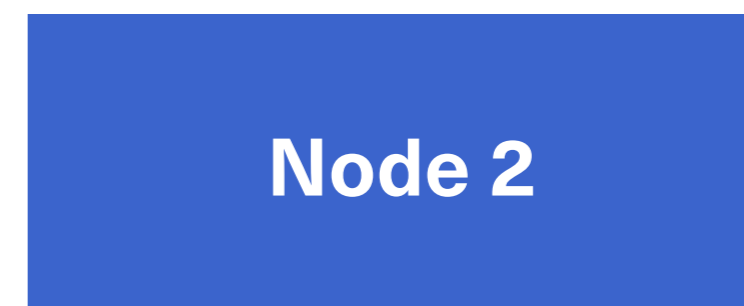


# Настраиваемая согласованность

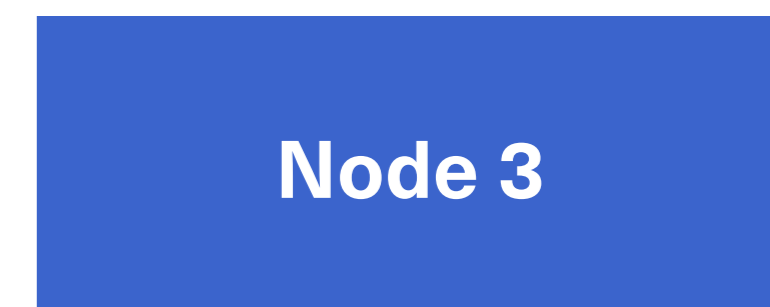
**Replication factor (RF)** —  
сколько узлов хранят  
копии одной партиции

**Consistency level (CL)** —  
сколько из этих узлов  
должны подтвердить  
операцию, чтобы вернуть  
клиенту «успех»

**RF=3**



Node 1		
key	jpoint	11:00



Node 5		
key	jpoint	11:00

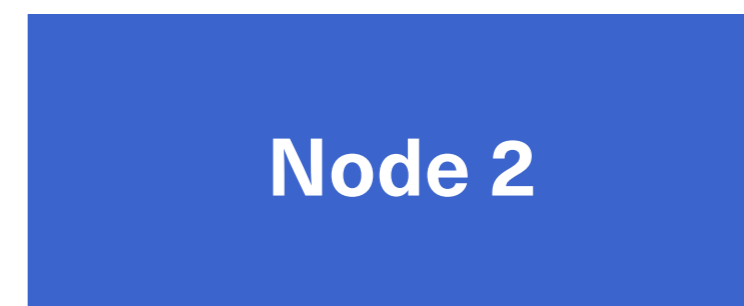
Node 4		
key	jocker	10:00

# Настраиваемая согласованность

**Replication factor (RF)** — сколько узлов хранят копии одной партии

**Consistency level (CL)** — сколько из этих узлов должны подтвердить операцию, чтобы вернуть клиенту «успех»

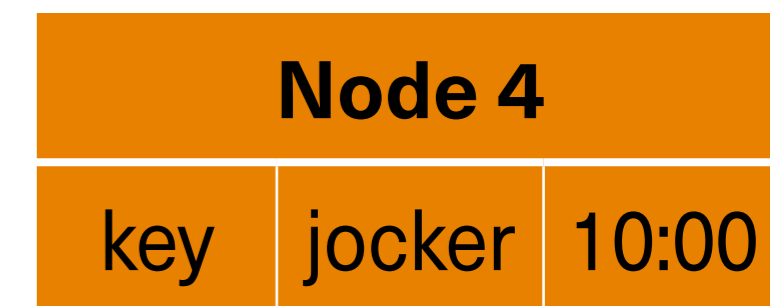
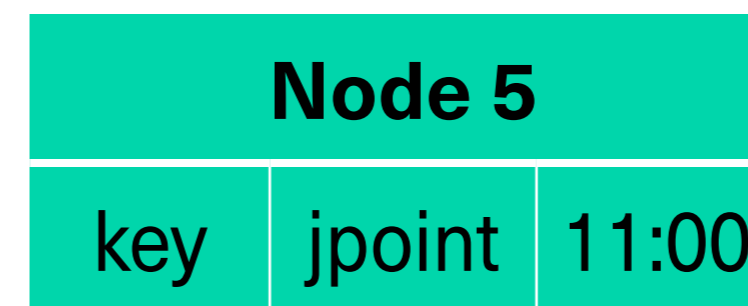
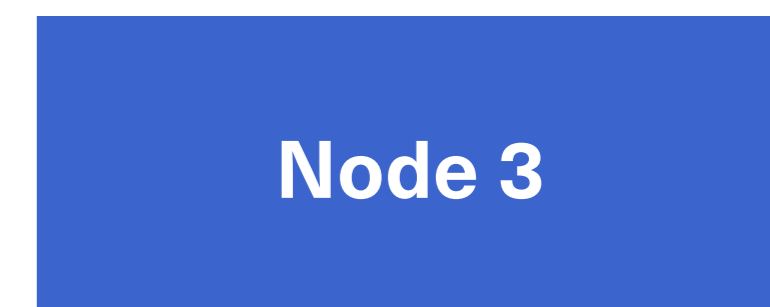
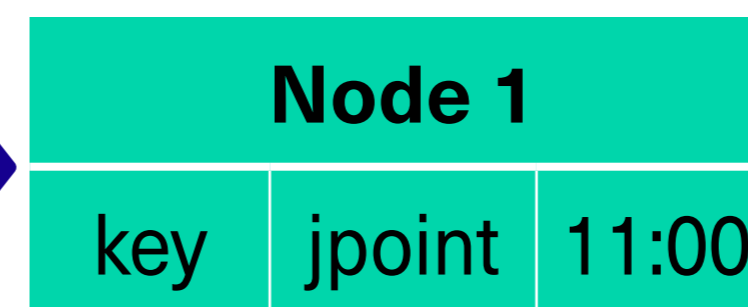
RF=3



SELECT

CL = QUORUM →

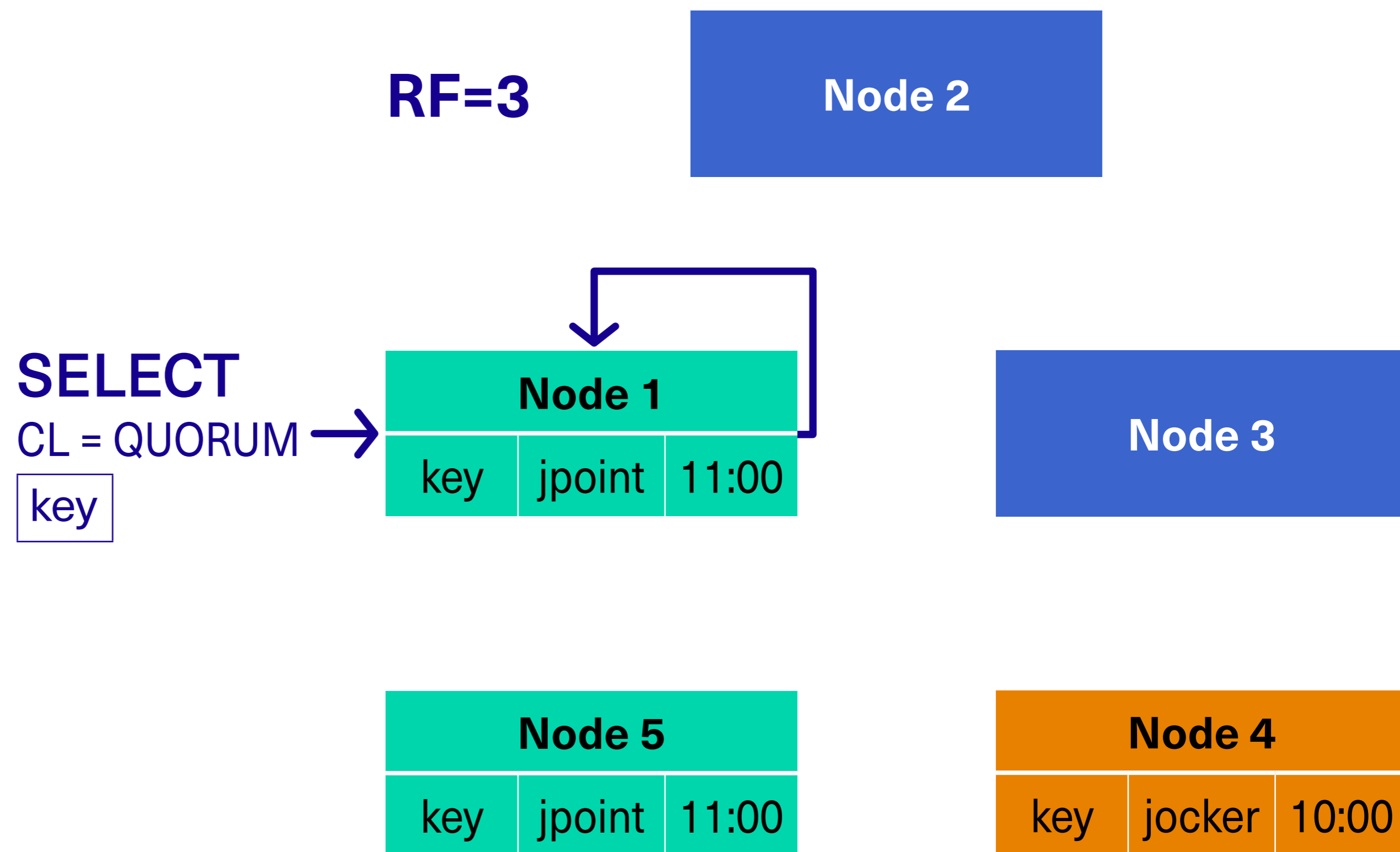
key



# Настраиваемая согласованность

**Replication factor (RF)** — сколько узлов хранят копии одной партии

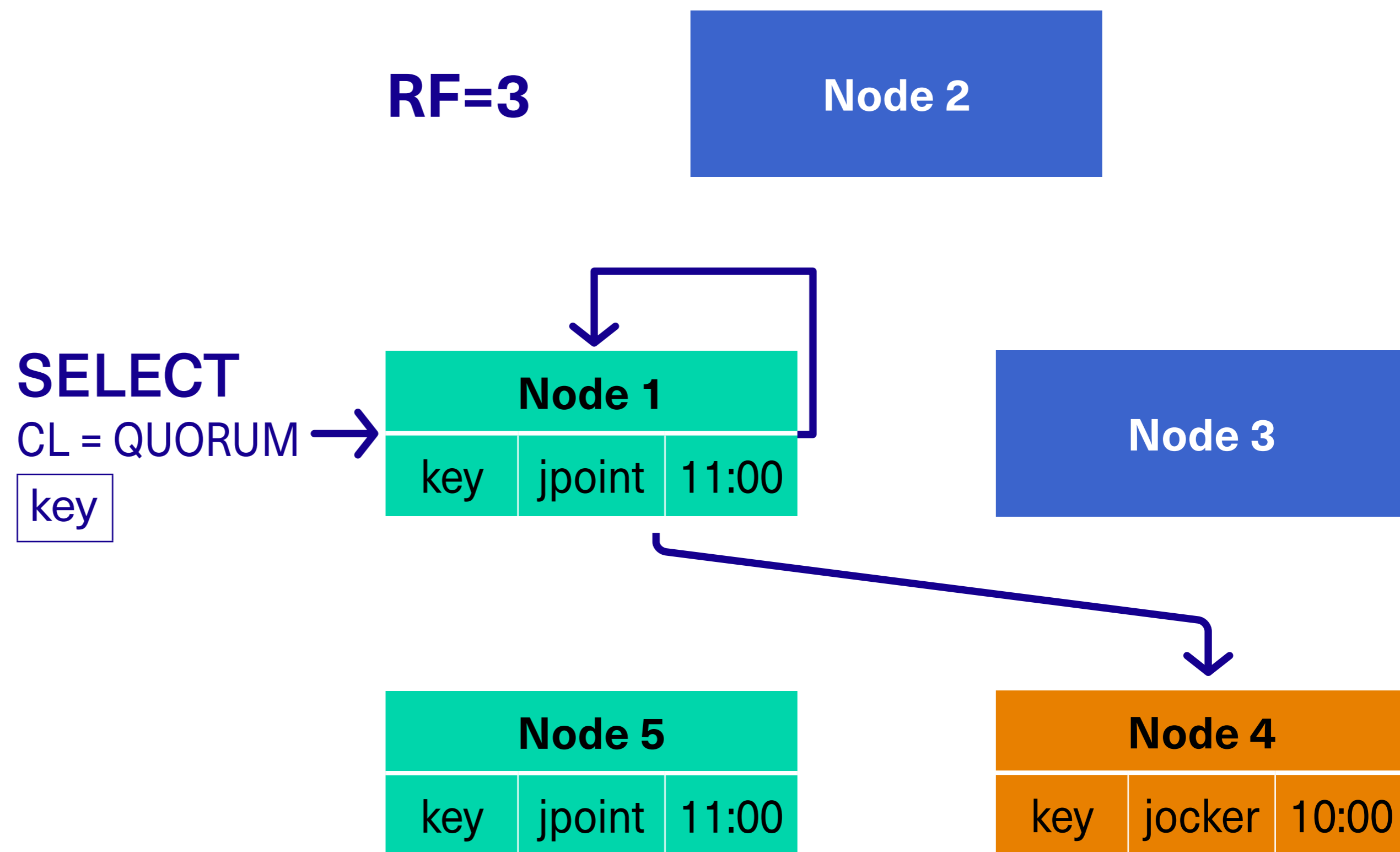
**Consistency level (CL)** — сколько из этих узлов должны подтвердить операцию, чтобы вернуть клиенту «успех»



# Настраиваемая согласованность

**Replication factor (RF)** — сколько узлов хранят копии одной партии

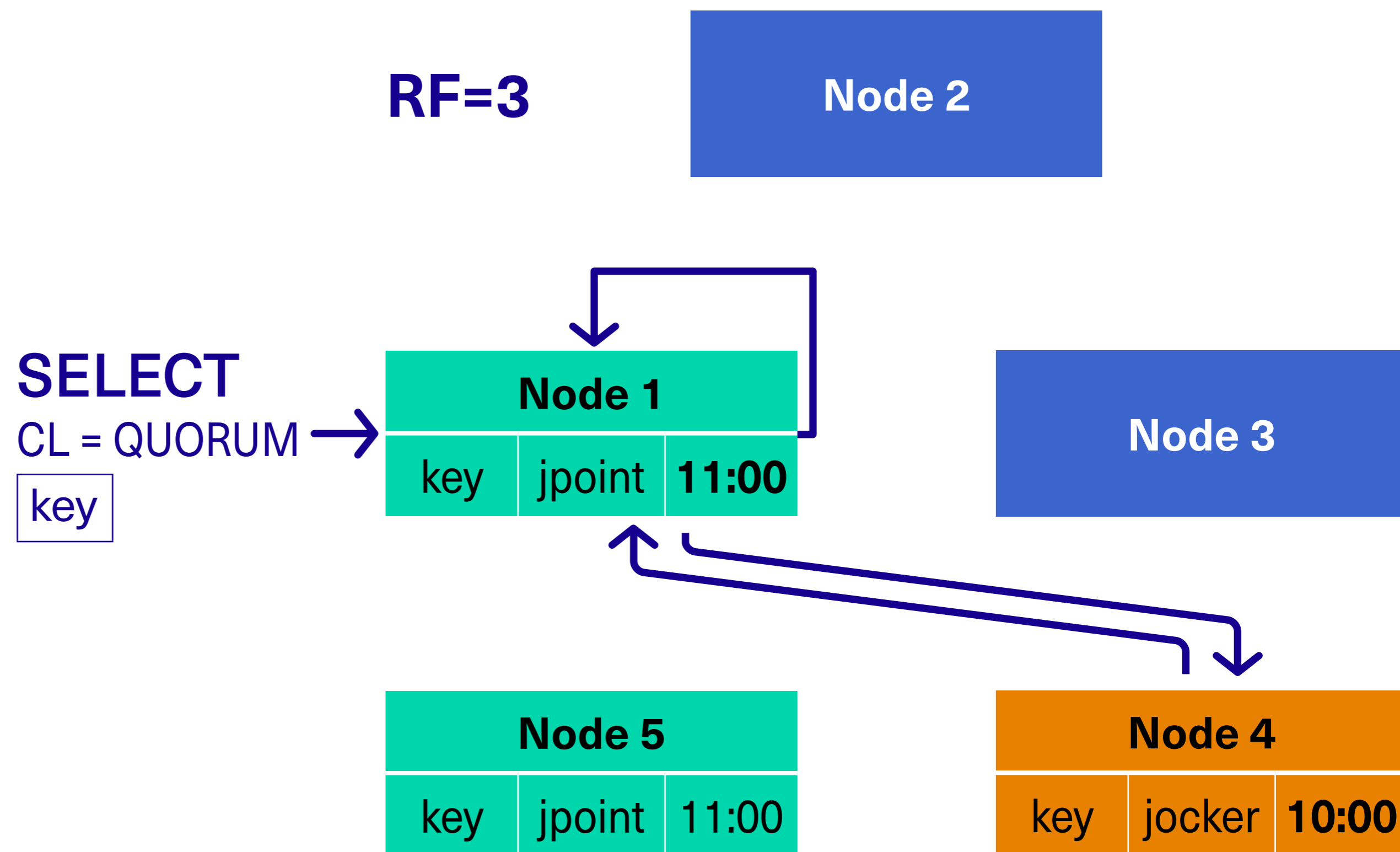
**Consistency level (CL)** — сколько из этих узлов должны подтвердить операцию, чтобы вернуть клиенту «успех»



# Настраиваемая согласованность

**Replication factor (RF)** — сколько узлов хранят копии одной партии

**Consistency level (CL)** — сколько из этих узлов должны подтвердить операцию, чтобы вернуть клиенту «успех»

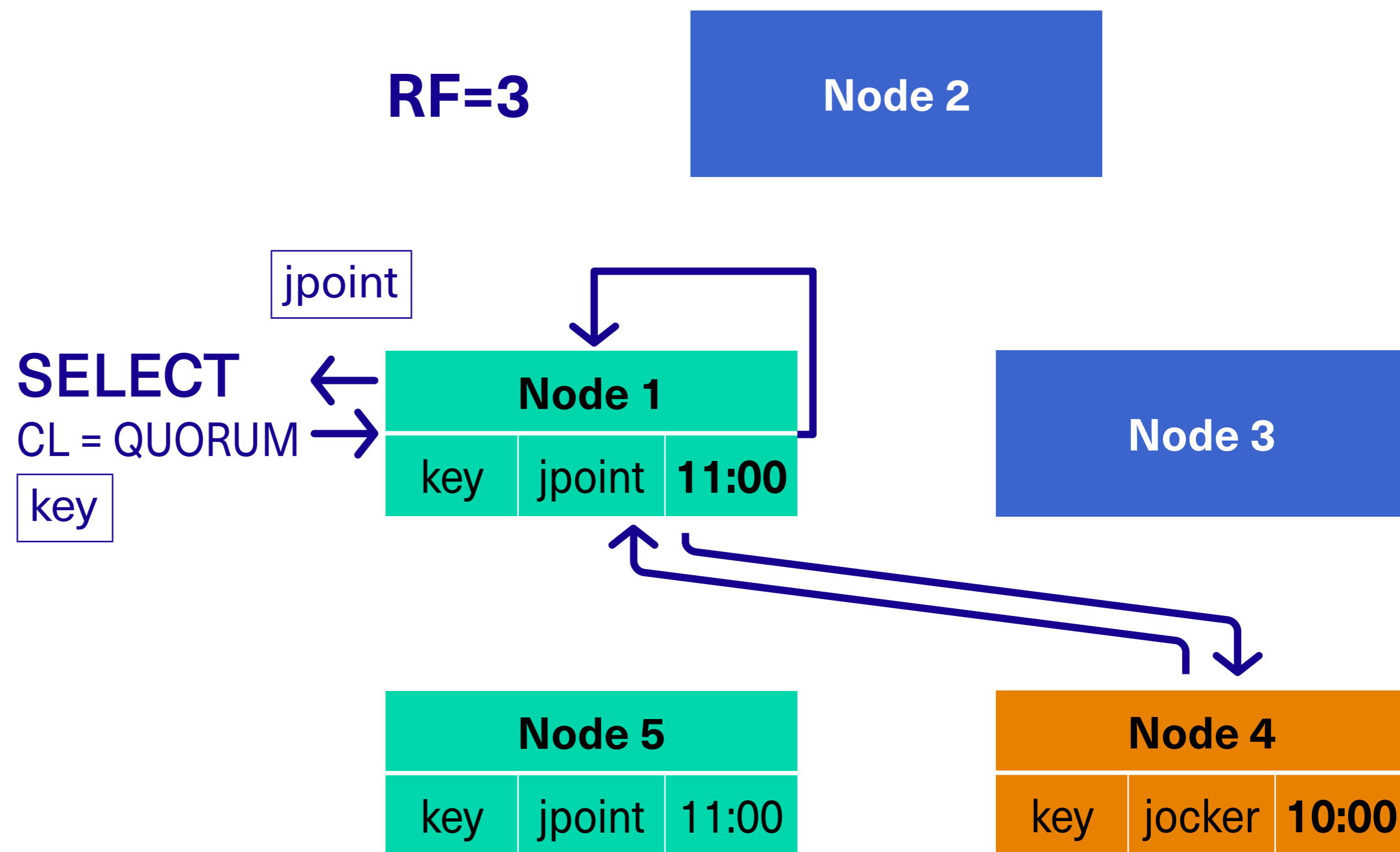




# Настраиваемая согласованность

**Replication factor (RF)** — сколько узлов хранят копии одной партии

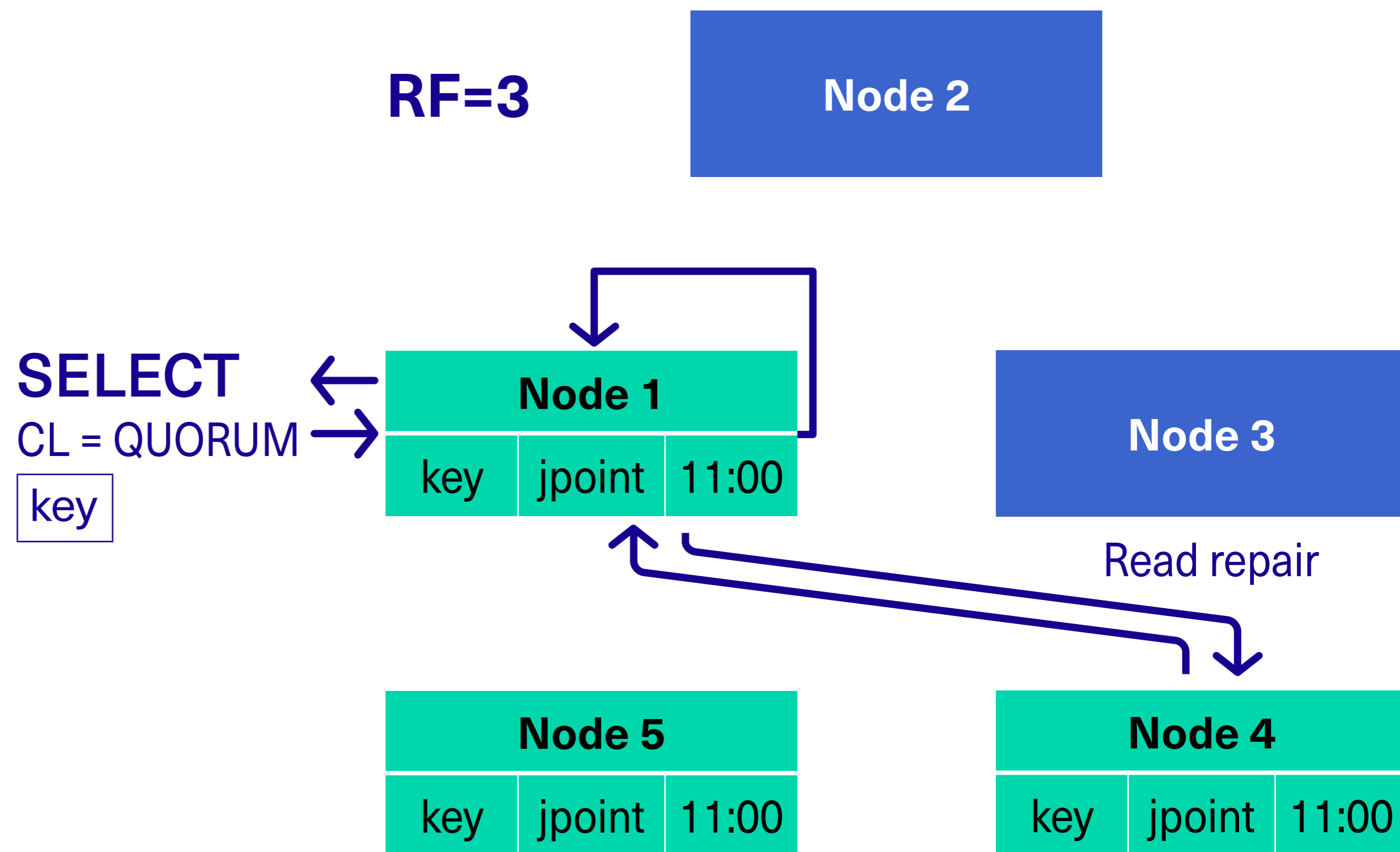
**Consistency level (CL)** — сколько из этих узлов должны подтвердить операцию, чтобы вернуть клиенту «успех»



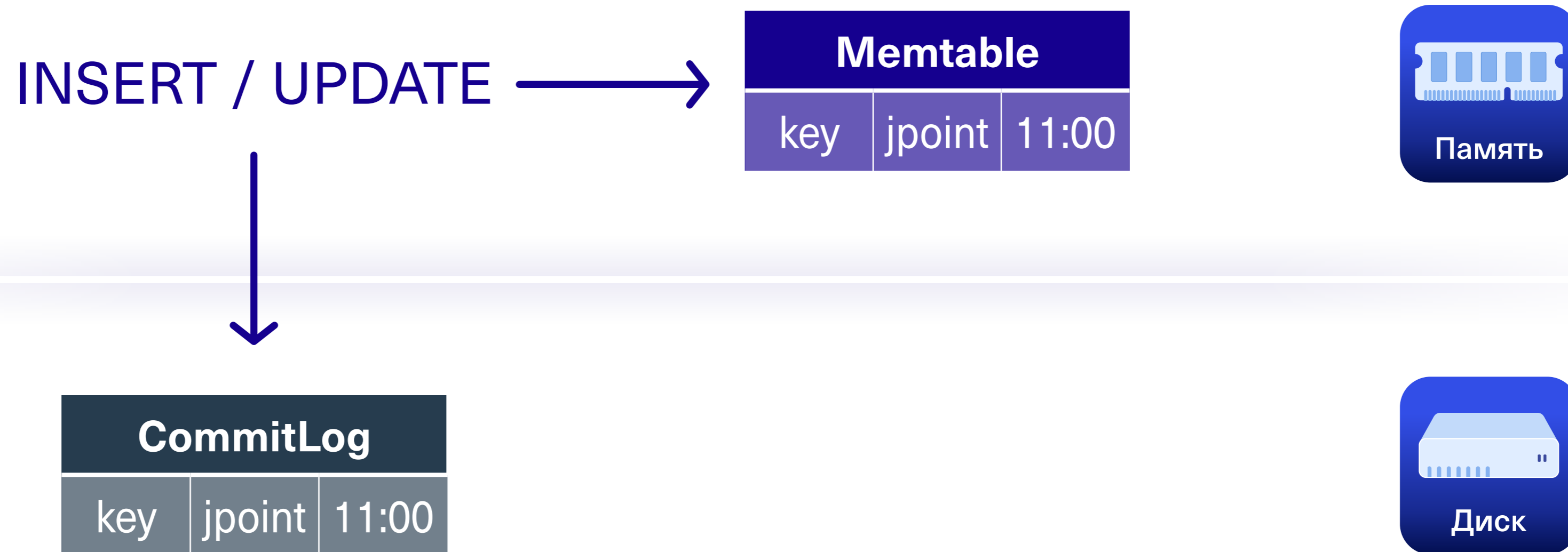
# Настраиваемая согласованность

**Replication factor (RF)** — сколько узлов хранят копии одной партии

**Consistency level (CL)** — сколько из этих узлов должны подтвердить операцию, чтобы вернуть клиенту «успех»

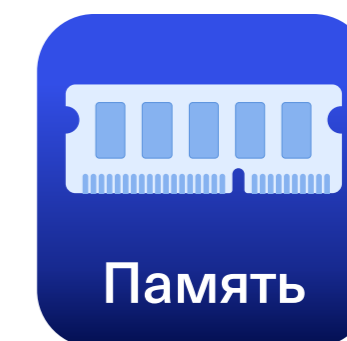


# Запись



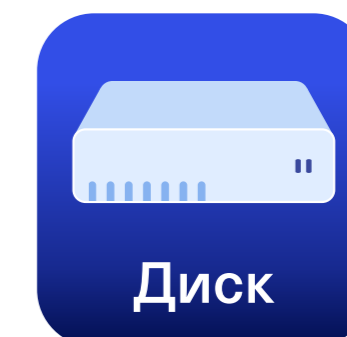
Сброс на диск

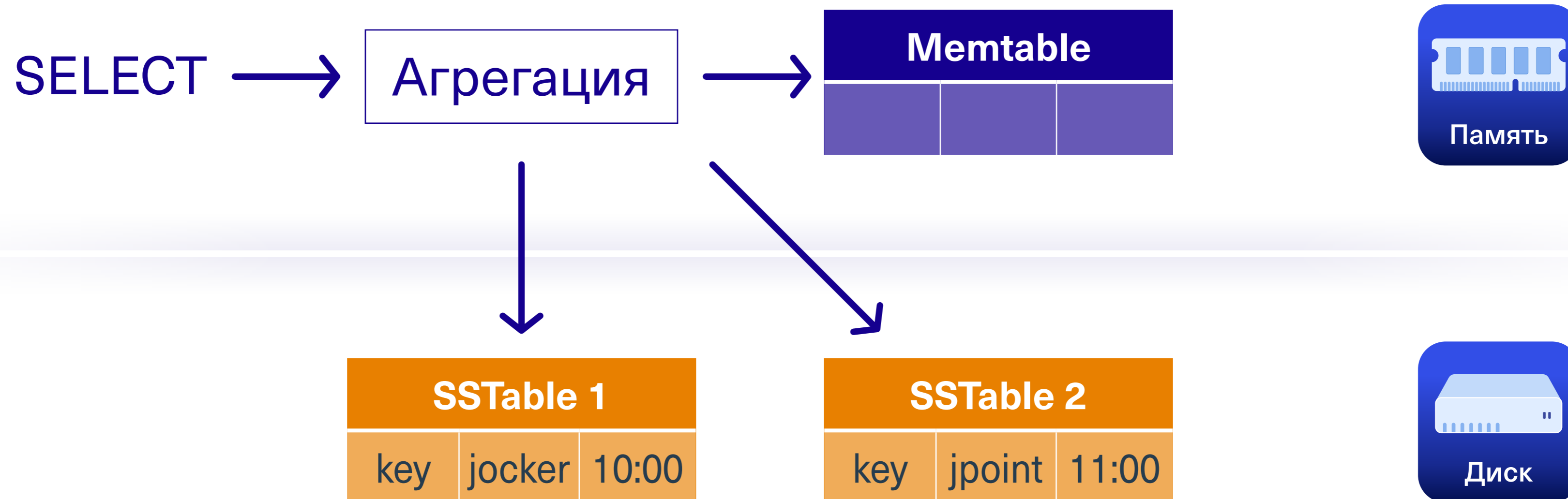
Memtable		
key	jpoint	11:00

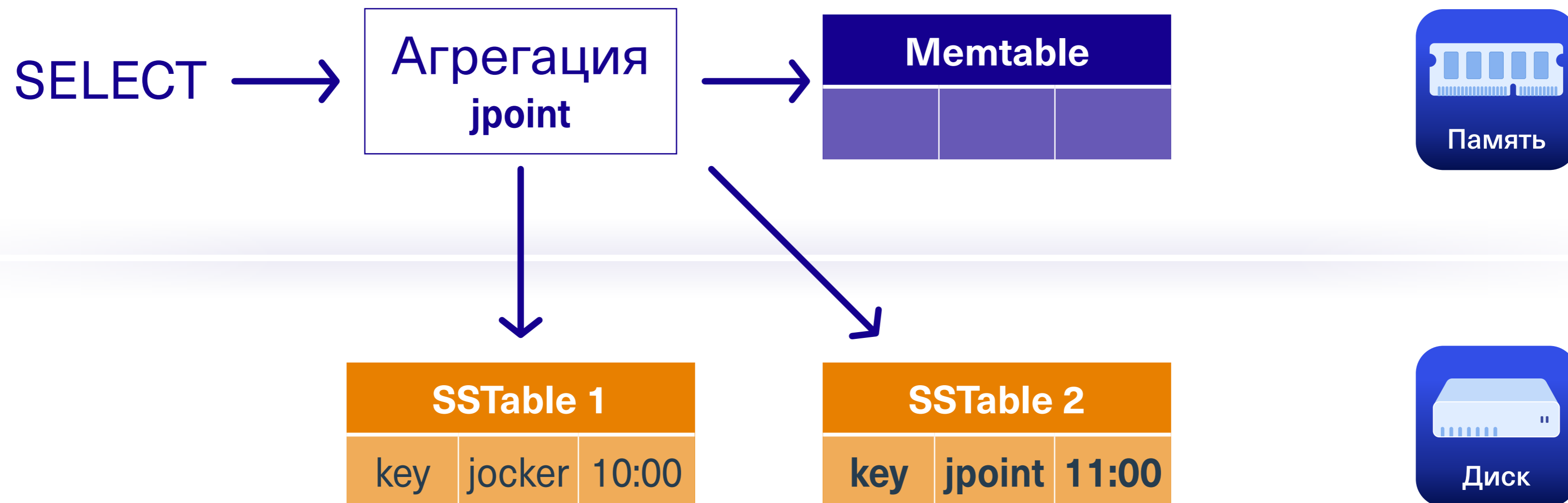


SSTable 1		
key	jocker	10:00

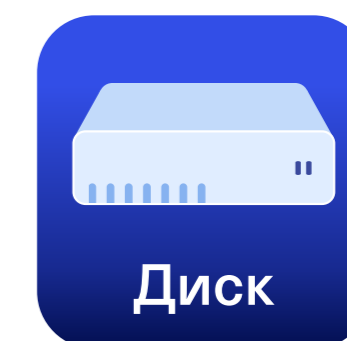
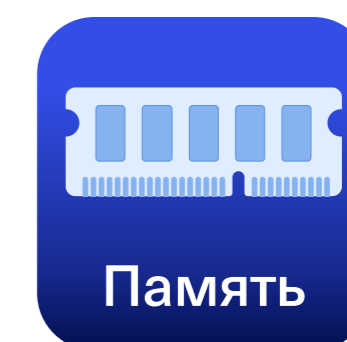
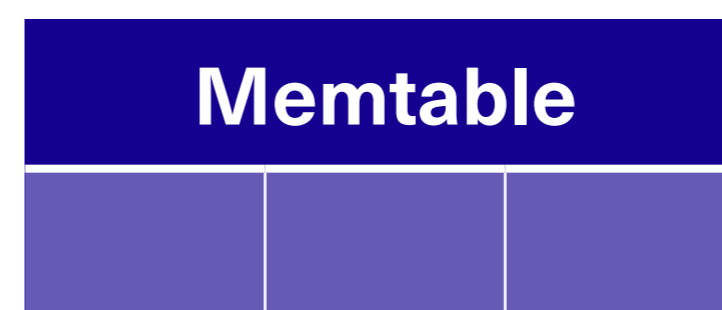
new SSTable		
key	jpoint	11:00







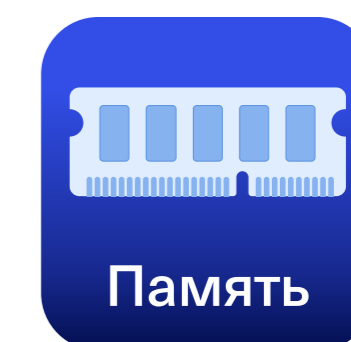
# Compaction



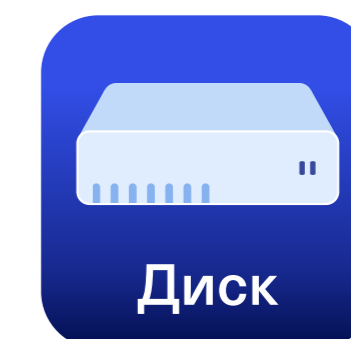
# Compaction



Memtable		

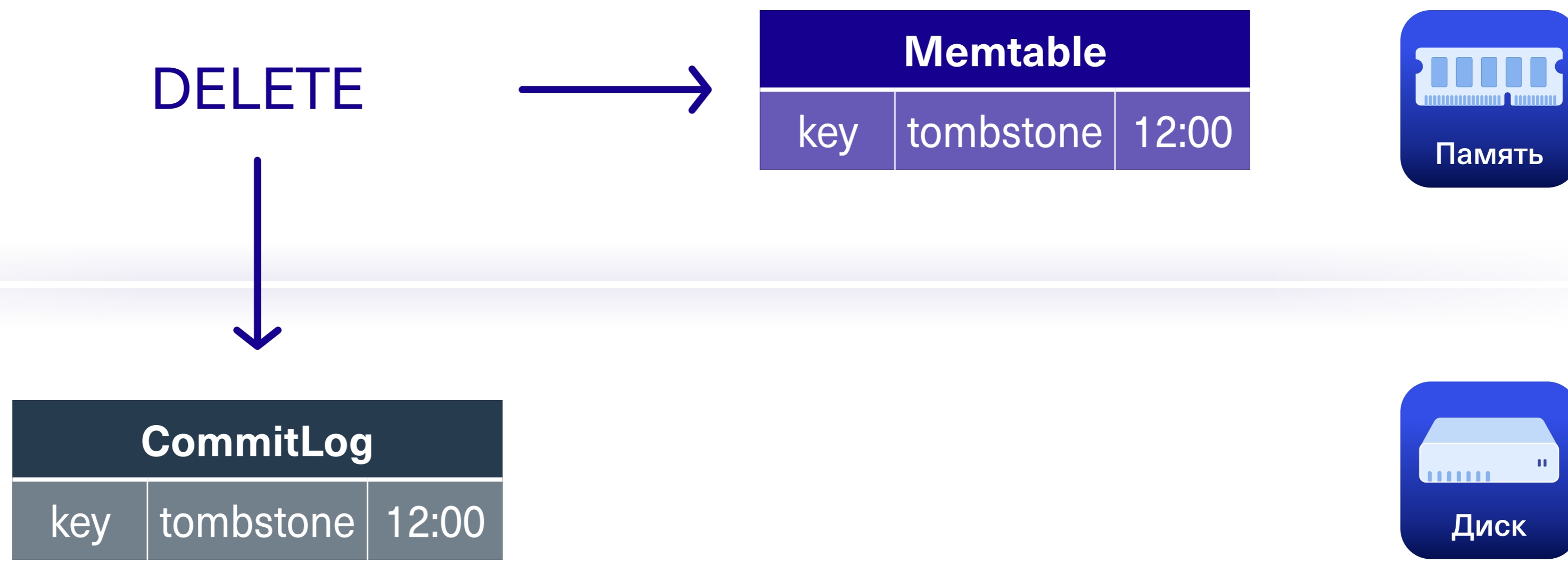


SSTable 2		
key	jpoint	11:00





# Удаление

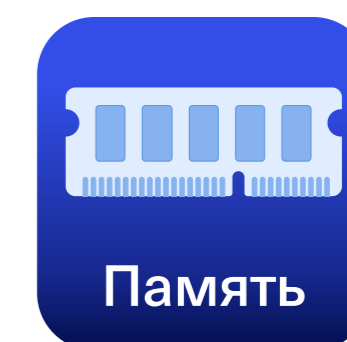


# Удаление



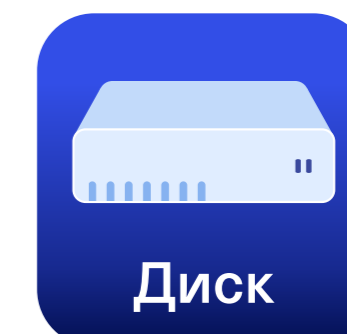
Сброс на диск

Memtable		
key	tombstone	12:00

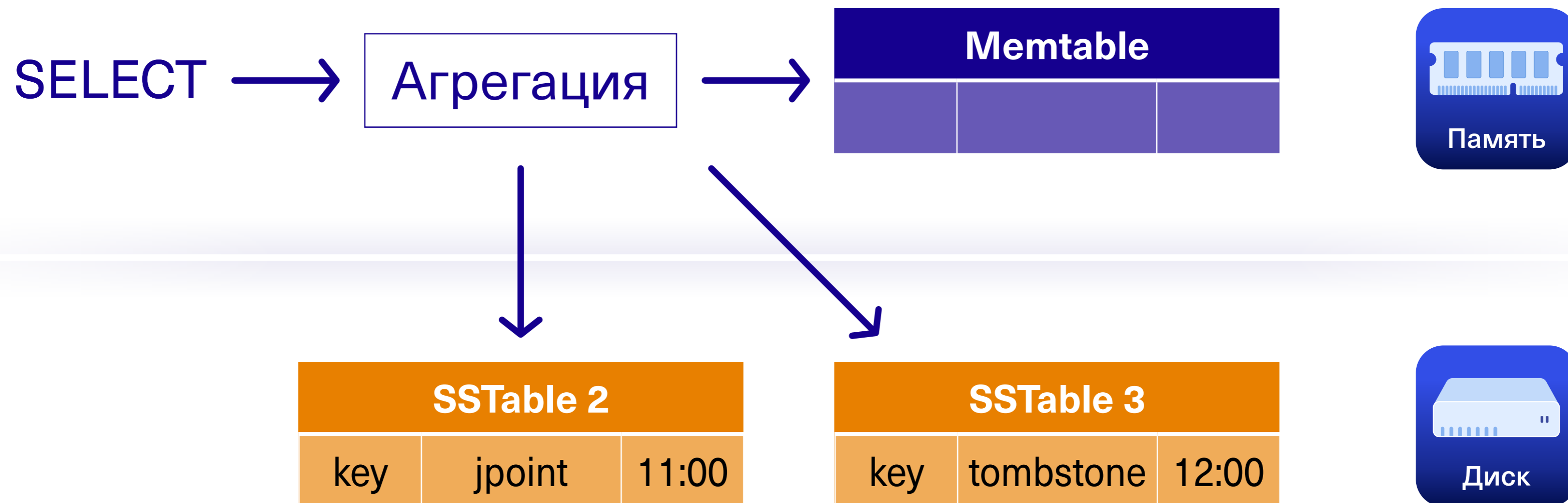


SSTable 2		
key	jpoint	11:00

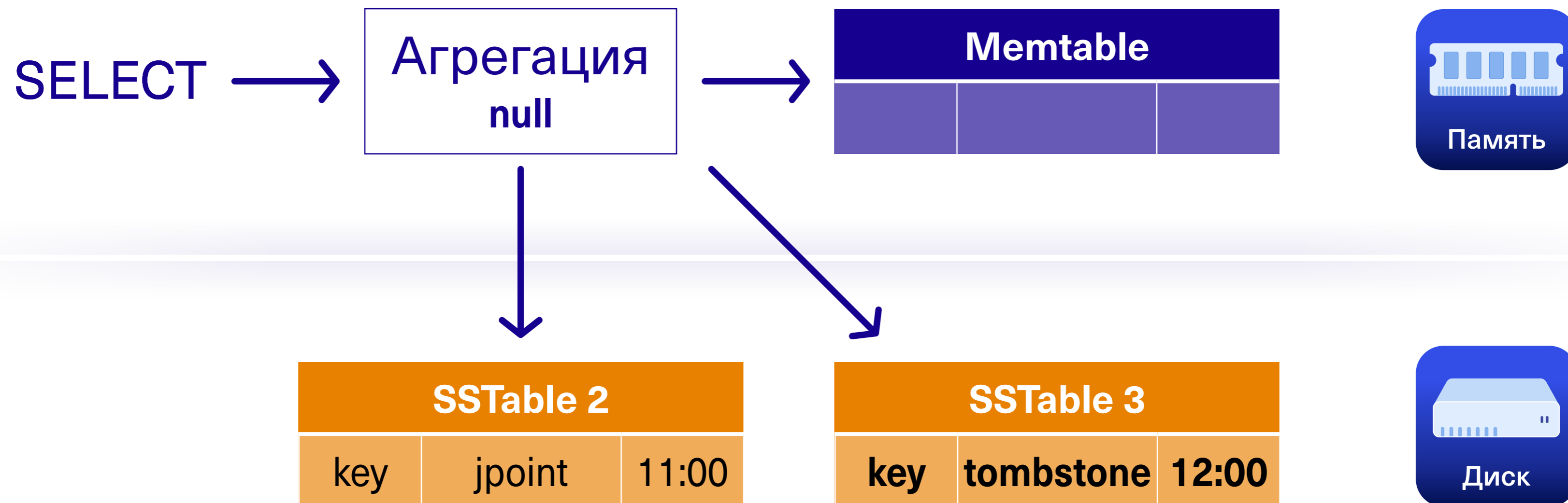
new SSTable		
key	tombstone	12:00



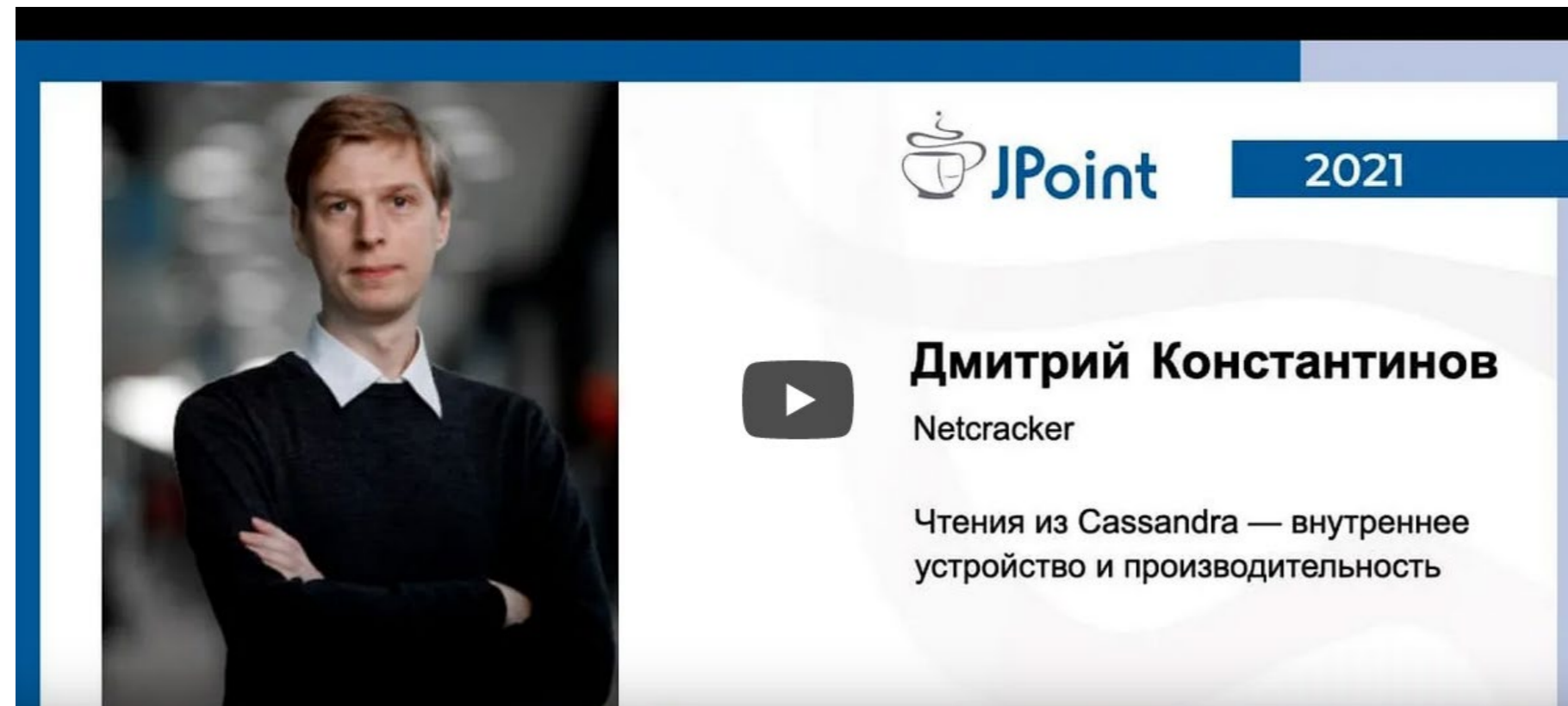
# Чтение удалённого значения



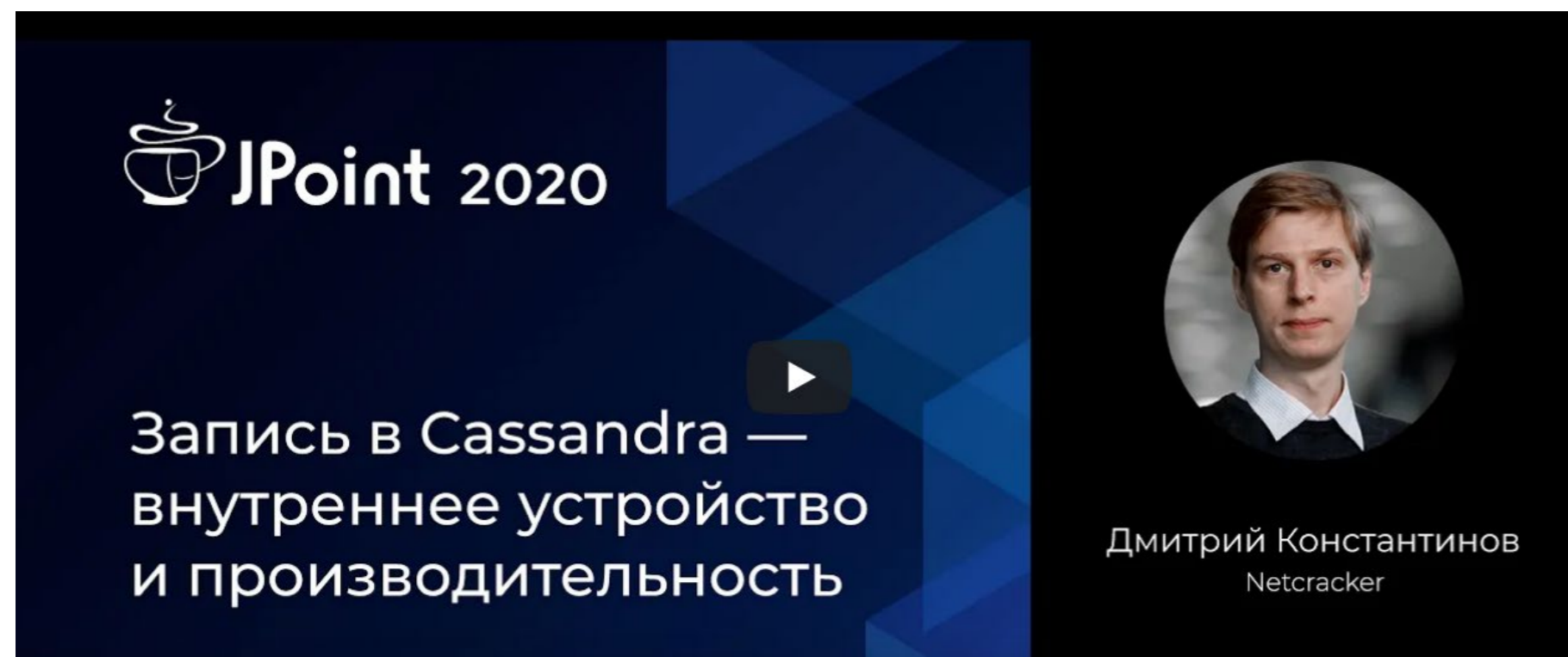
# Чтение удалённого значения



# Чтение и запись подробнее

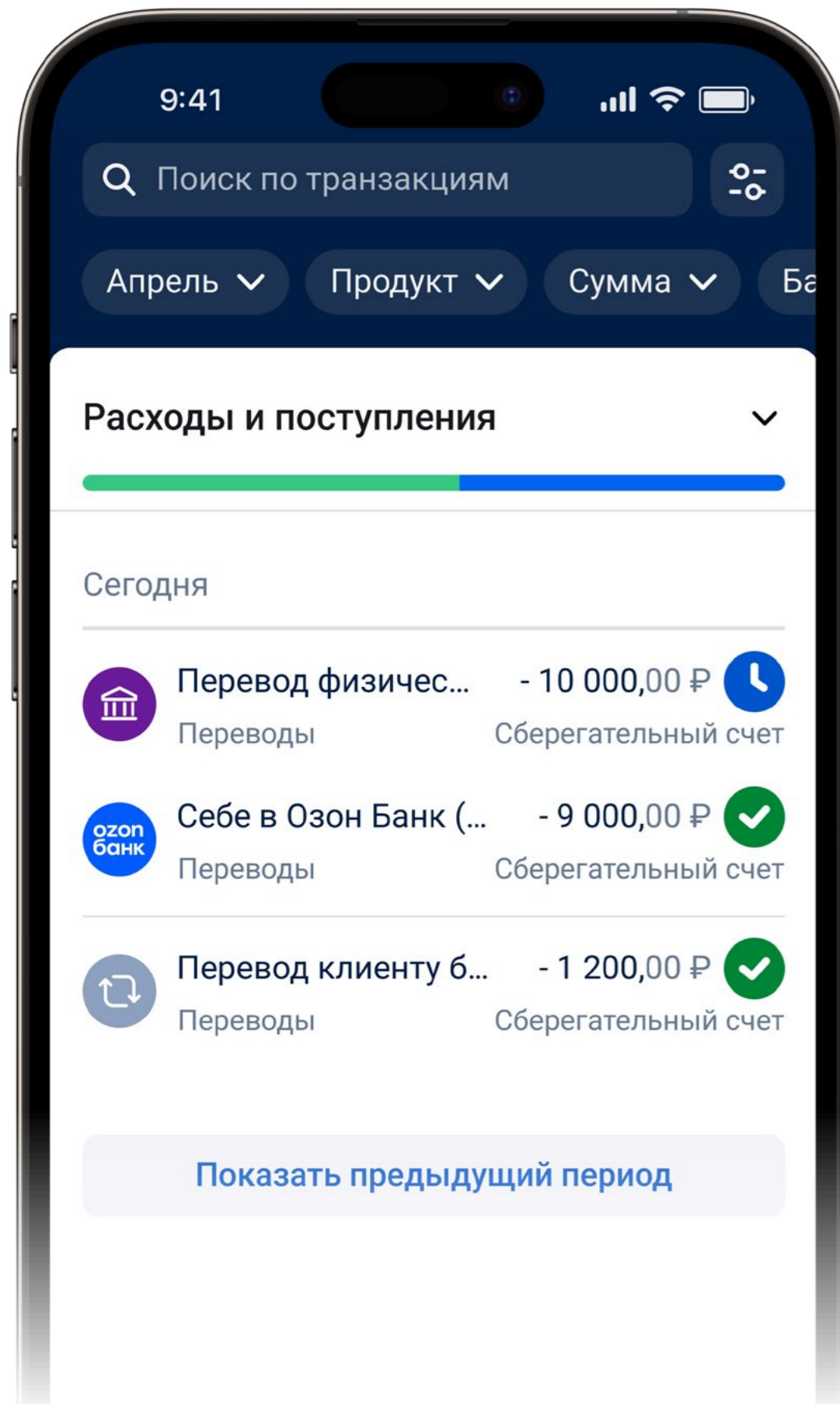


## «Чтение из Cassandra — внутреннее устройство и производительность» Jpoint 2021

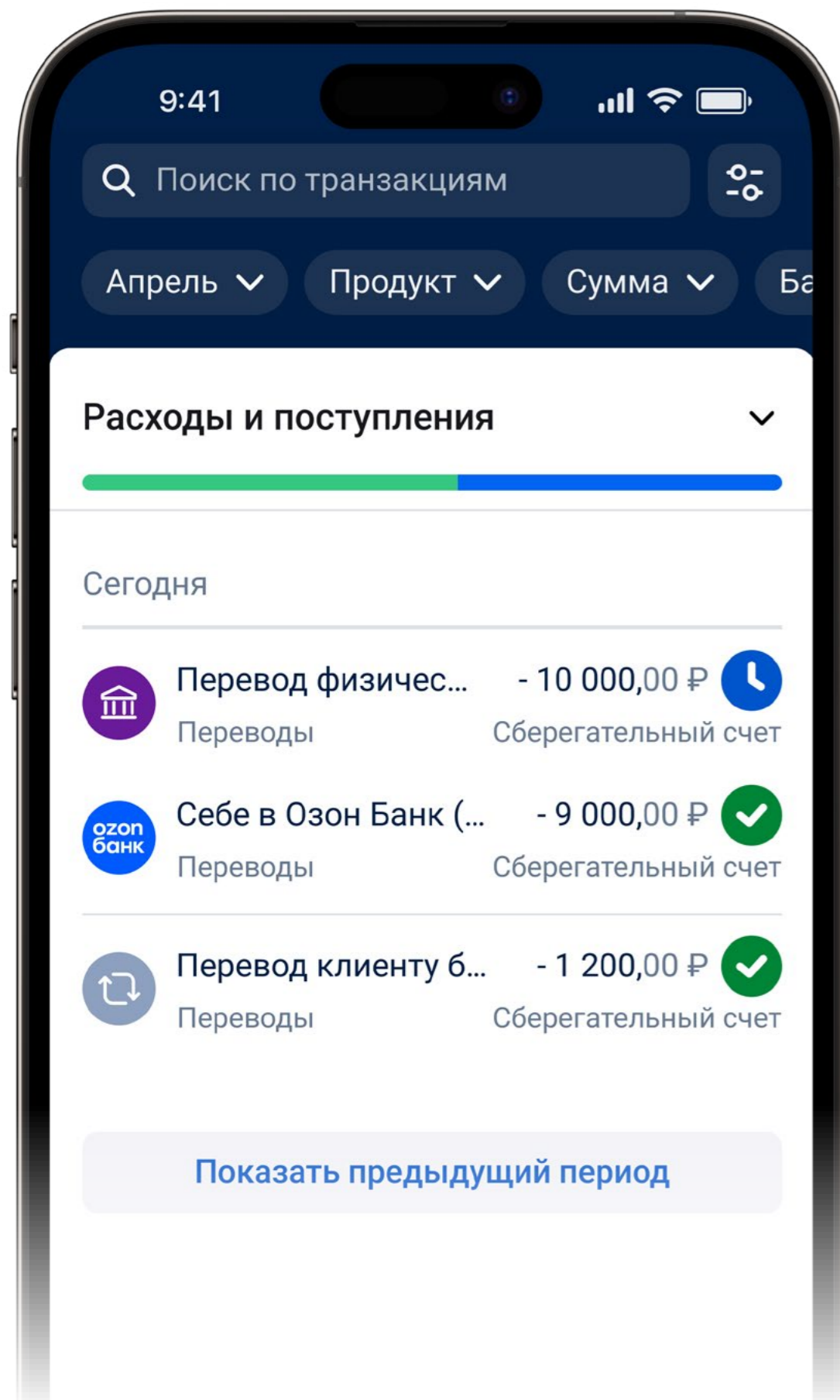


## «Запись в Cassandra — внутреннее устройство и производительность» Jpoint 2020

# Пишем ленту операций

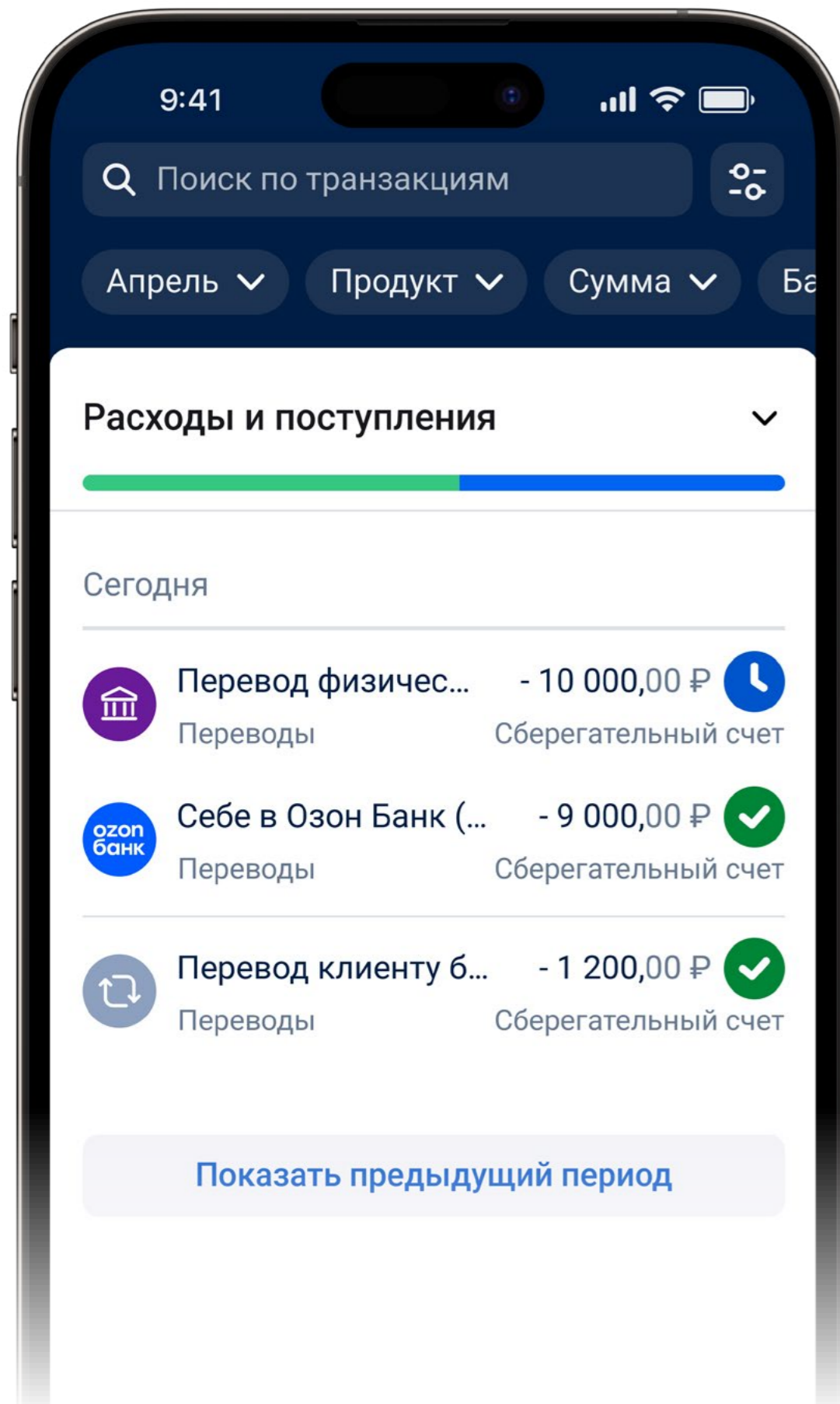


# Пишем ленту операций

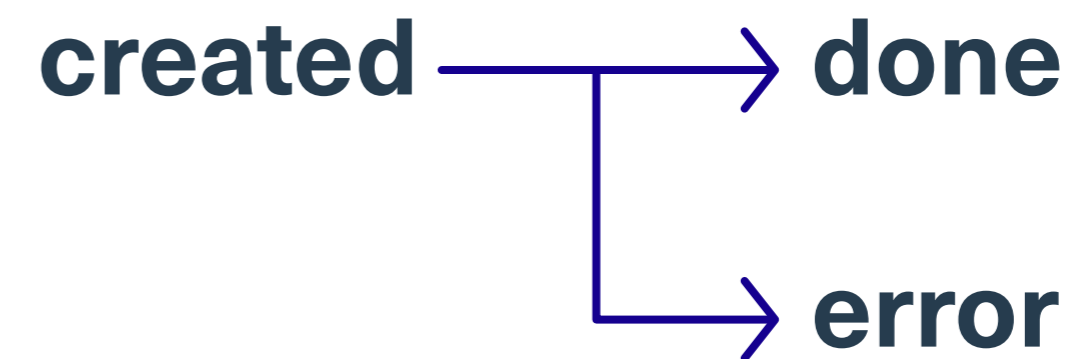


- Много записей, удалений нет, поиск ID клиента и дате

# Пишем ленту операций

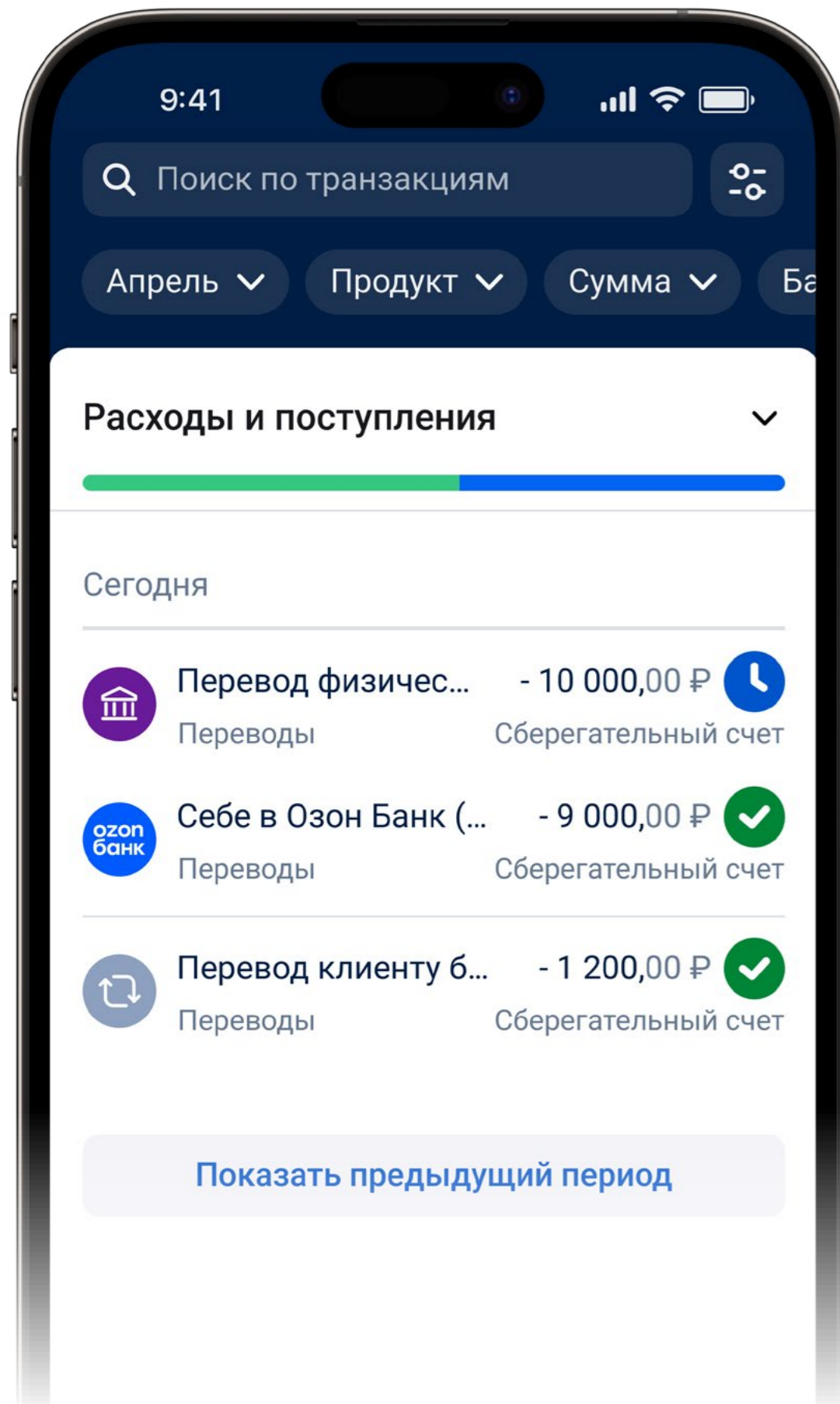


- Много записей, удалений нет, поиск ID клиента и дате
- Статус операции может меняться





# Пишем ленту операций

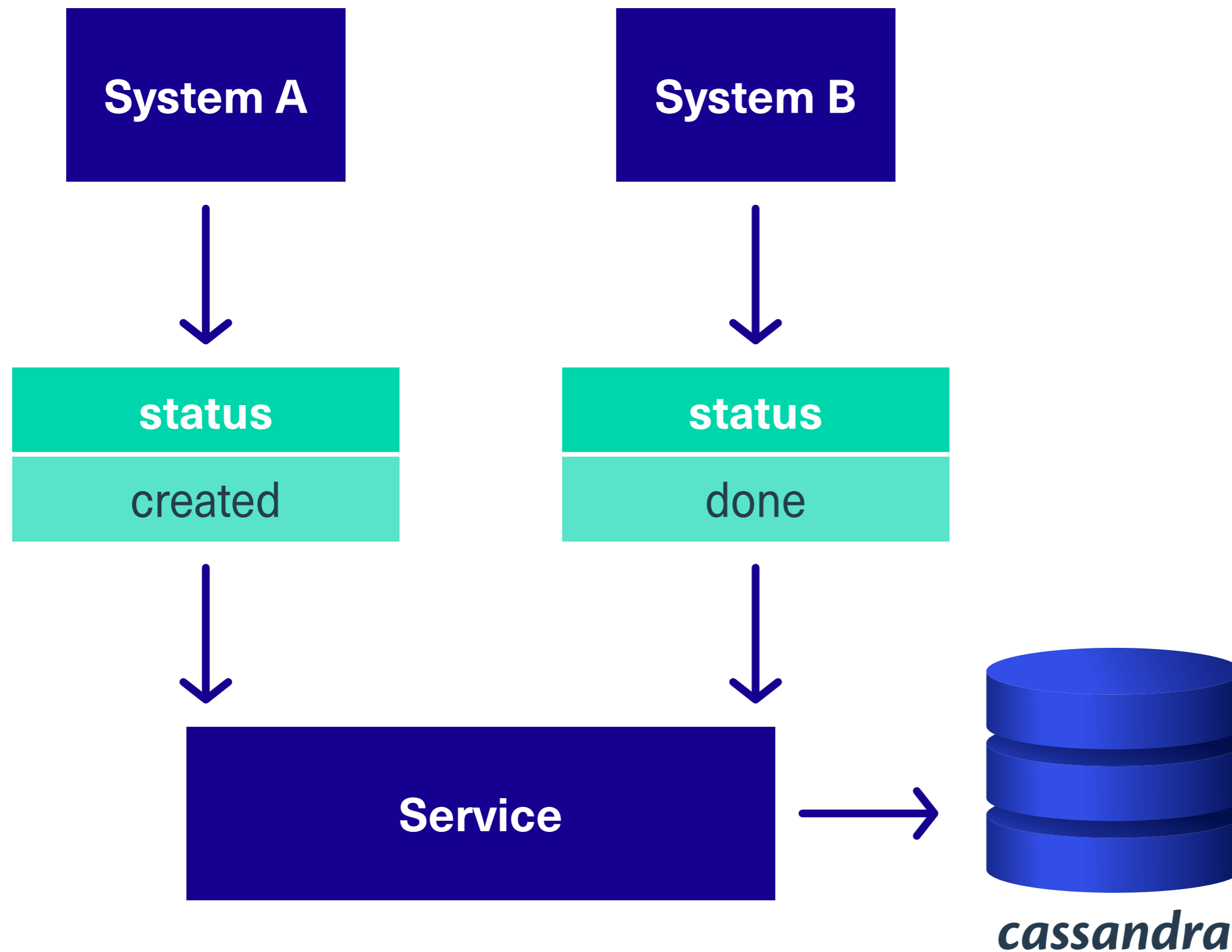


- Много записей, удалений нет, поиск ID клиента и дате
- Статус операции может меняться

«С SQL на Cassandra»,  
Jocker 2023

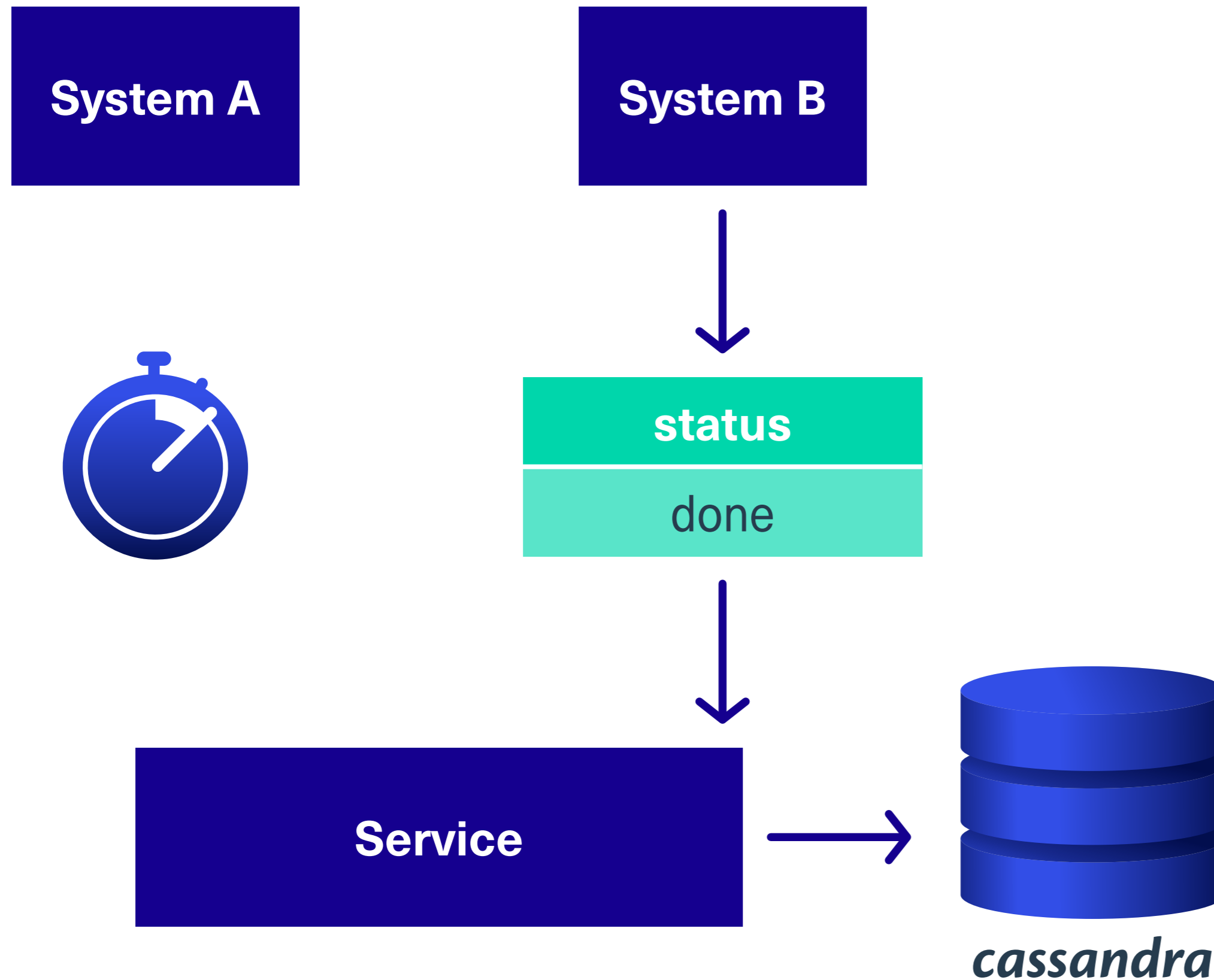
Таблицы необходимо проектировать от запросов!

# Пишем ленту операций



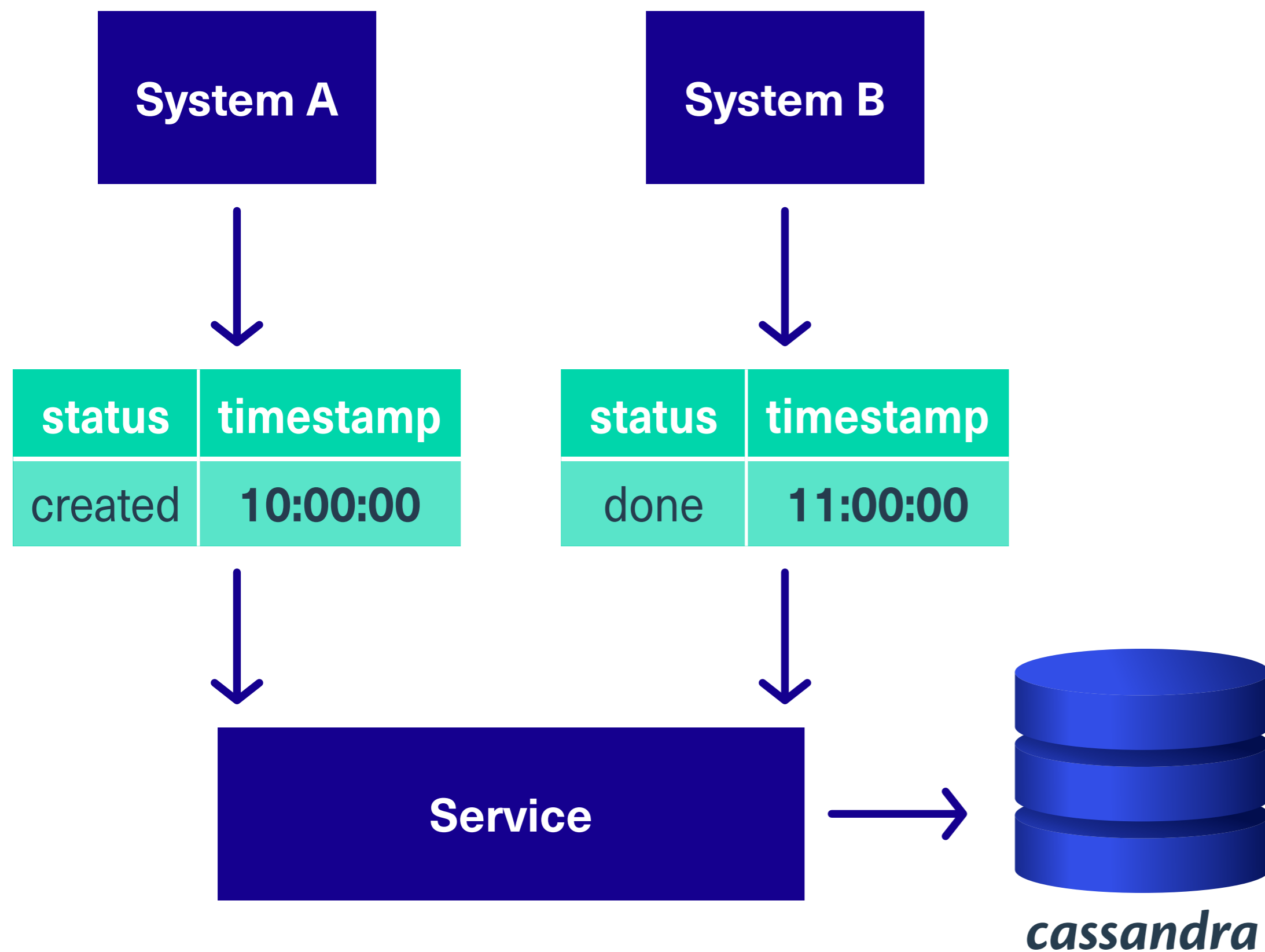
- Много записей, удалений нет, поиск ID клиента и дате
- Статус операции может меняться
- Порядок событий обновления статуса не гарантирован

# Пишем ленту операций



- Много записей, удалений нет, поиск ID клиента и дате
- Статус операции может меняться
- Порядок событий обновления статуса не гарантирован

# Пишем ленту операций



- Много записей, удалений нет, поиск ID клиента и дате
- Статус операции может меняться
- Порядок событий обновления статуса не гарантирован
- Есть временная метка по которой можно понять порядок

# Пишем ленту операций



## 1. Lightweight transactions (LWT) <, <=, >, >=, !=, IN, EXISTS, NOT EXISTS

```
UPDATE conf.timeline
SET status = 'done'
WHERE client_id = 1
AND year = 2024
AND month = 4
AND date = '2024-04-25'
AND op_id = 333
IF time < '11:00:00'
```

status	timestamp
done	11:00:00

# Пишем ленту операций



1. Lightweight transactions (LWT)  $<$ ,  $<=$ ,  $>$ ,  $>=$ ,  $!=$ , IN, EXISTS, NOT EXISTS  
но это дорого т. к. 4 запроса вместо одного (и один из них чтение)

```
UPDATE conf.timeline
SET status = 'done'
WHERE client_id = 1
AND year = 2024
AND month = 4
AND date = '2024-04-25'
AND op_id = 333
IF time < '11:00:00'
```

status	timestamp
done	11:00:00

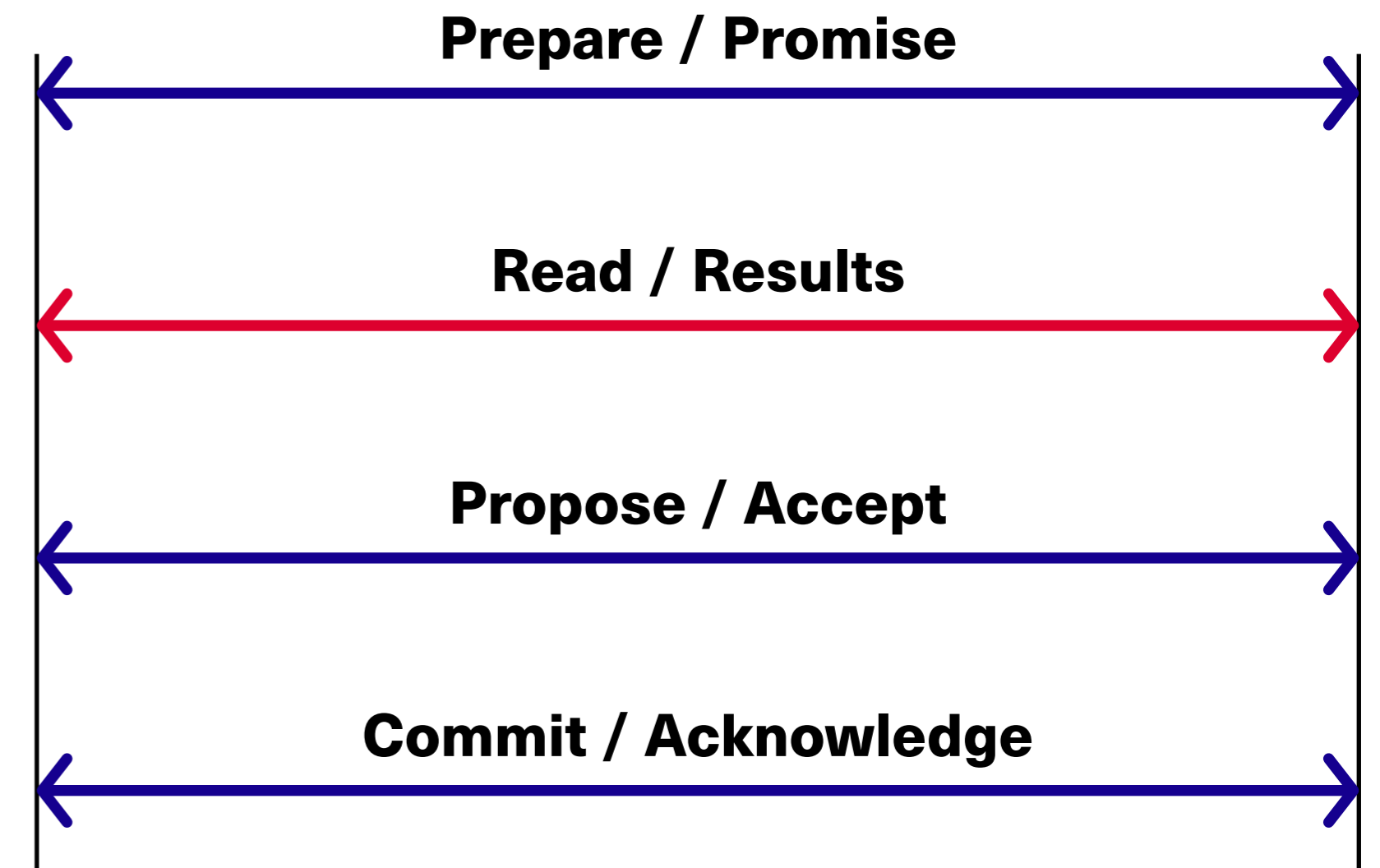
# Пишем ленту операций



1. Lightweight transactions (LWT)  $<$ ,  $<=$ ,  $>$ ,  $>=$ ,  $!=$ , IN, EXISTS, NOT EXISTS  
но это дорого т. к. 4 запроса вместо одного (и один из них чтение)

Compare and set (CAS) на базе PAXOS

[Light-weight transactions and Paxos algorithm](#)



# Пишем ленту операций



**1. Lightweight transactions (LWT) <, <=, >, >=, !=, IN, EXISTS, NOT EXISTS**  
но это дорого т. к. 4 запроса вместо одного (и один из них чтение)

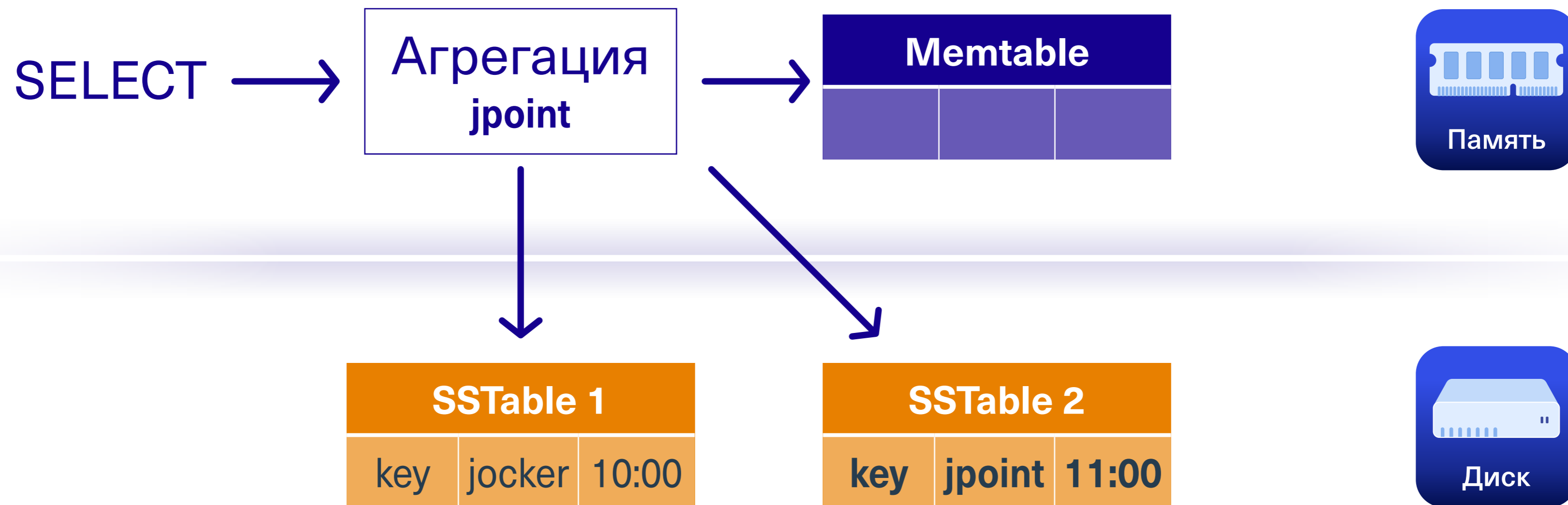
## **2. USING TIMESTAMP**

это дешево т.к. не добавляет новой логики

```
UPDATE conf.timeline
USING TIMESTAMP 1714028400000
SET status = 'done'
WHERE client_id = 1
and year = 2024
AND month = 4
and date = '2024-04-25'
and op_id = 333
```

status	timestamp
done	11:00:00





# Пишем ленту операций



**1. Lightweight transactions (LWT) <, <=, >, >=, !=, IN, EXISTS, NOT EXISTS**  
но это дорого т. к. 4 запроса вместо одного (и один из них чтение)

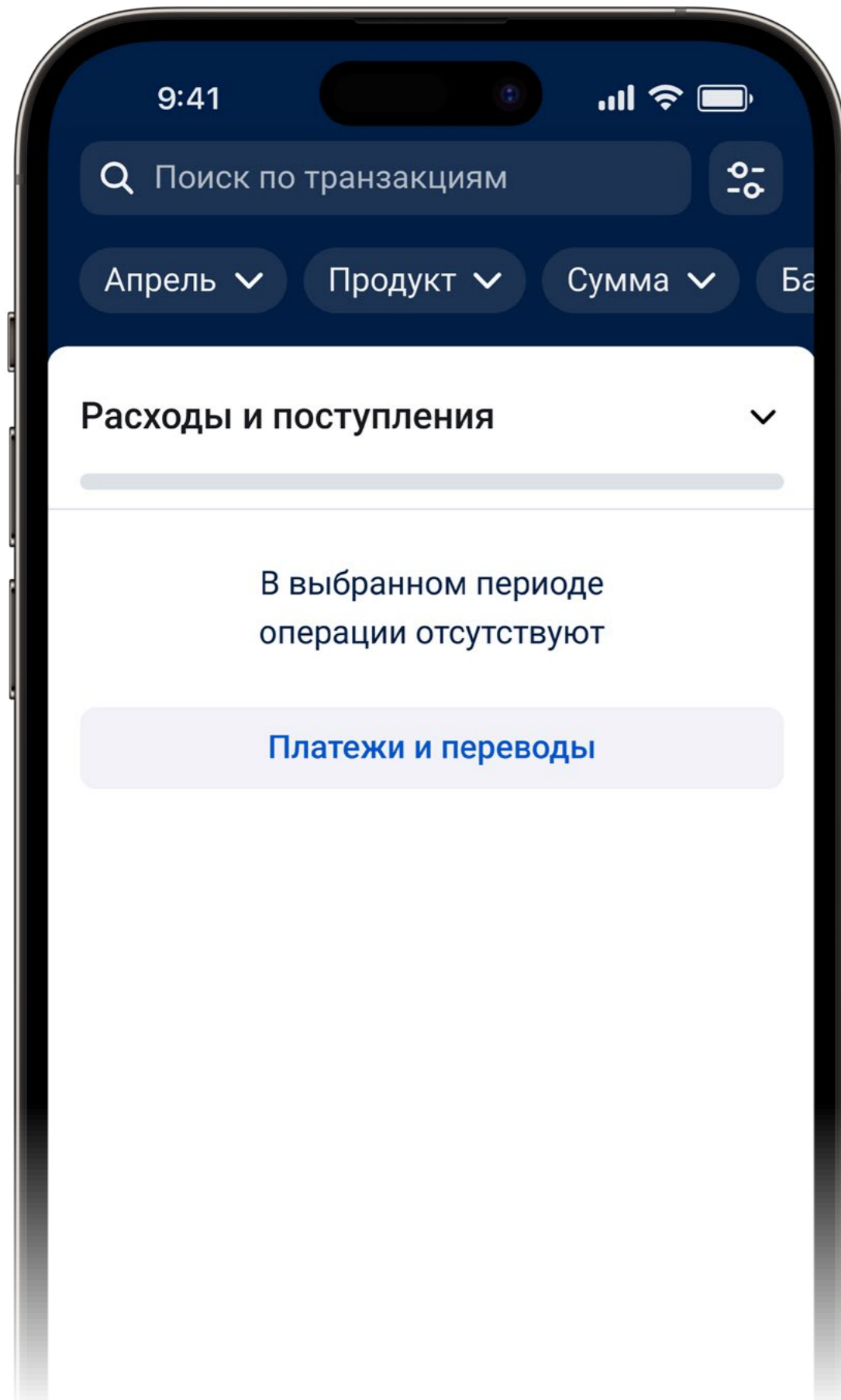
## **2. USING TIMESTAMP**

это дешево т.к. не добавляет новой логики

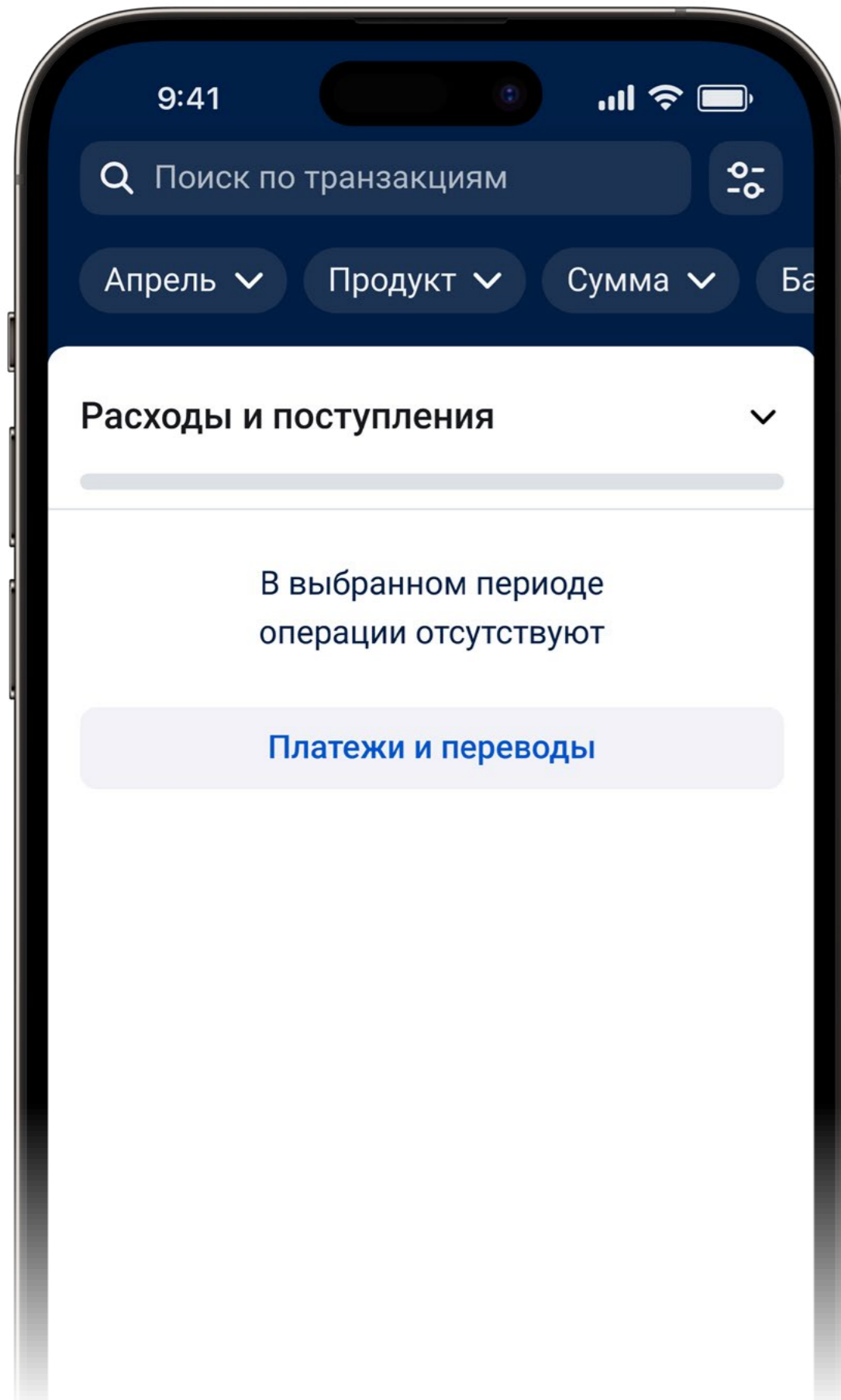
```
UPDATE conf.timeline
USING TIMESTAMP 1714028400000
SET status = 'done'
WHERE client_id = 1
and year = 2024
AND month = 4
and date = '2024-04-25'
and op_id = 333
```

status	timestamp
done	11:00:00

# Обновляем с timestamp

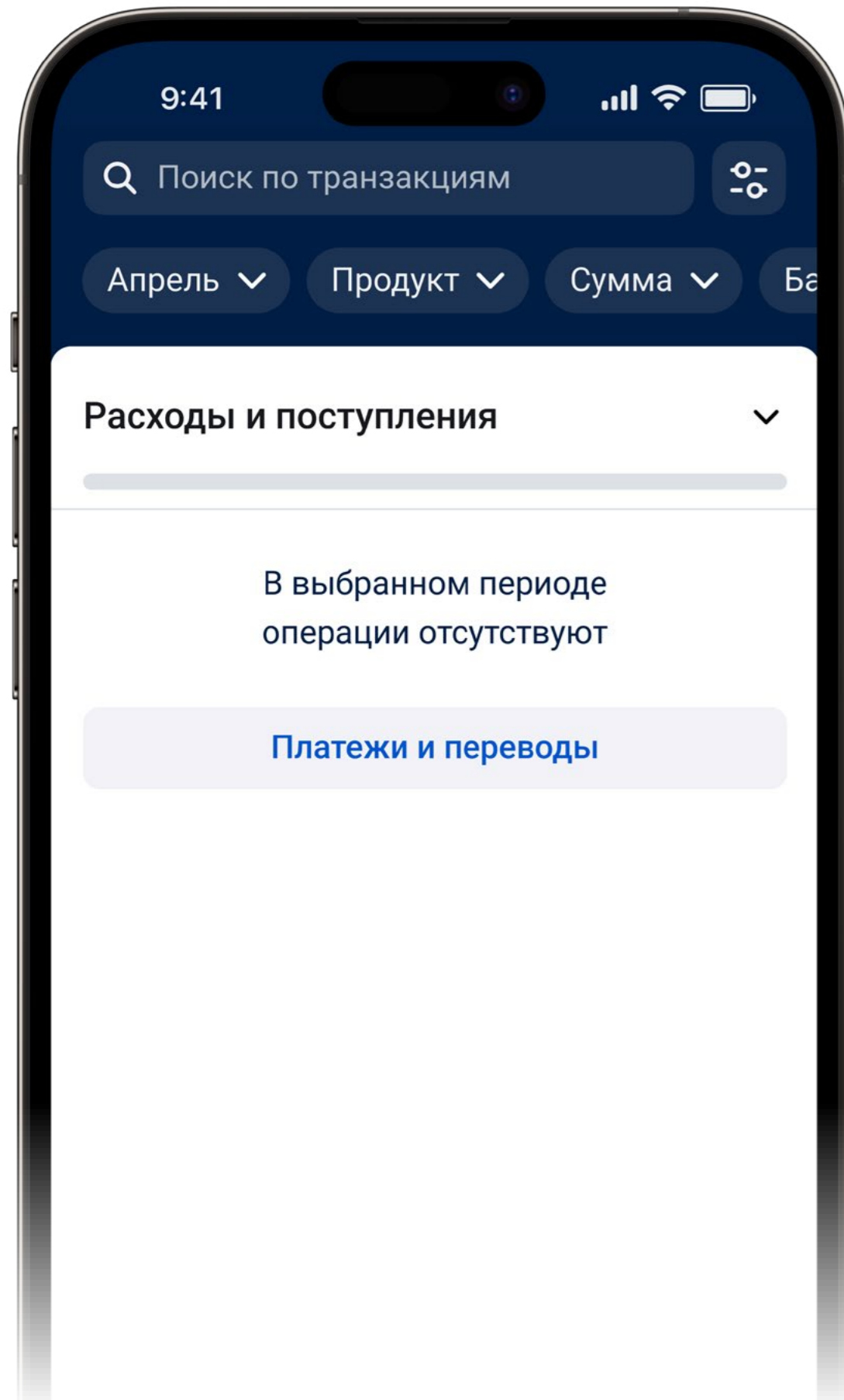


# Обновляем с timestamp

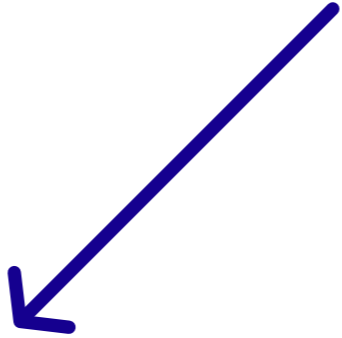


SSTable 1	
created	10:00

# Обновляем с timestamp

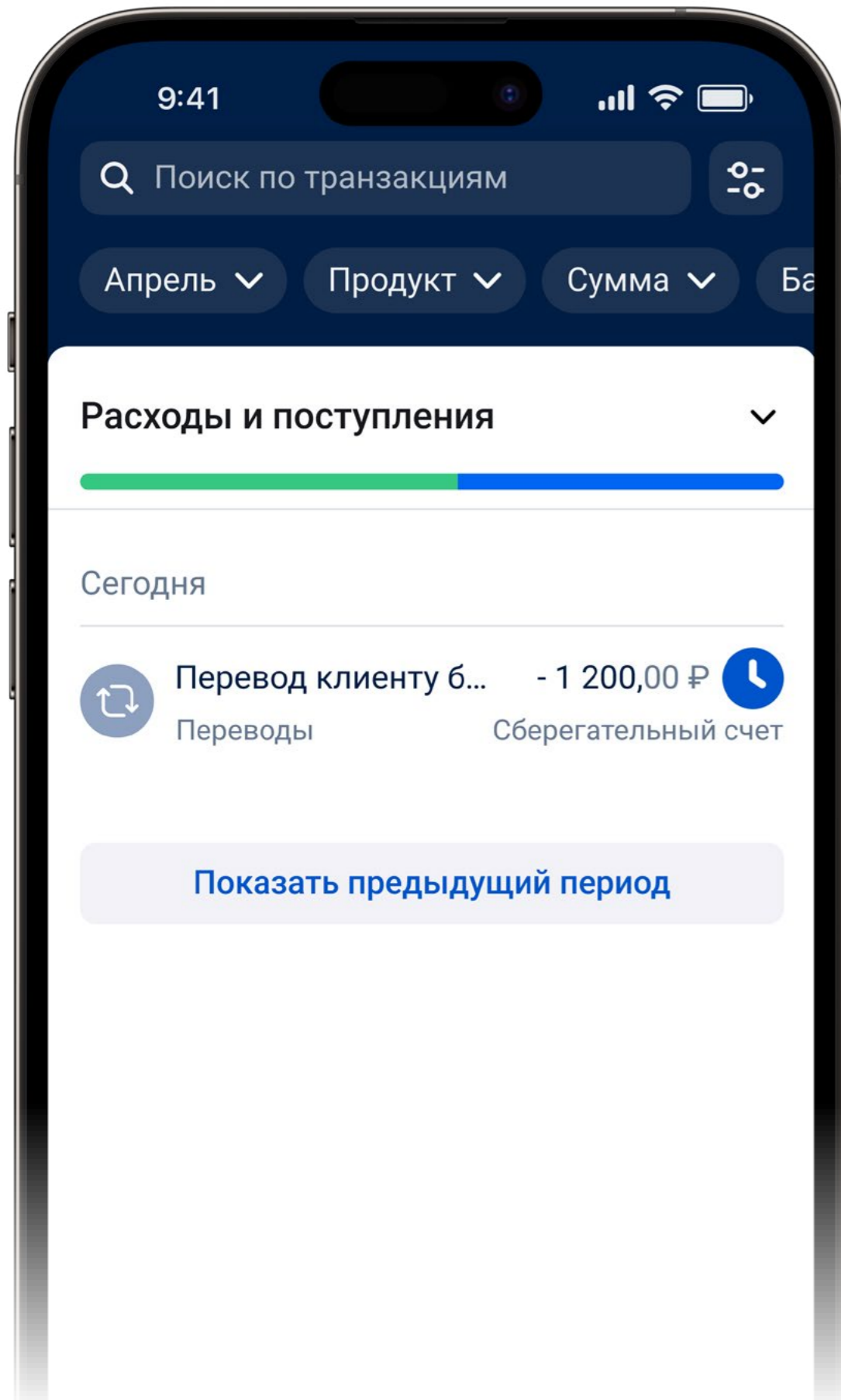


Агрегация  
status = **created**



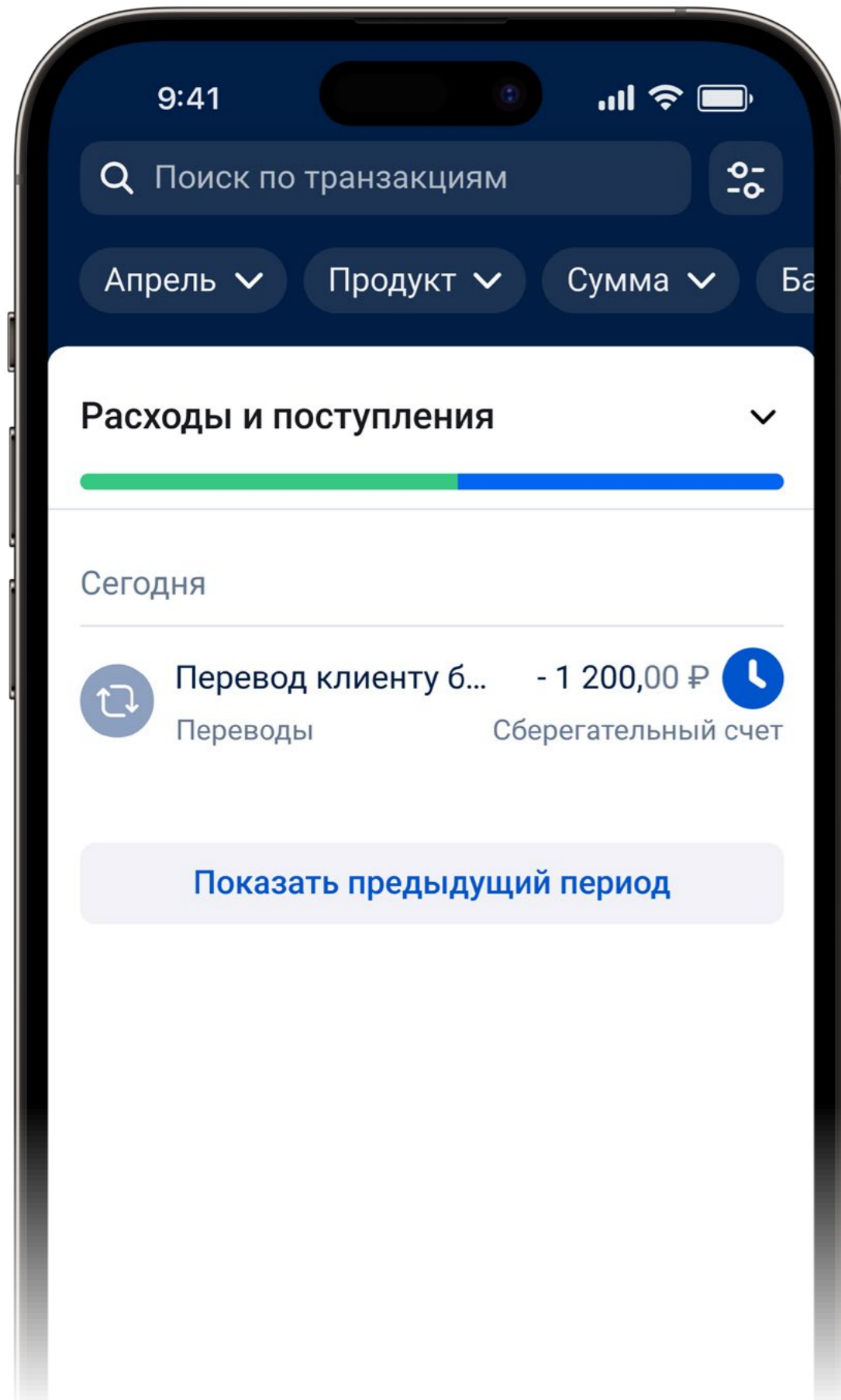
SStable 1	
created	10:00

# Обновляем с timestamp



SSTable 1	
created	10:00

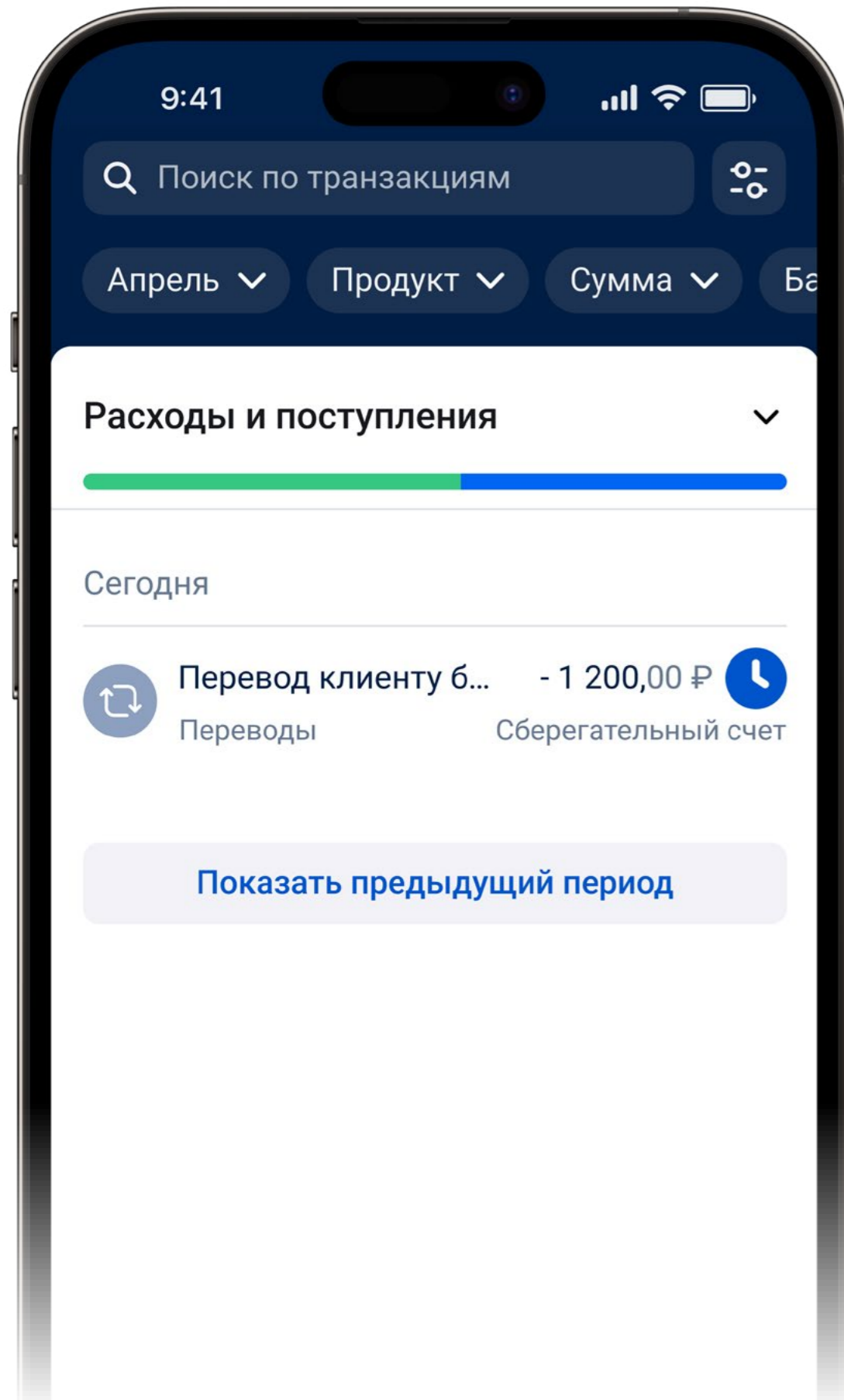
# Обновляем с timestamp



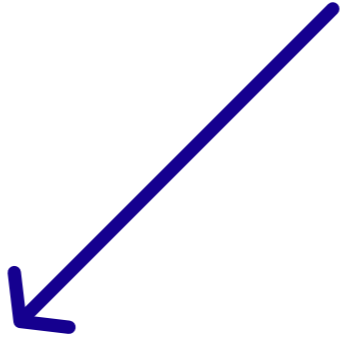
SSTable 1	
created	10:00

SSTable 2	
done	10:05

# Обновляем с timestamp



Агрегация  
status = **done**

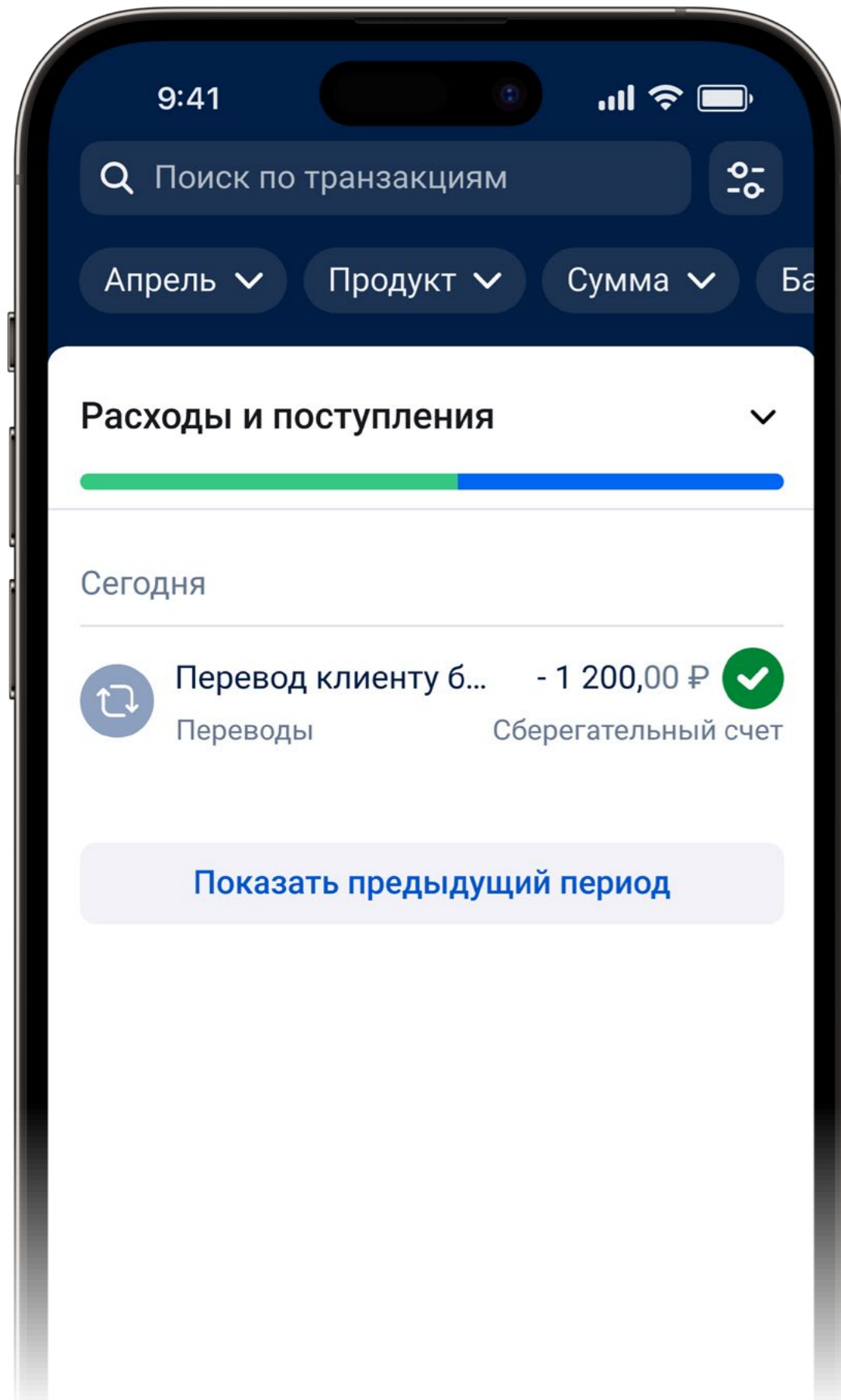


SSTable 1	
created	10:00

SSTable 2	
done	<b>10:05</b>



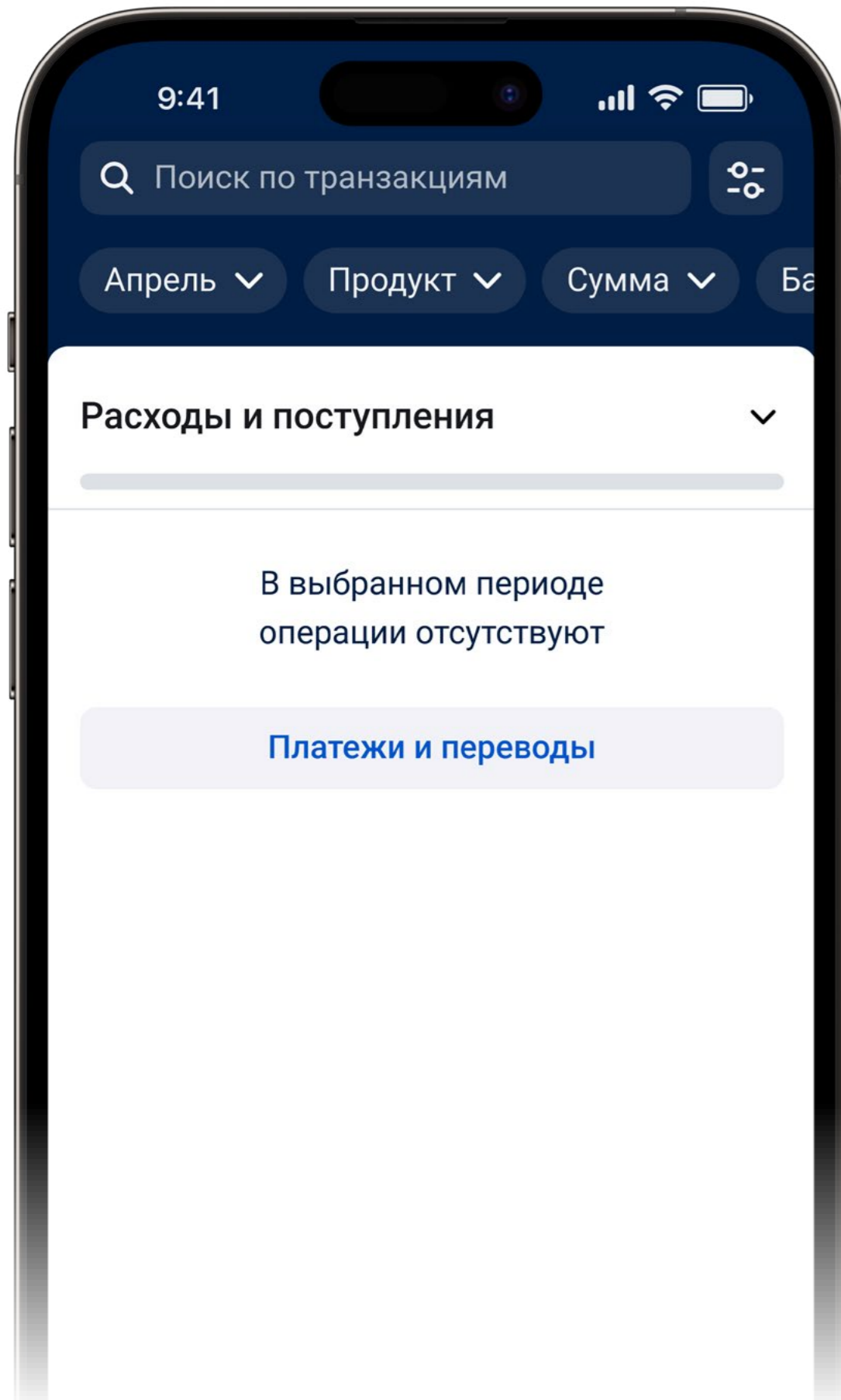
# Обновляем с timestamp



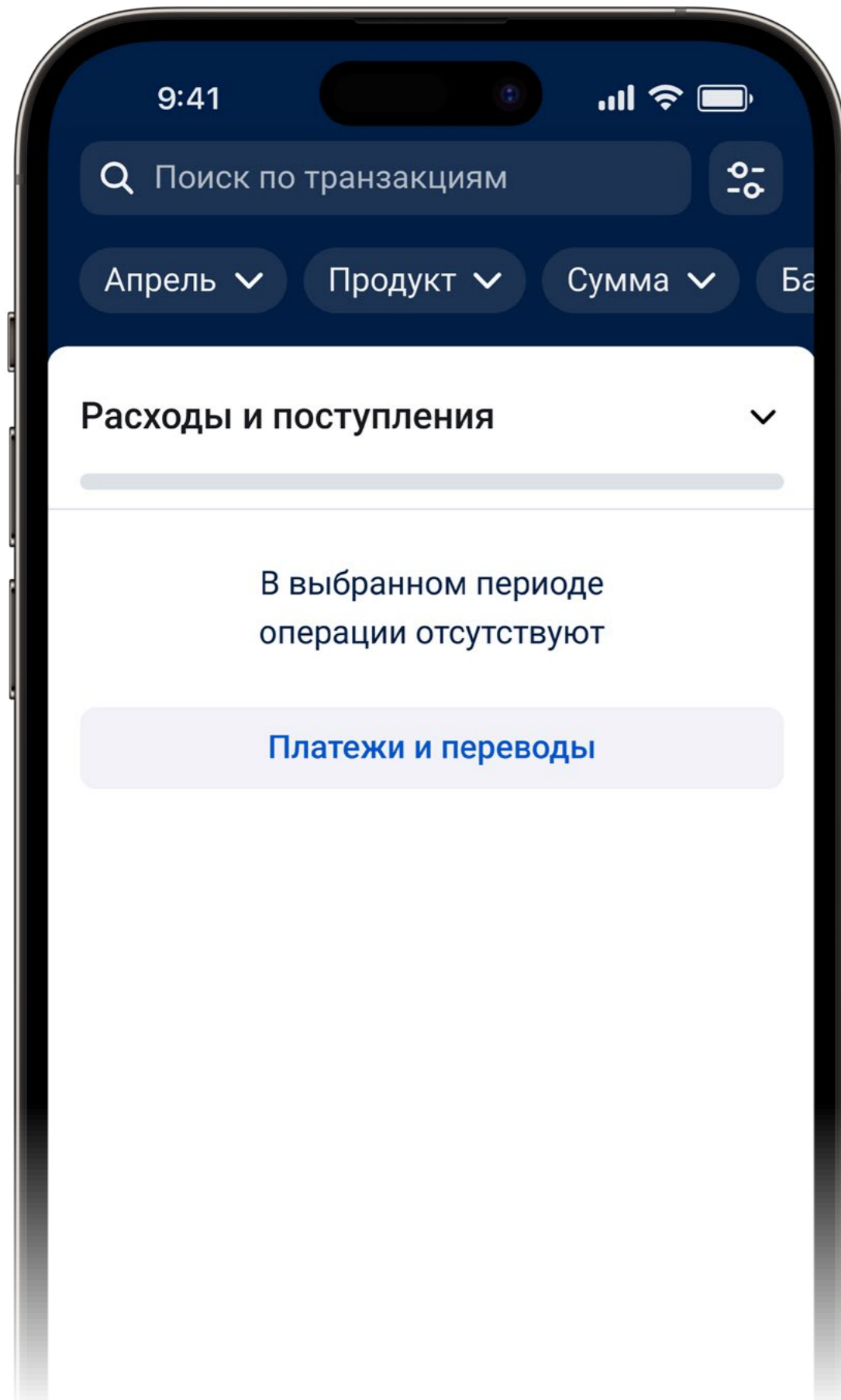
SSTable 1	
created	10:00

SSTable 2	
done	10:05

# Обновляем с timestamp

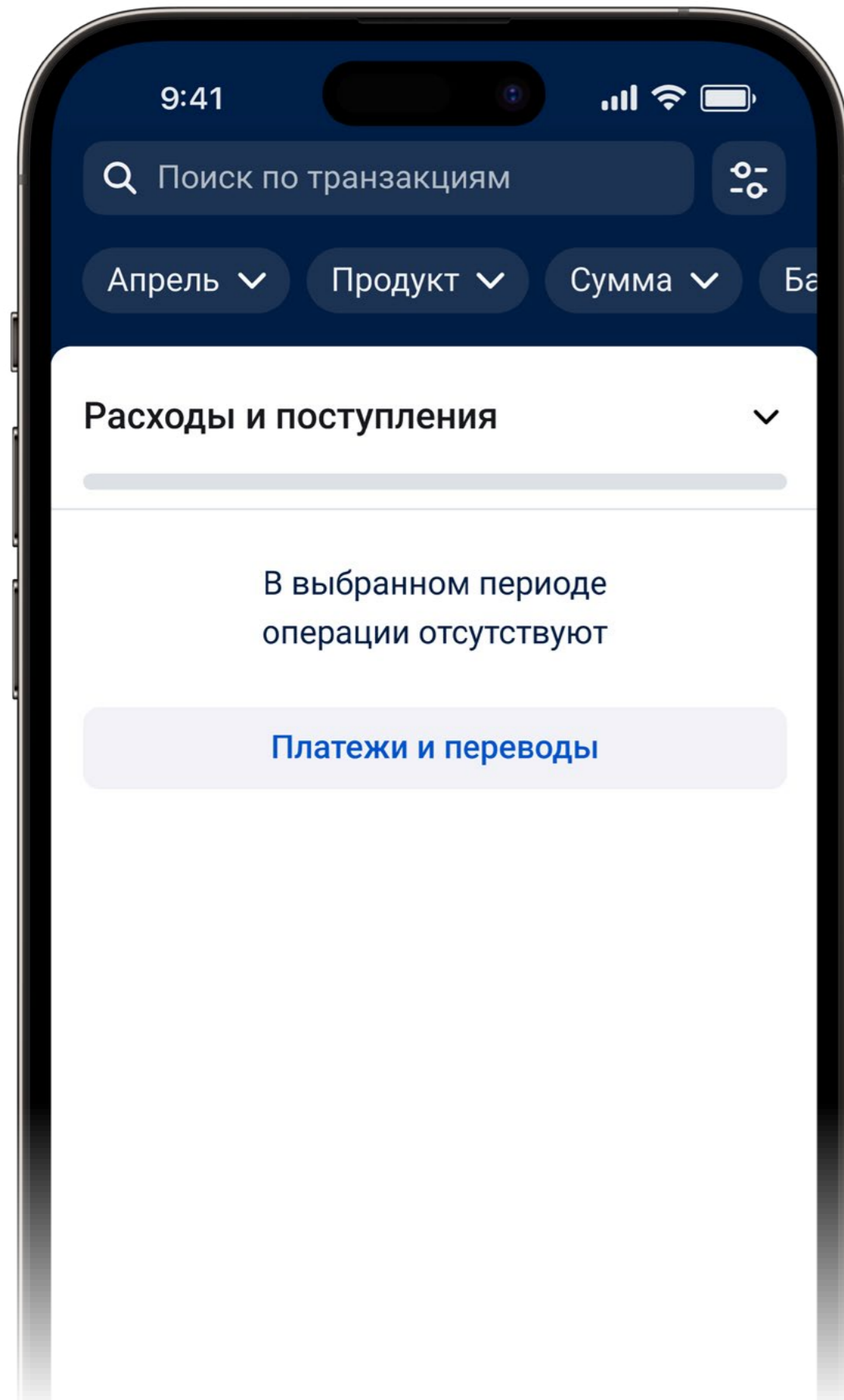


# Обновляем с timestamp

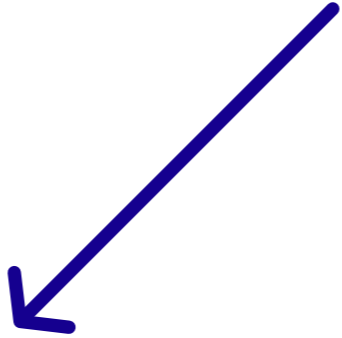


SSTable 1	
done	10:05

# Обновляем с timestamp

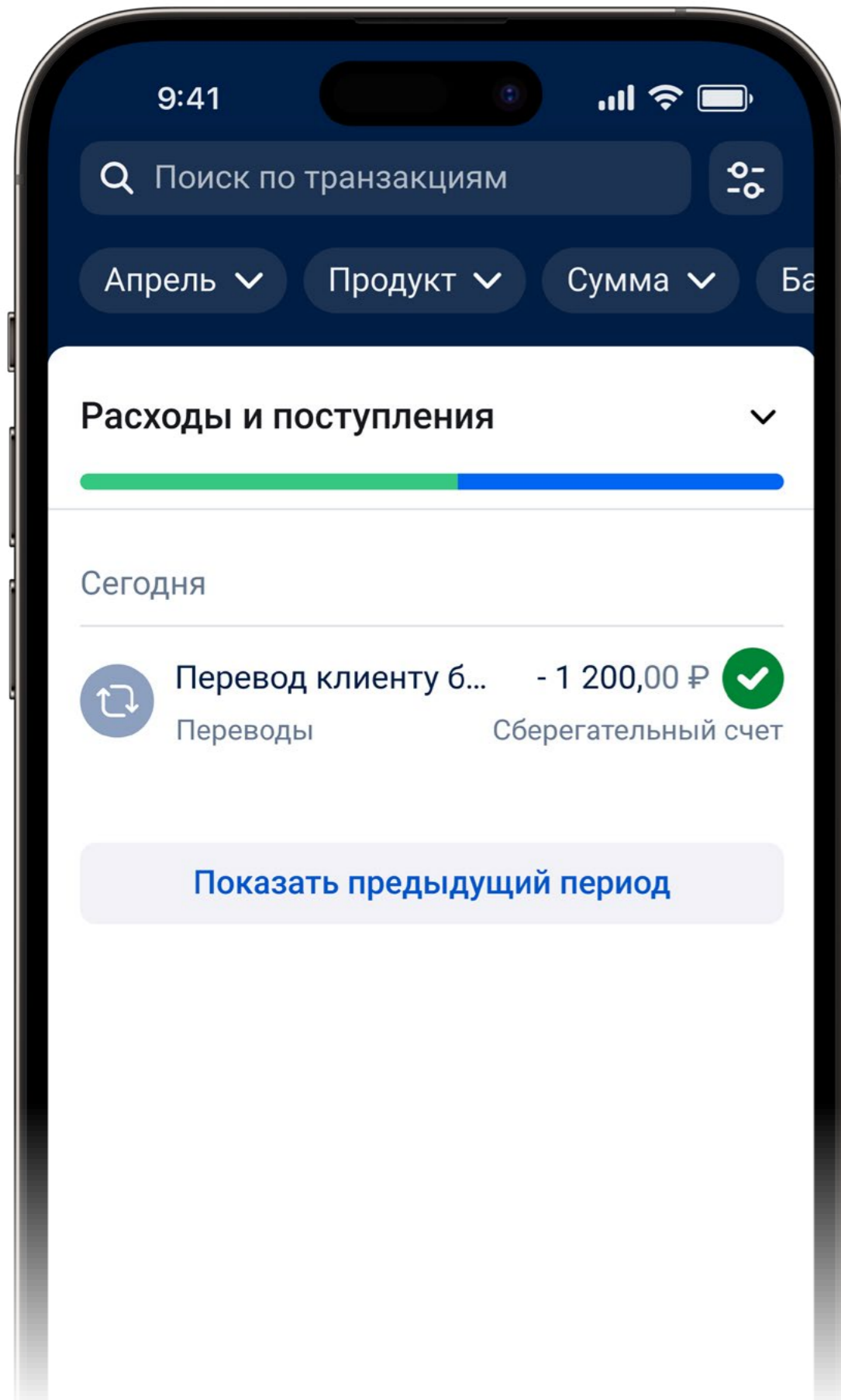


Агрегация  
status = **done**

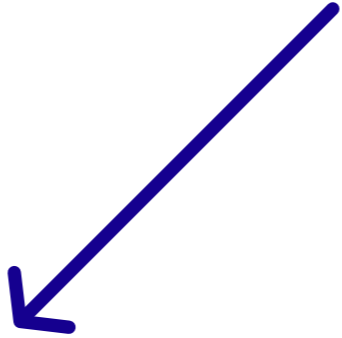


SSTable 1	
done	10:05

# Обновляем с timestamp

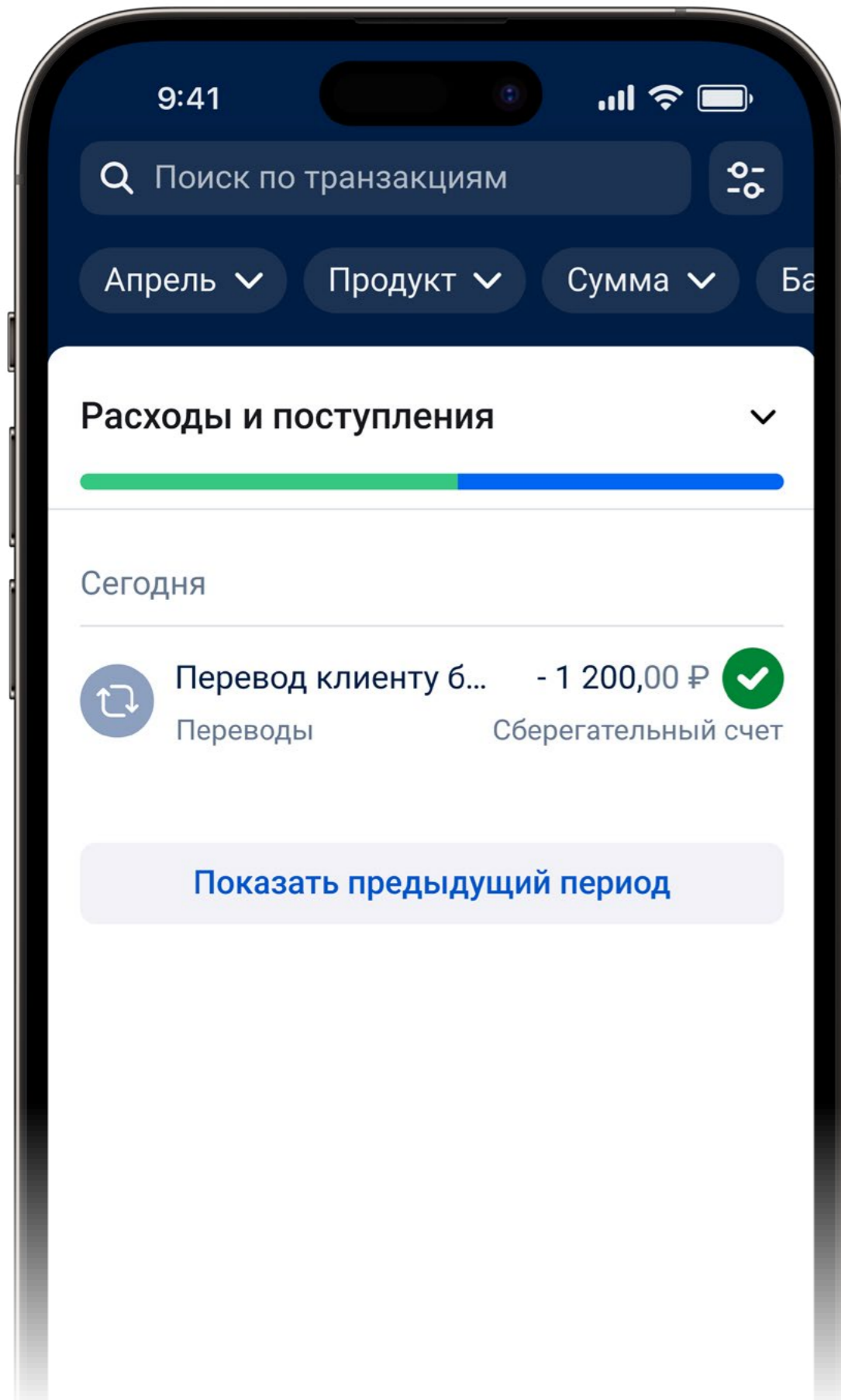


Агрегация  
status = **done**



SSTable 1	
done	10:05

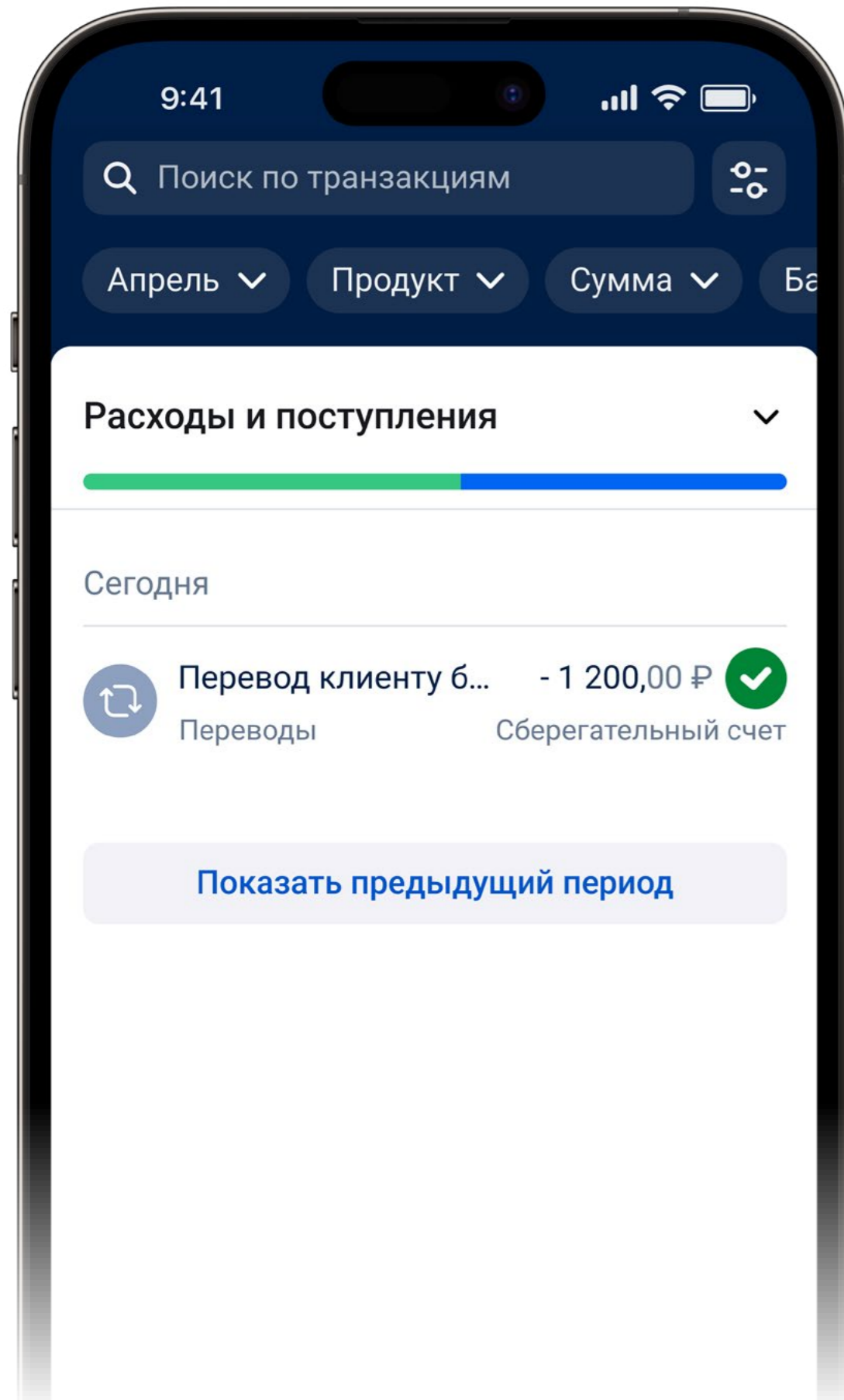
# Обновляем с timestamp



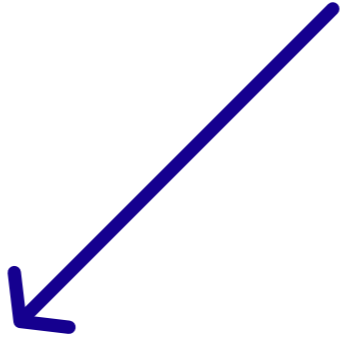
SSTable 1	
done	10:05

SSTable 2	
created	10:00

# Обновляем с timestamp



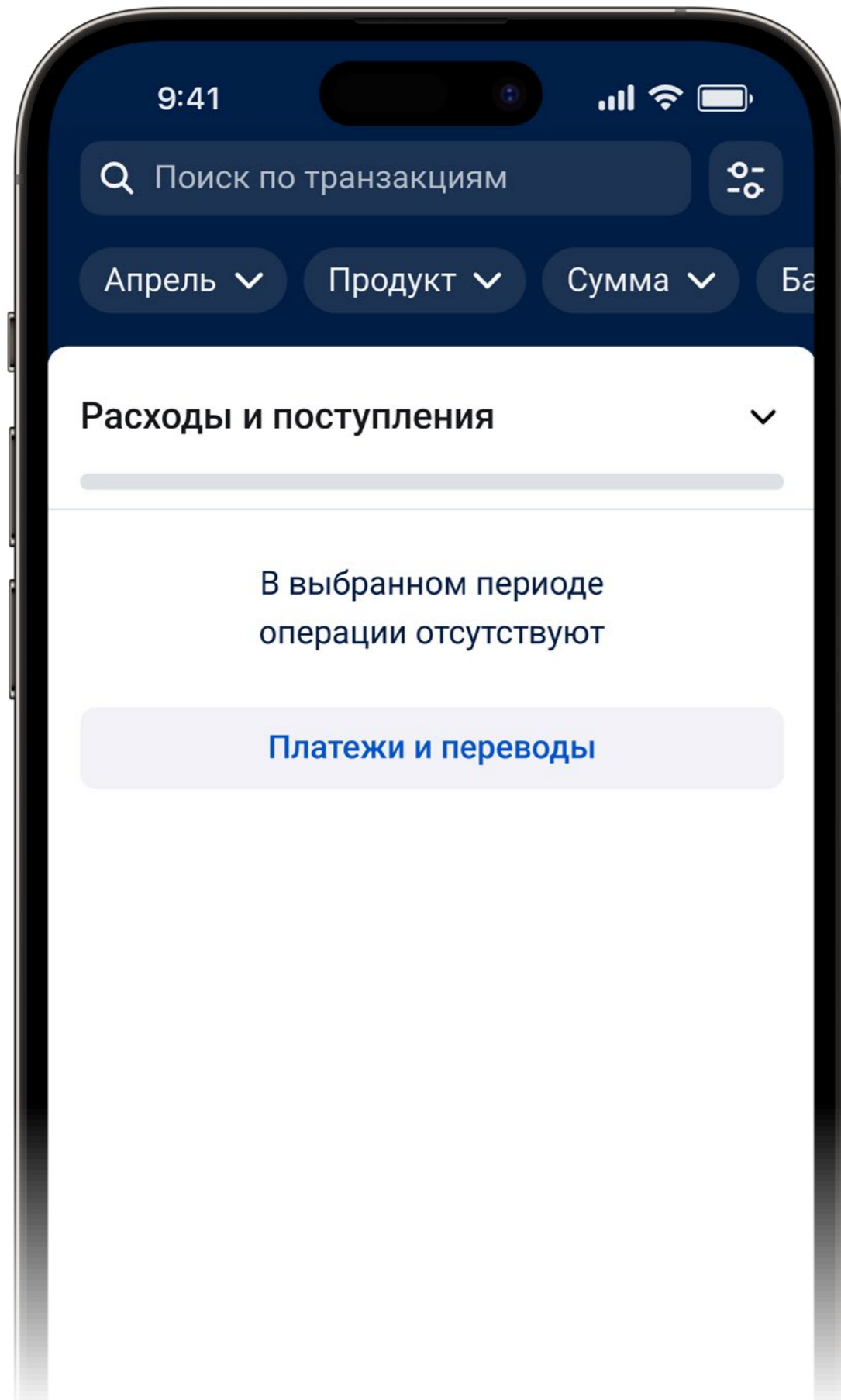
Агрегация  
status = **done**



SSTable 1	
done	10:05

SSTable 2	
created	10:00

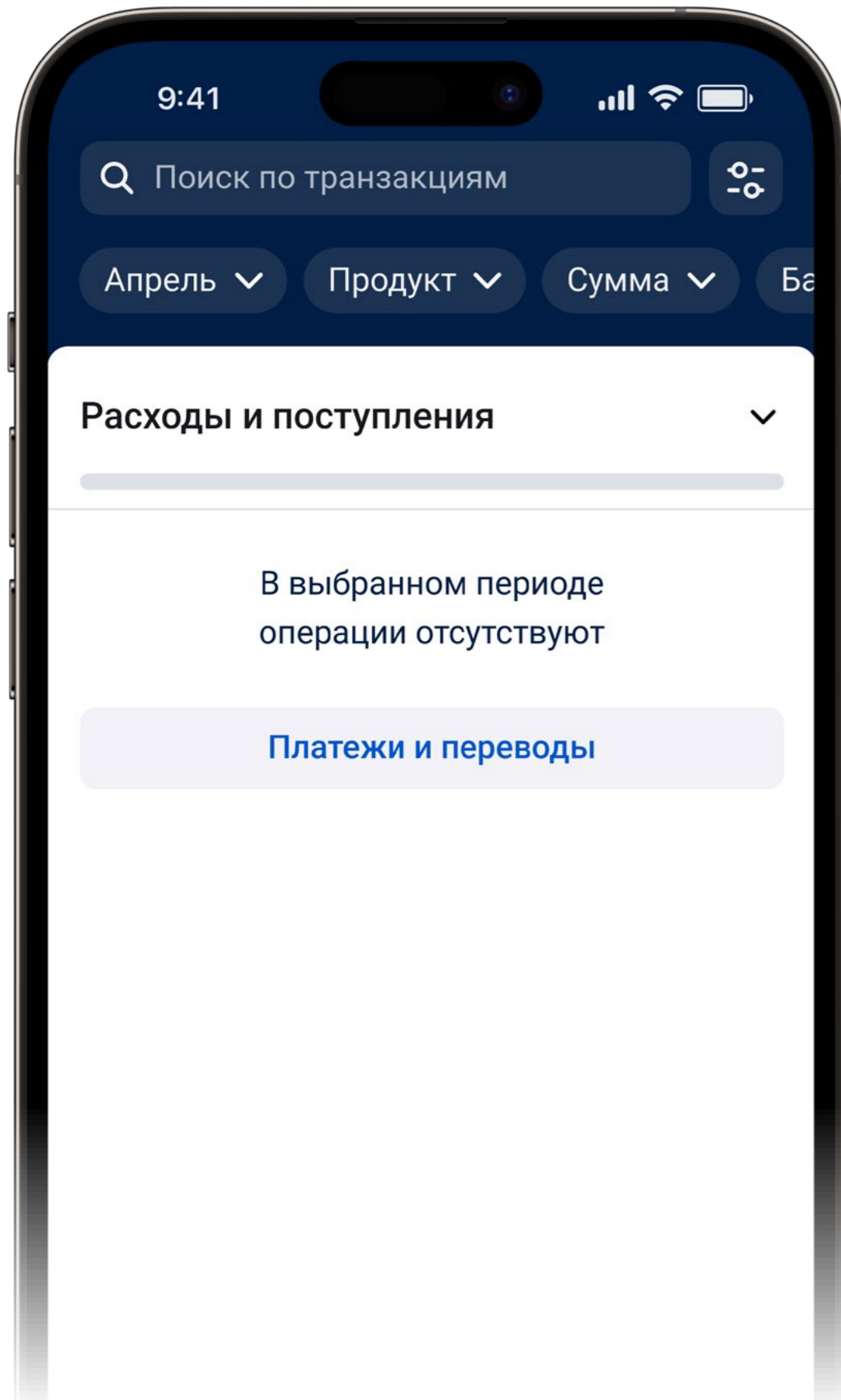
# Удаляем ошибочную запись



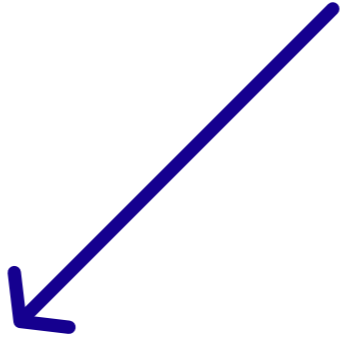
SSTable 1	
error	10:00



# Удаляем ошибочную запись

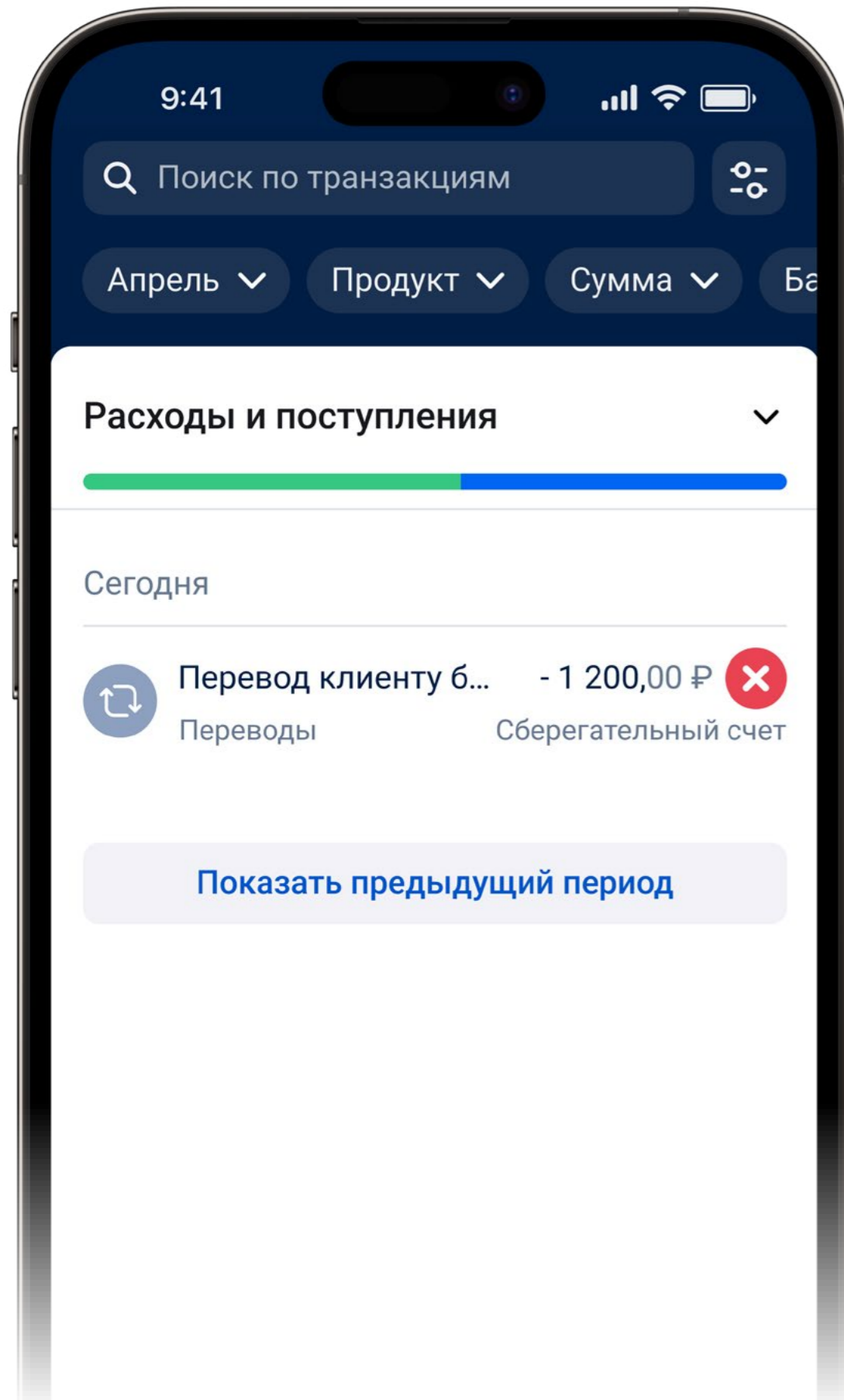


Агрегация  
status = **error**

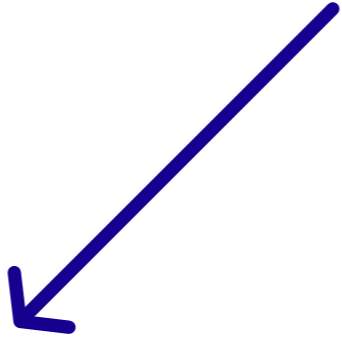


SSTable 1	
error	10:00

# Удаляем ошибочную запись

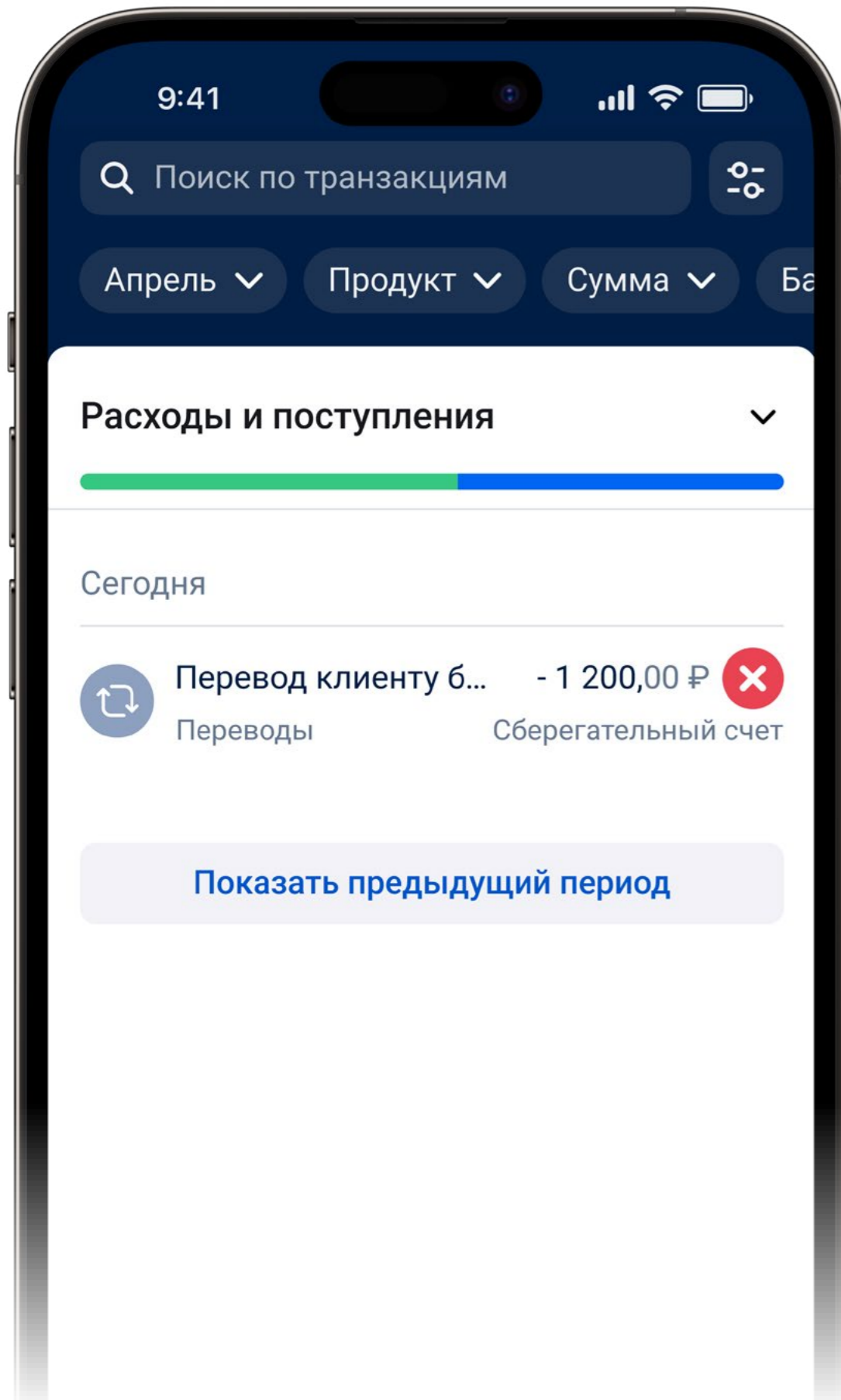


Агрегация  
status = **error**



SSTable 1	
error	10:00

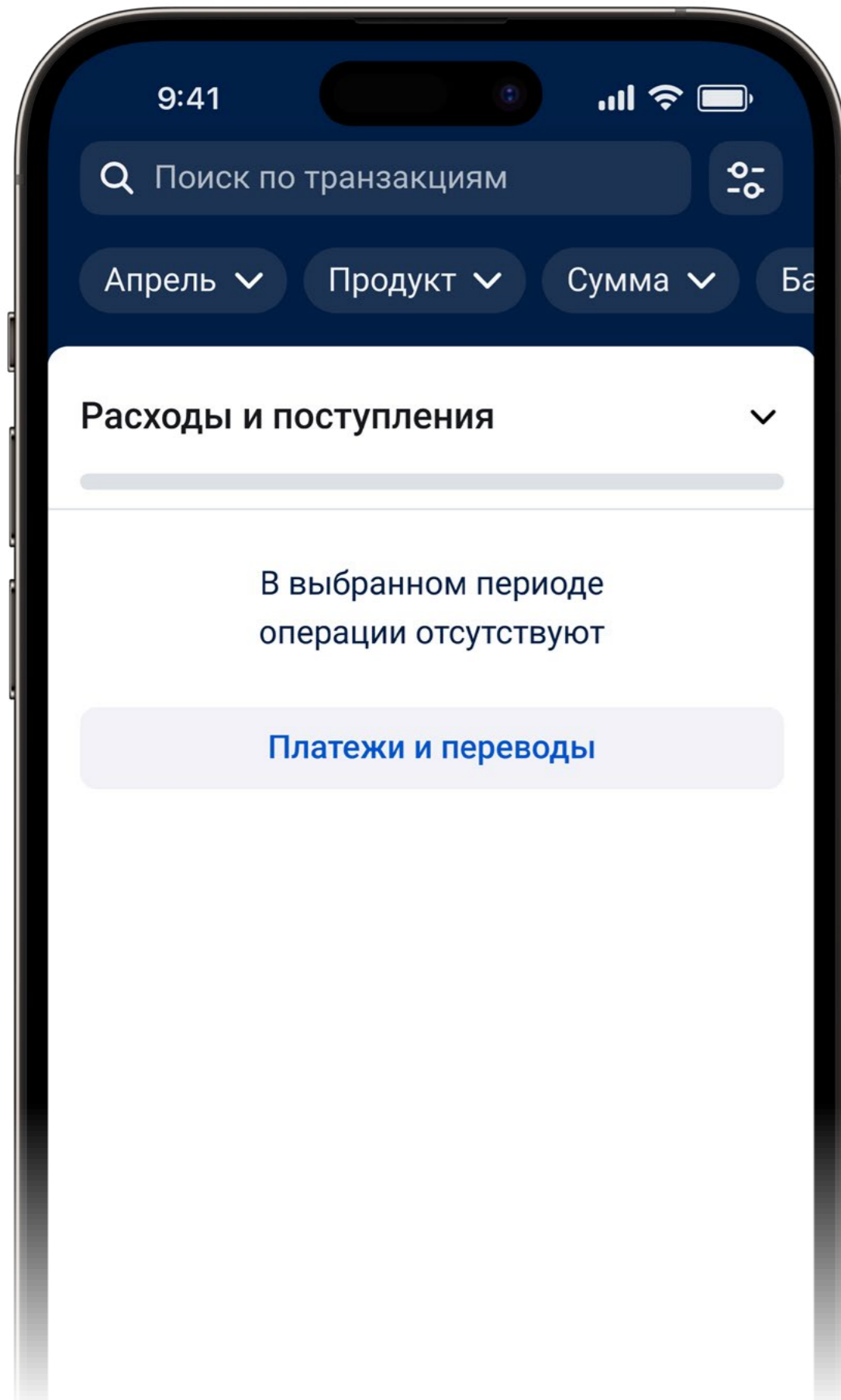
# Удаляем ошибочную запись



SSTable 1	
error	10:00

SSTable 2	
tombstone	11:00

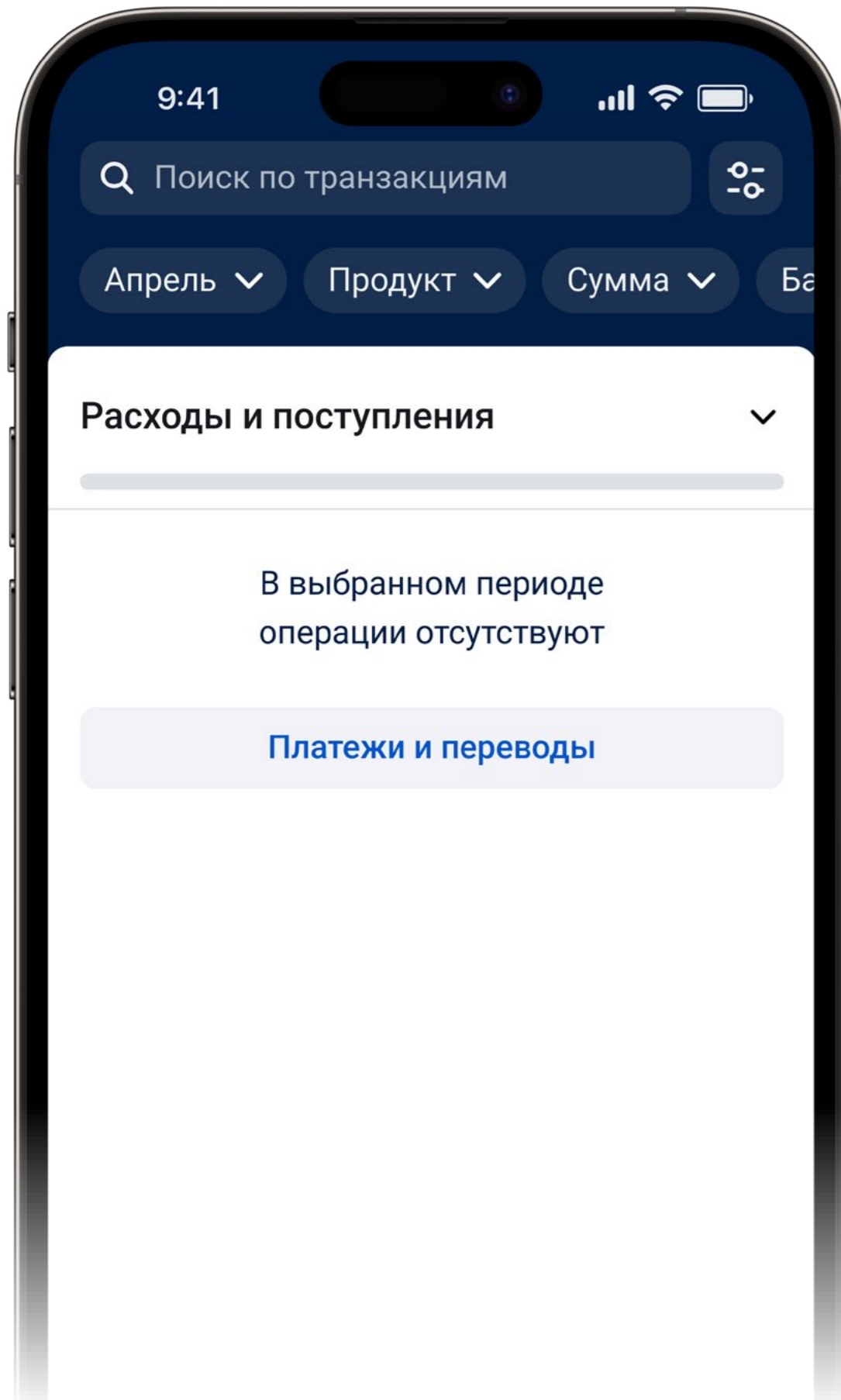
# Удаляем ошибочную запись



SSTable 1	
error	10:00

SSTable 2	
tombstone	11:00

# Удаляем ошибочную запись

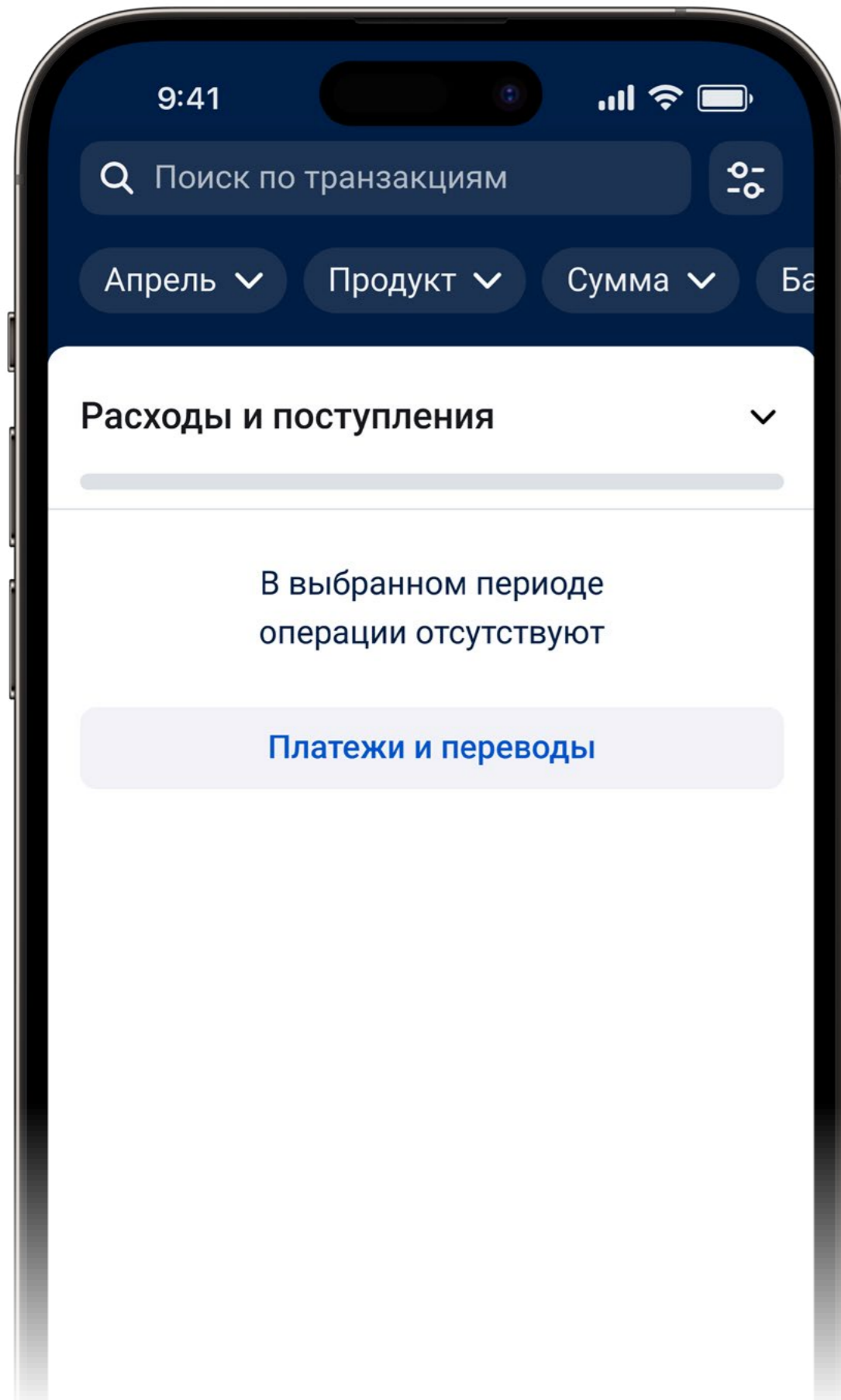


SSTable 1	
error	10:00

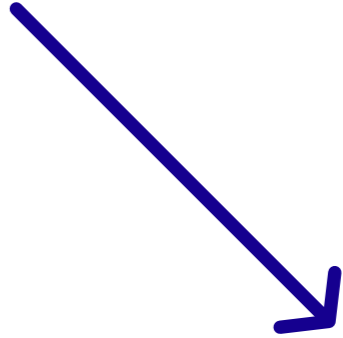
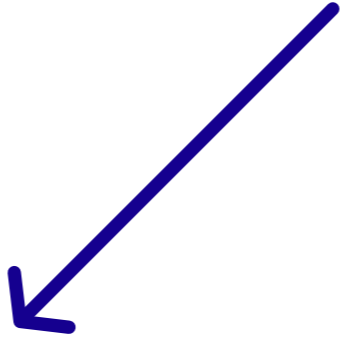
SSTable 2	
tombstone	11:00

SSTable 3	
created	10:05

# Удаляем ошибочную запись



Агрегация  
status = **null**

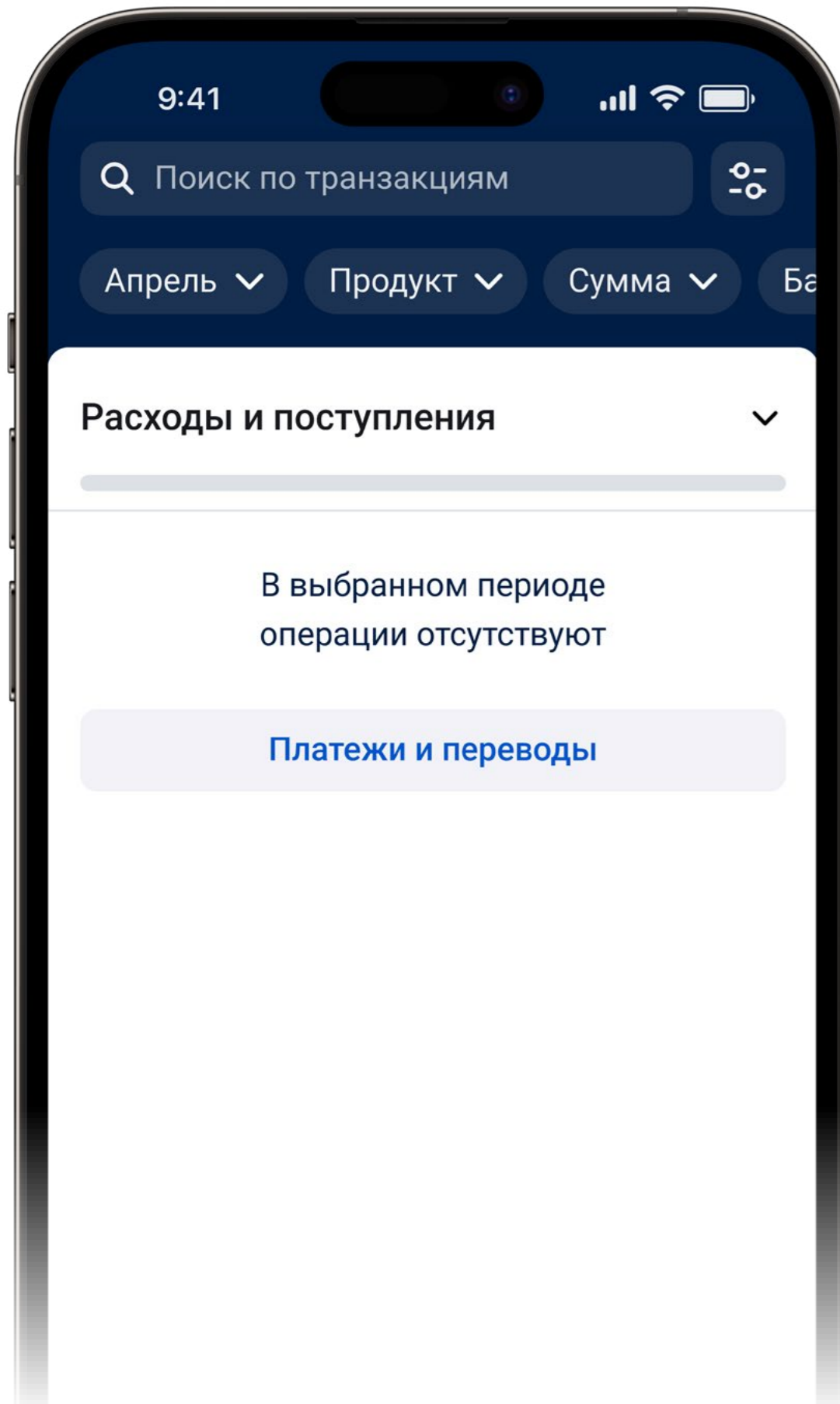


SSTable 1	
error	10:00

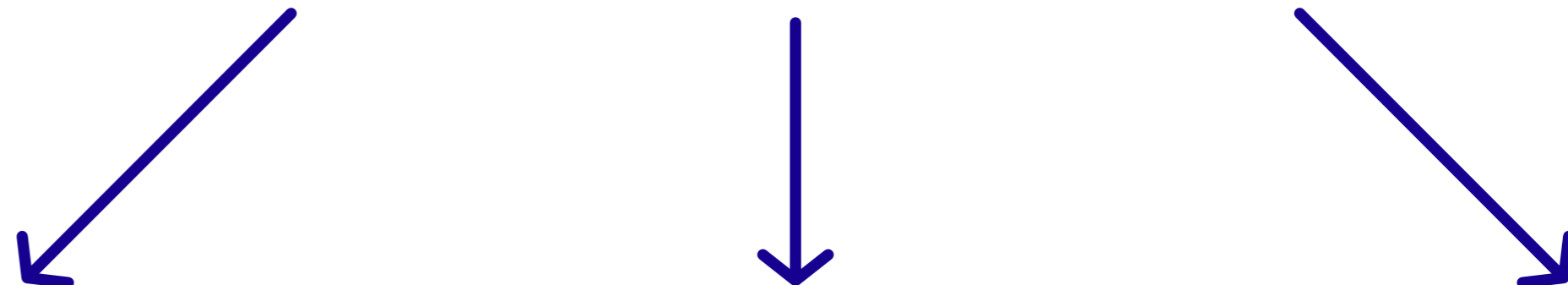
SSTable 2	
tombstone	<b>11:00</b>

SSTable 3	
created	10:05

# Удаляем ошибочную запись



Агрегация  
status = **null**



SSTable 1	
error	10:00

SSTable 2	
tombstone	<b>11:00</b>

SSTable 3	
created	10:05

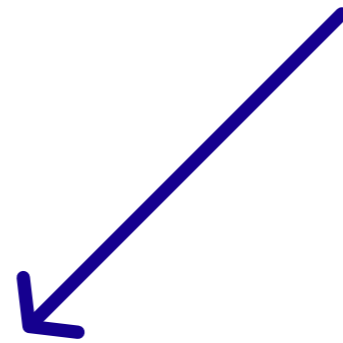
<https://medium.com/@jeeyoungk/why-i-love-databases-1d4cc433685f>

“You may be able to delete a future write with an old **doomstone**. This is colloquially referred as a doomstone. Hilarious as it is, it’s a real problem.”

# Удаляем ошибочную запись правильно



Агрегация  
status = **error**



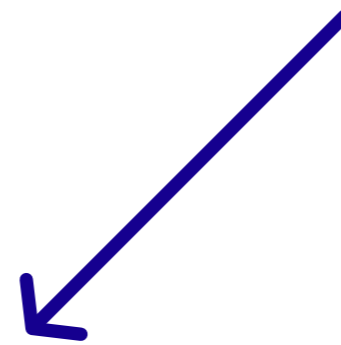
SSTable 1	
error	<b>10:00</b>



# Удаляем ошибочную запись правильно



Агрегация  
status = **null**



```
DELETE status
FROM conf.timeline
USING TIMESTAMP
1711868400001
WHERE client_id = 1
AND year = 2024
AND month = 4
AND date = '2024-04-25'
AND op_id = 333
```

SSTable 1	
error	10:00

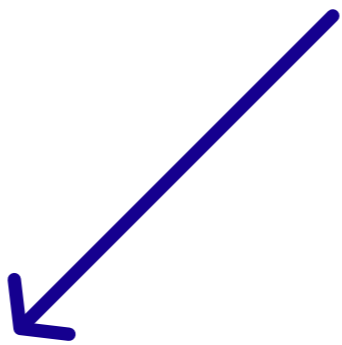
SSTable 2	
tombstone	<b>10:00</b>

При совпадении временных меток приоритет отдаётся tombstone

# Как узнать timestamp



Агрегация  
status = **error**

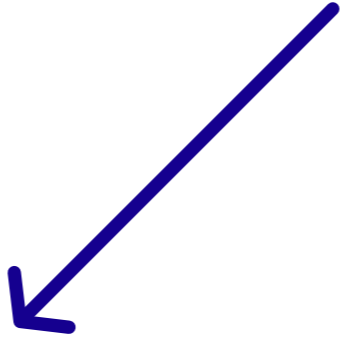


SSTable 1	
error	???

# Как узнать timestamp



Агрегация  
status = **error**



```
SELECT status,  
WRITETIME(status)  
FROM conf.timeline  
WHERE client_id = 1  
AND year = 2024  
AND month = 4  
AND date = '2024-04-25'  
AND op_id = 333
```

SSTable 1	
error	???

status	writetime(status)
error	1713456291673174

# Как узнать timestamp

А если уже удалили?

```
SELECT status,  
WRITETIME(status)  
FROM conf.timeline  
WHERE client_id = 1  
AND year = 2024  
AND month = 4  
AND date = '2024-04-25'  
AND op_id = 333
```

Агрегация  
status = **null**

SSTable 1	
error	10:00

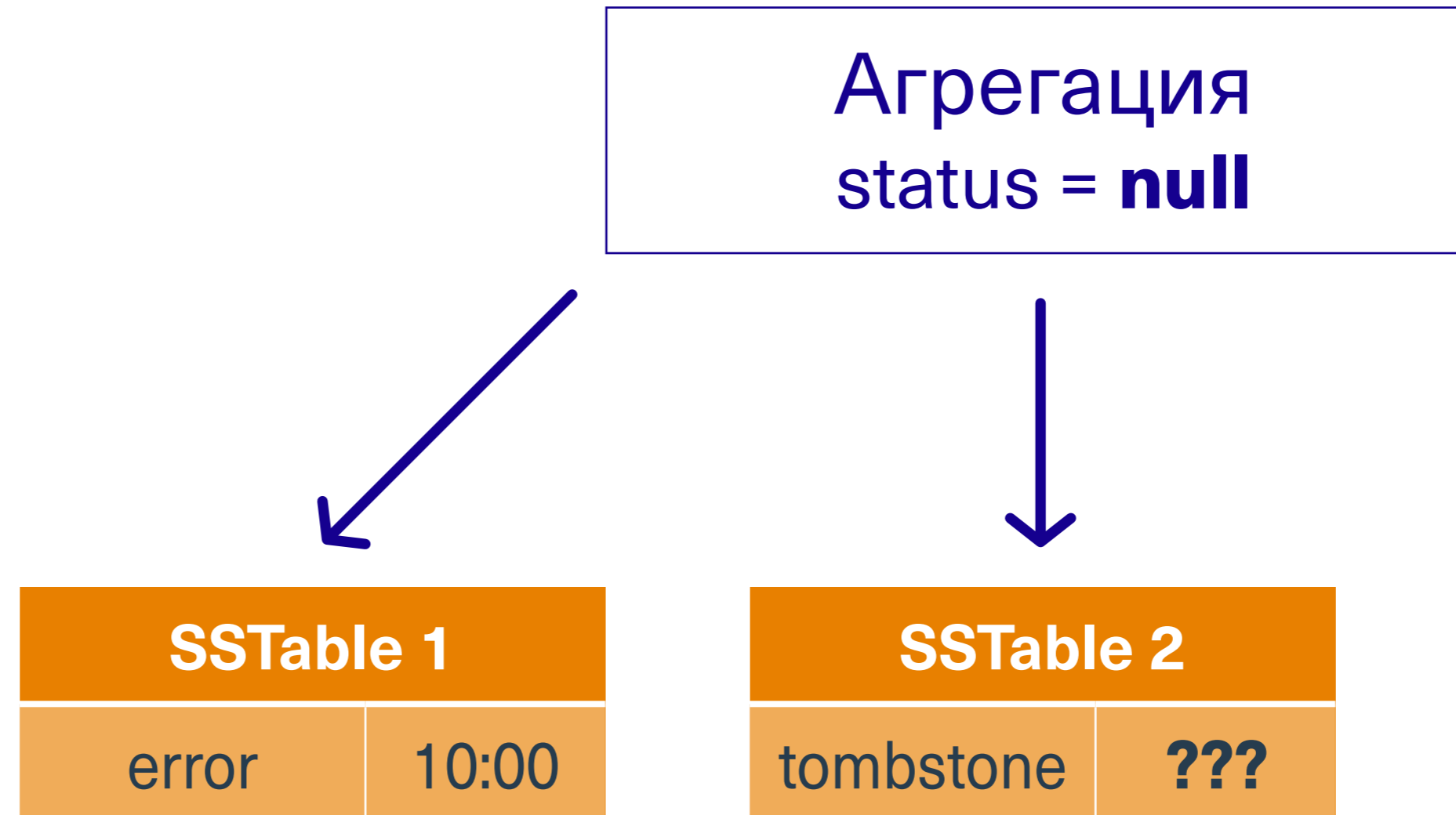
SSTable 2	
tombstone	???

status	writetime(status)
null	null

# Как узнать timestamp



А если уже удалили?



```
sstabledump lb-6-big-Data.db  
[  
  {"key": "1:2024",  
   "cells": [[ "2024-04-25:1:status»,1711962663,1711868400001, "d" ] ] }  
]
```

# Tombstone требуют внимания



## Создаются при

- Удалении строки
- Удалении ячейки
- Удалении диапазона ячеек
- Удалении элемента коллекций
- Вставке значения null

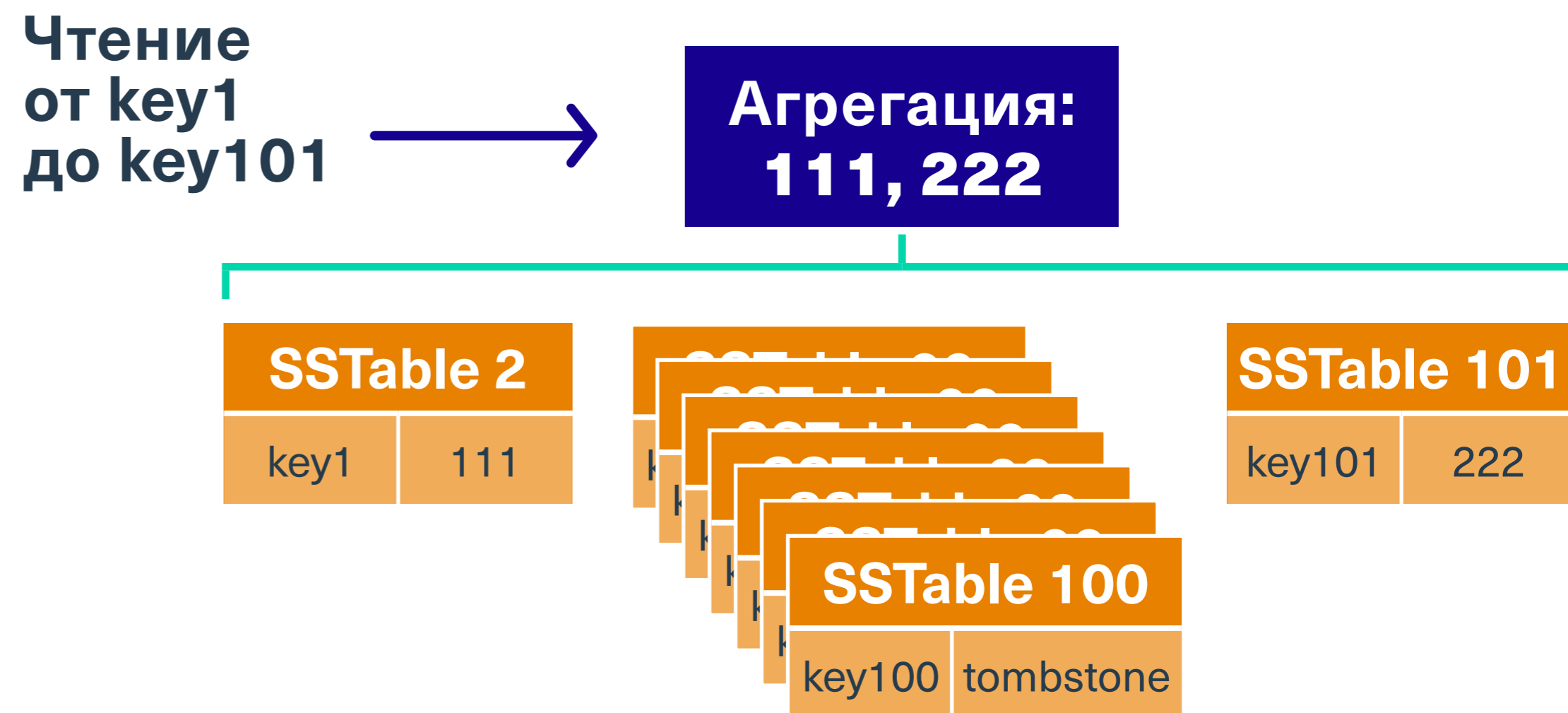
# Tombstone требуют внимания

## Создаются при

- Удалении строки
- Удалении ячейки
- Удалении диапазона ячеек
- Удалении элемента коллекций
- Вставке значения null

## Занимают место на диске

Влияют на IO и загружаются в память при запросах по диапазону ключей и full scan



# Tombstone требуют внимания



tombstone\_warn\_threshold: 1000

tombstone\_failure\_threshold: 100000



# Tombstone требуют внимания

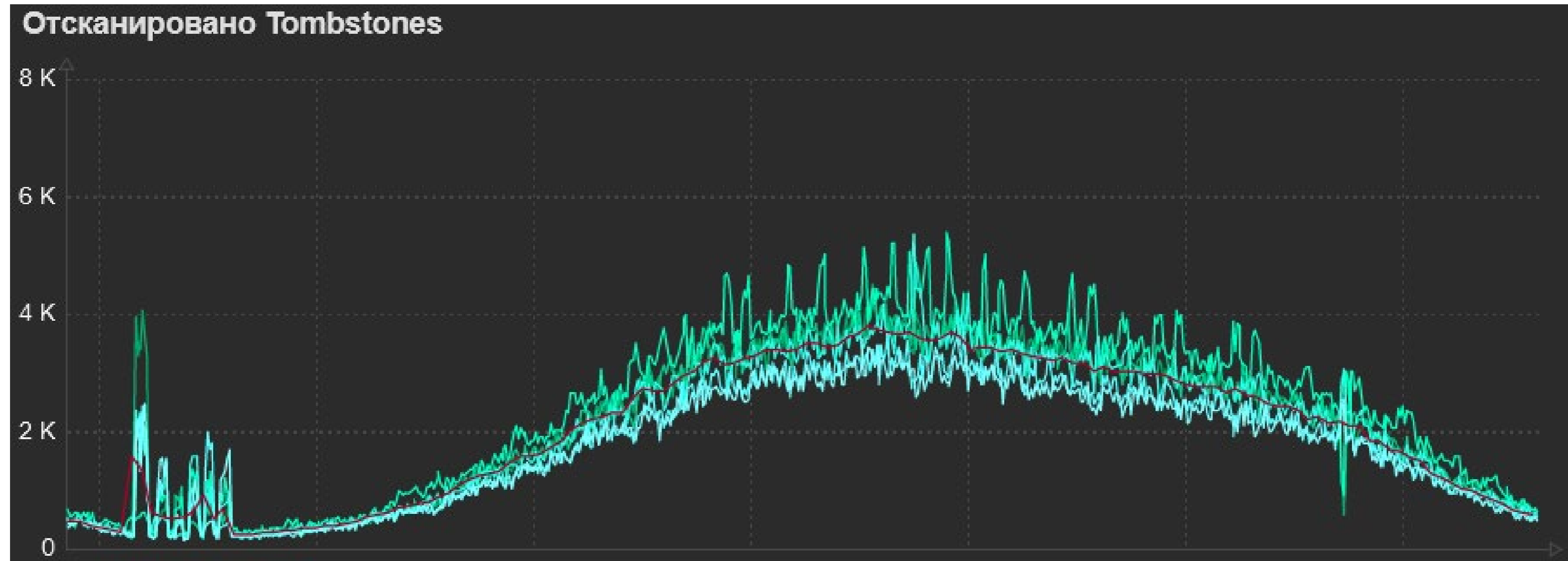


tombstone\_warn\_threshold: 1000

tombstone\_failure\_threshold: 100000

```
nodetool tablestats conf.timeline
...
Average tombstones per slice (last five minutes): 1.0
Maximum tombstones per slice (last five minutes): 1
```

# Tombstone требуют внимания

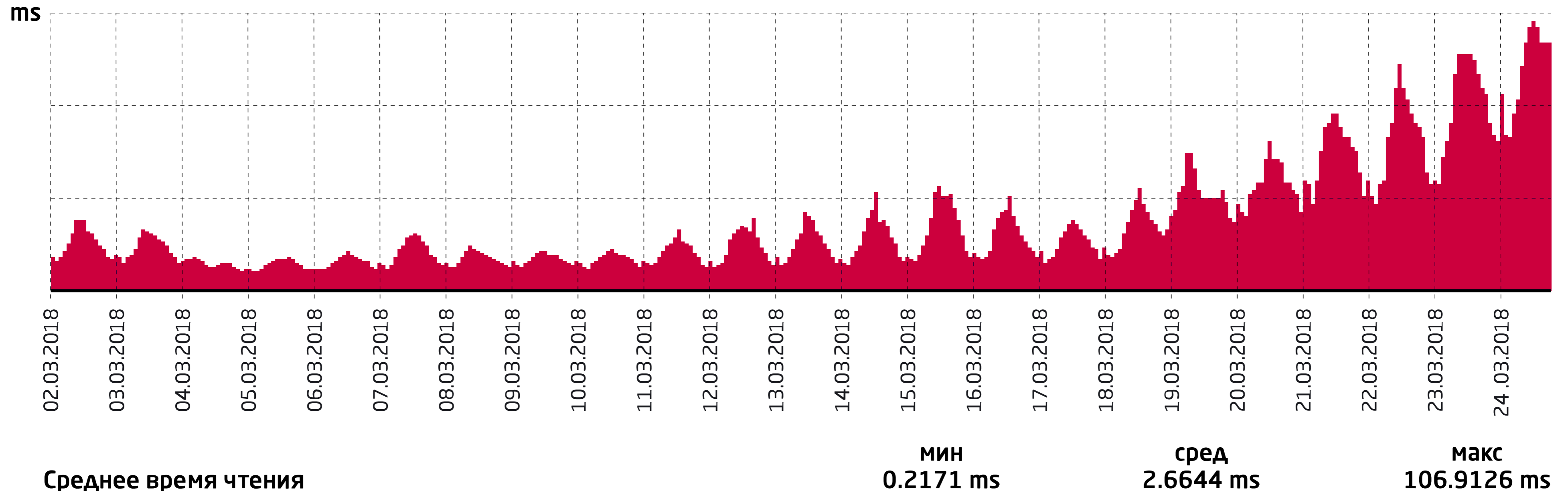


TombstoneScannedHistogram — количество tombstone обработанных при чтении из таблицы

# Новая проблема



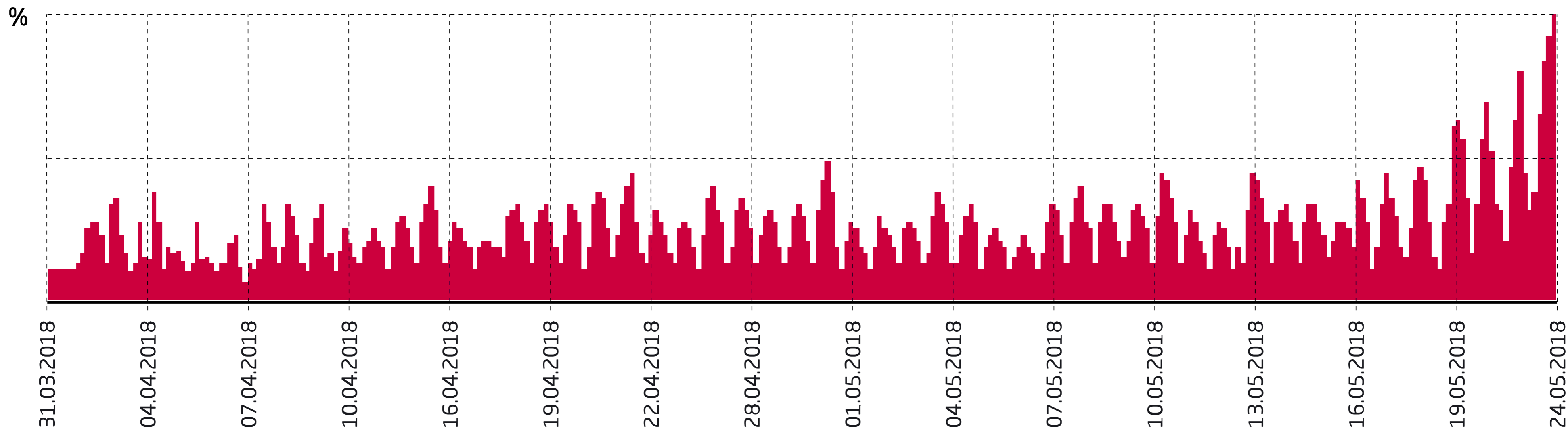
## ■ Время чтения увеличивается



# Новая проблема



- Время чтения увеличивается
- А ещё растёт iowait



CPU iowait time

МИН  
0.01668 %

сред  
6.3244 %

макс  
48.06 %

# Новая проблема



- Время чтения увеличивается
- А ещё растёт iowait
- Упёрлись в ограничения дисков



# Новая проблема



- Время чтения увеличивается
- А ещё растёт `iowait`
- Упёрлись в ограничения дисков
- Времени на обновление железа нет
- Ожидается пиковая нагрузка



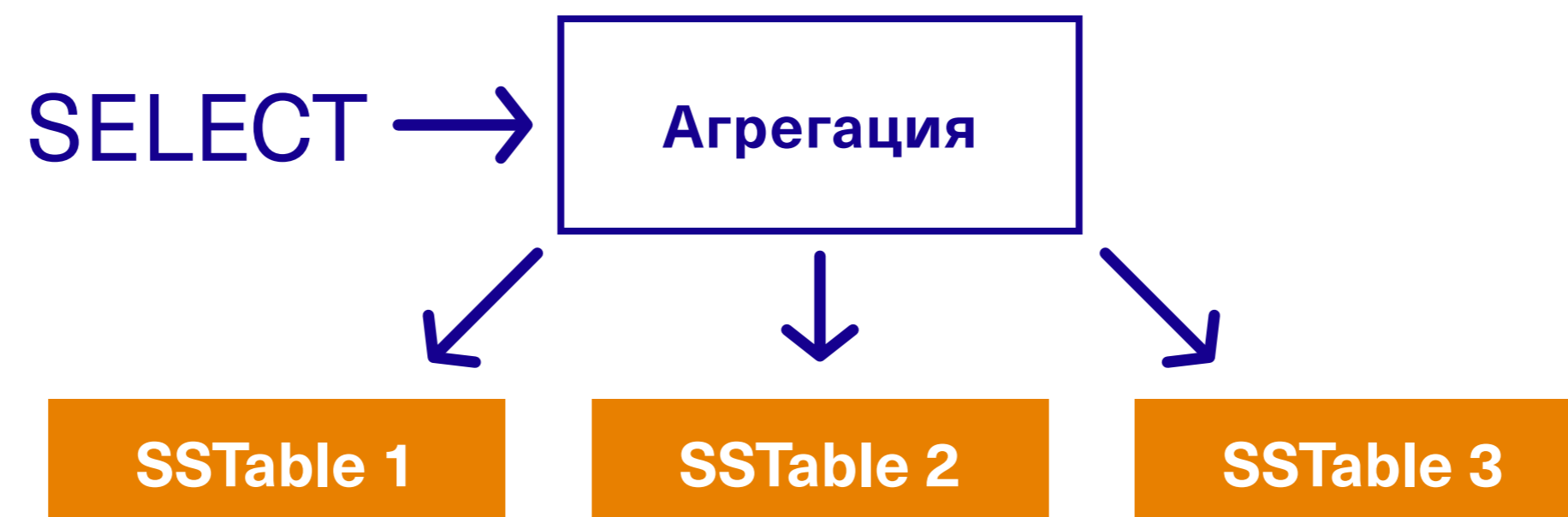
# Как разгрузить диск?



- Full scan-ов нет
- Большого количества tombstone нет
- Есть кейсы когда мы много обновляем и читаем

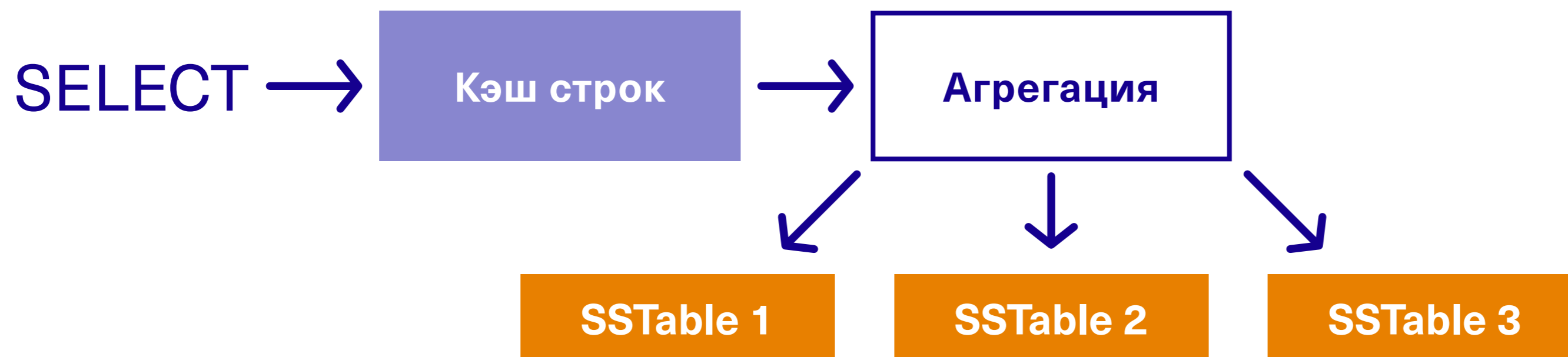
# Как разгрузить диск?

- Full scan-ов нет
- Большого количества tombstone нет
- Есть кейсы когда мы много обновляем и читаем

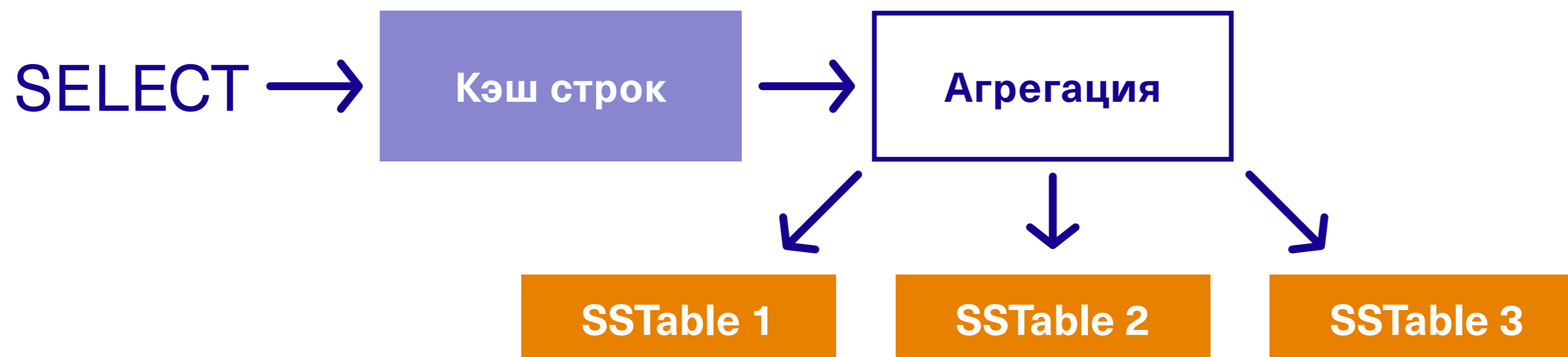




# Кэш строк



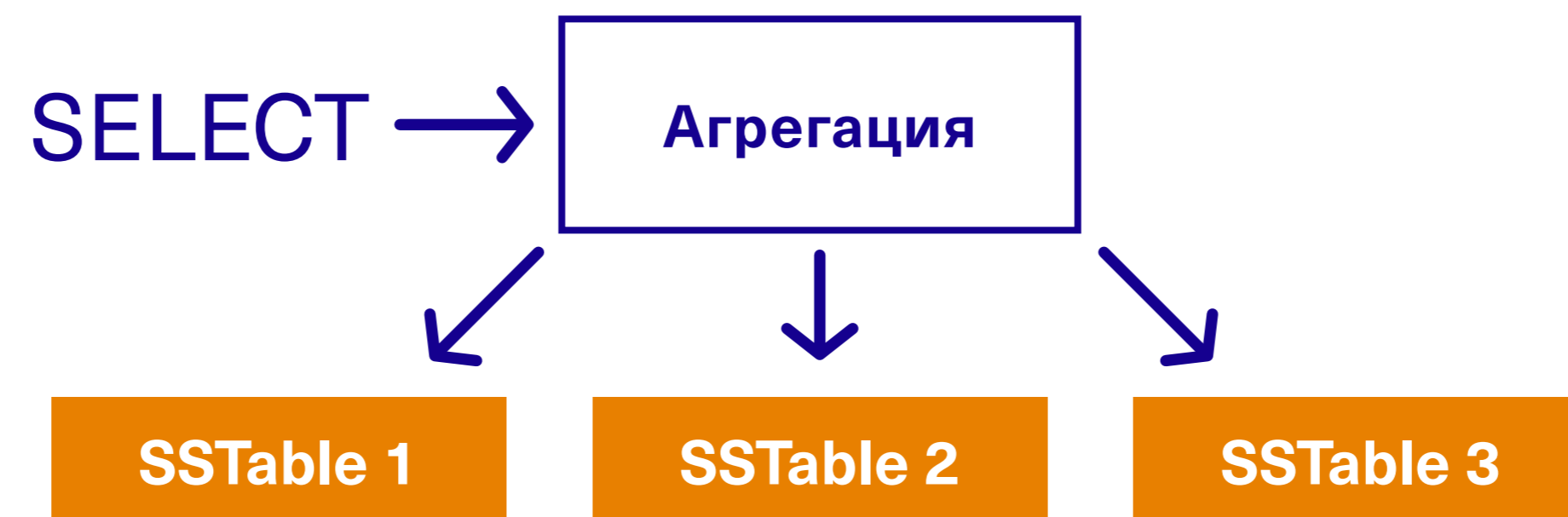
**rows\_per\_partition** — сколько первых строк партиции кэшировать



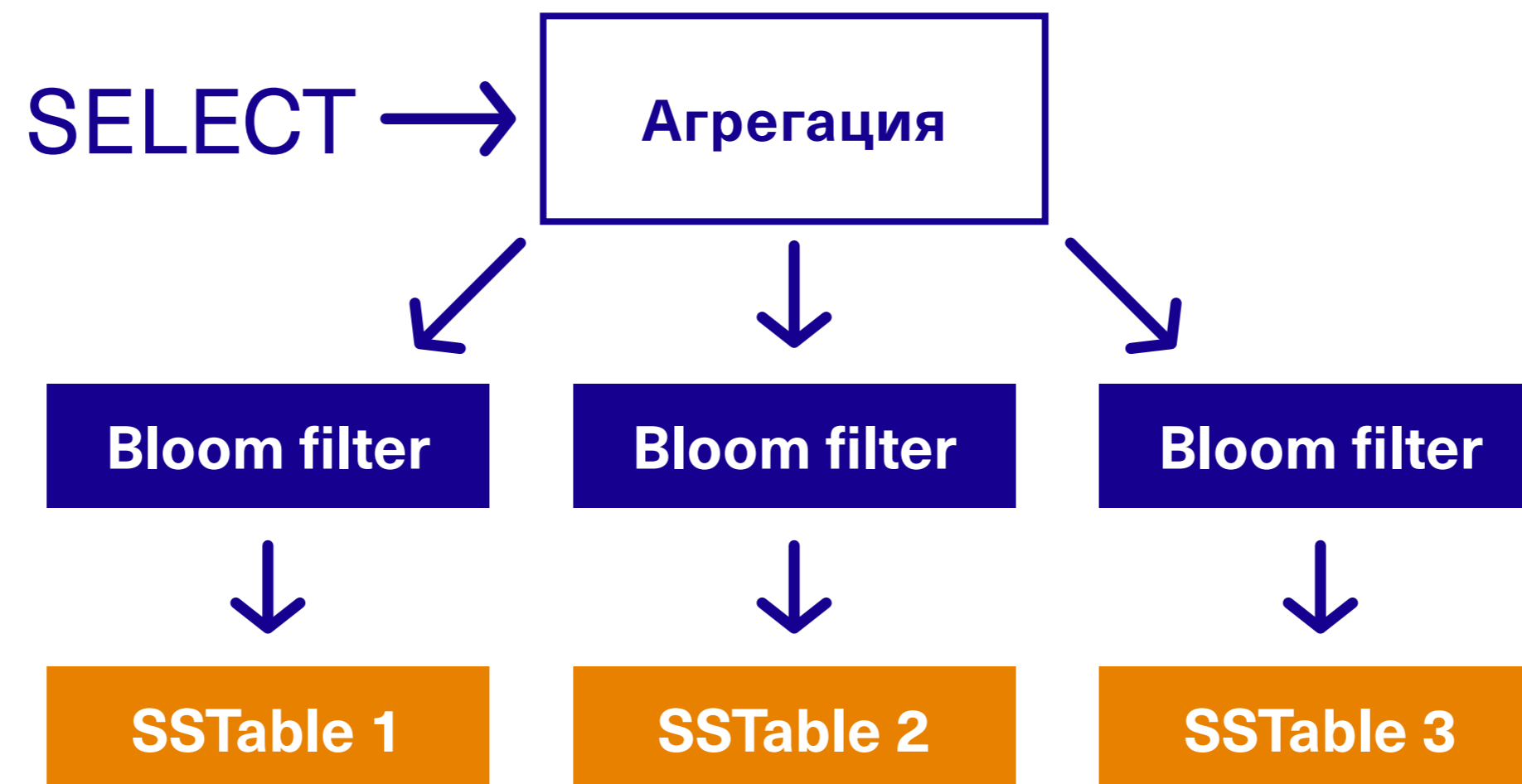
**rows\_per\_partition** — сколько первых строк партиции кэшировать

- Кэширует только первые строки
- Сбрасывается при записи
- Utilizing appropriate OS page cache will result in better performance than using row caching

# Bloom filter

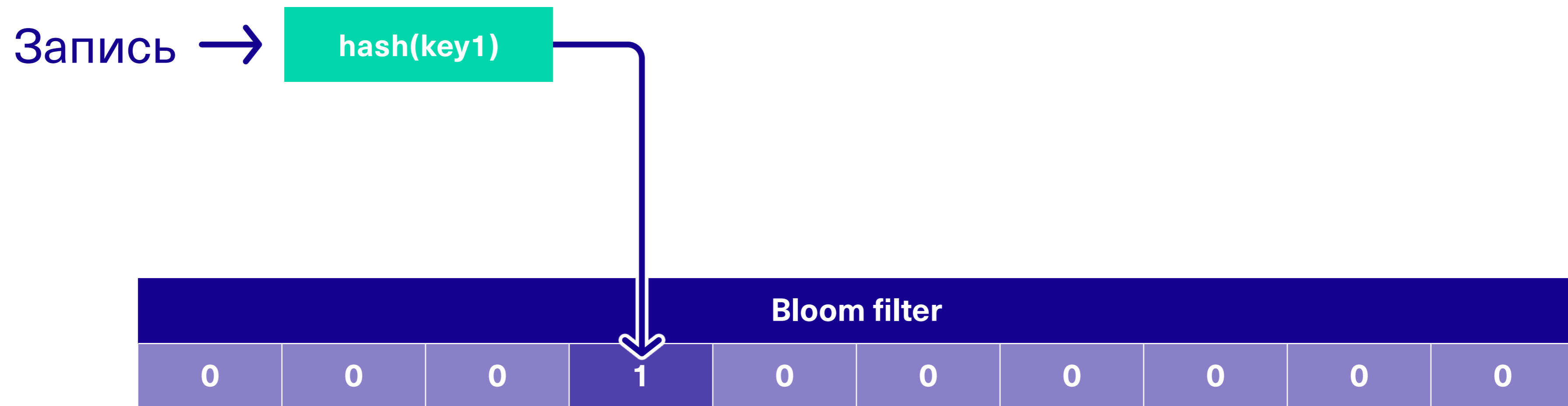


# Bloom filter

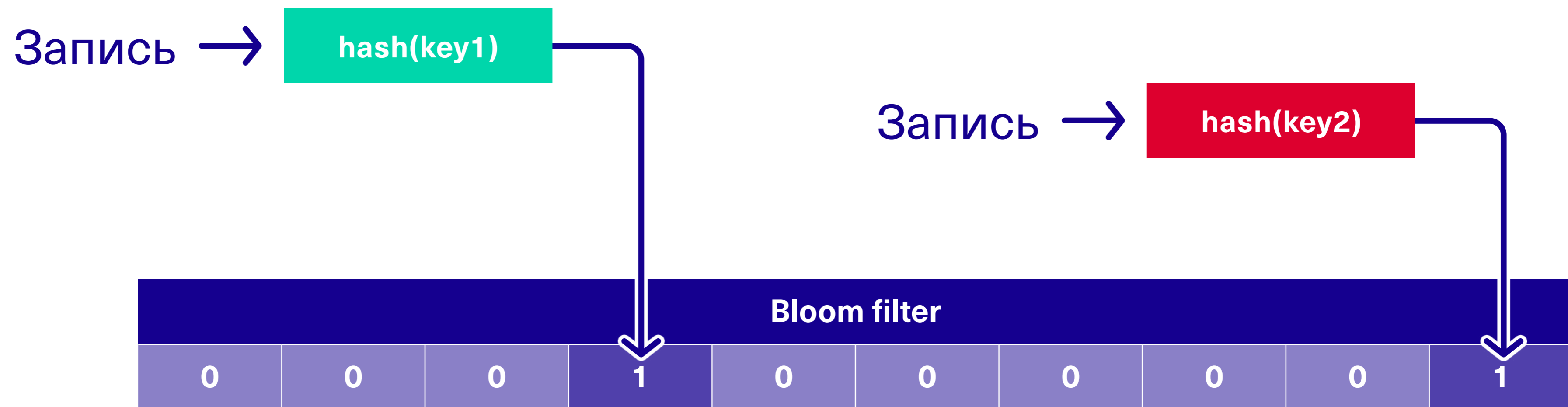


**Bloom filter говорит «ключа точно нет» или «ключ возможно присутствует»**

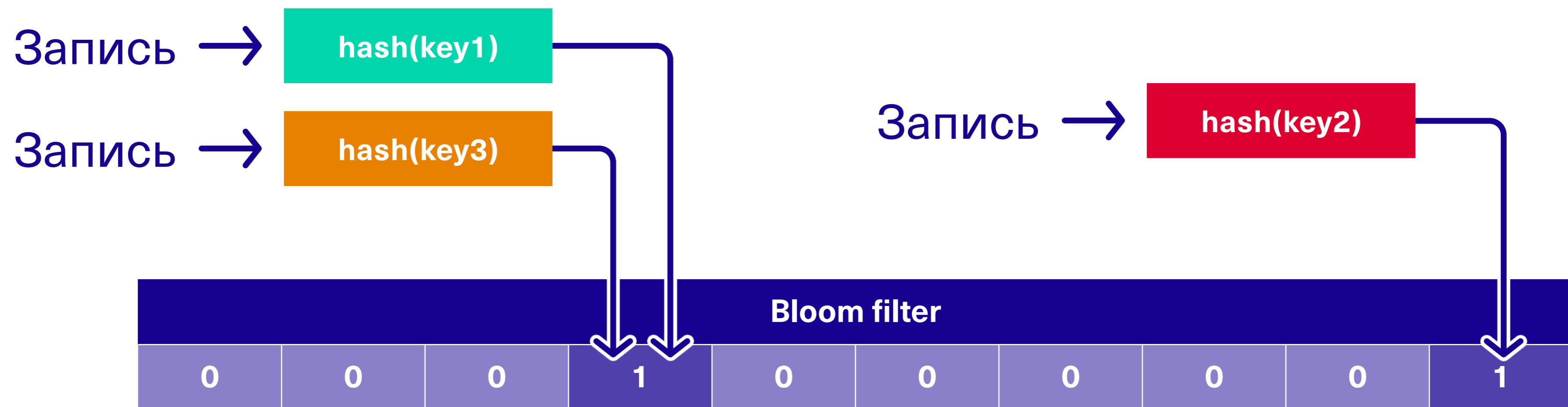
# Bloom filter



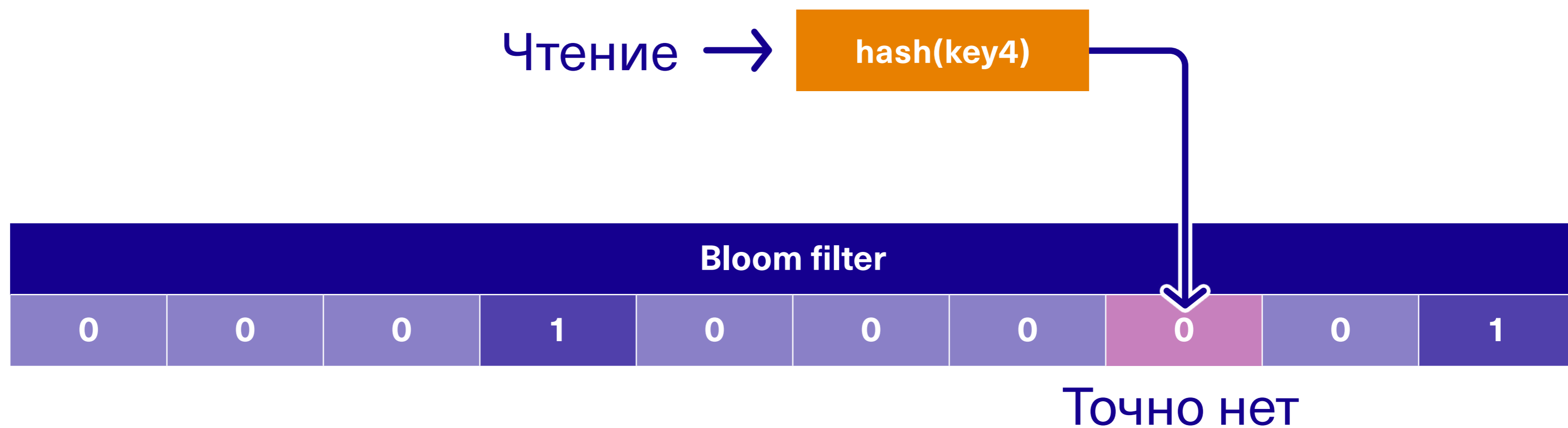
# Bloom filter



# Bloom filter

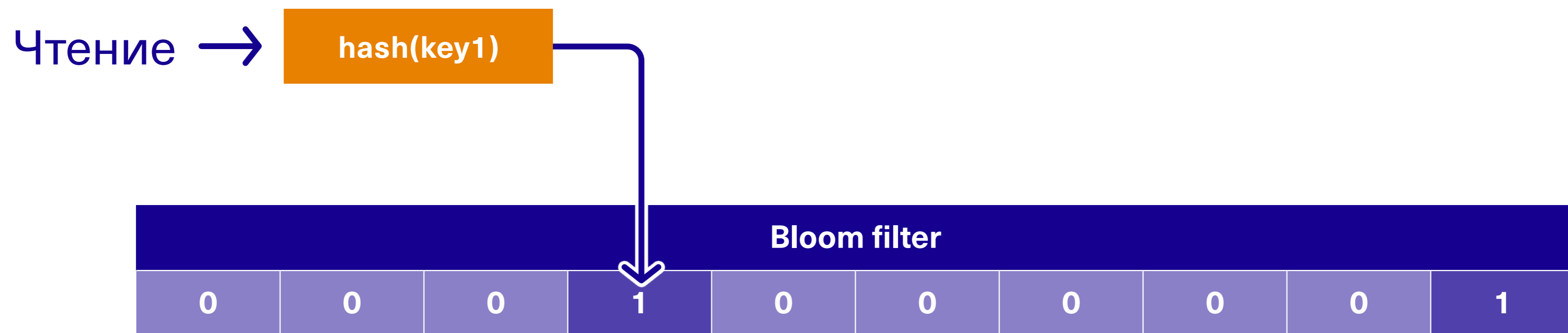


# Bloom filter



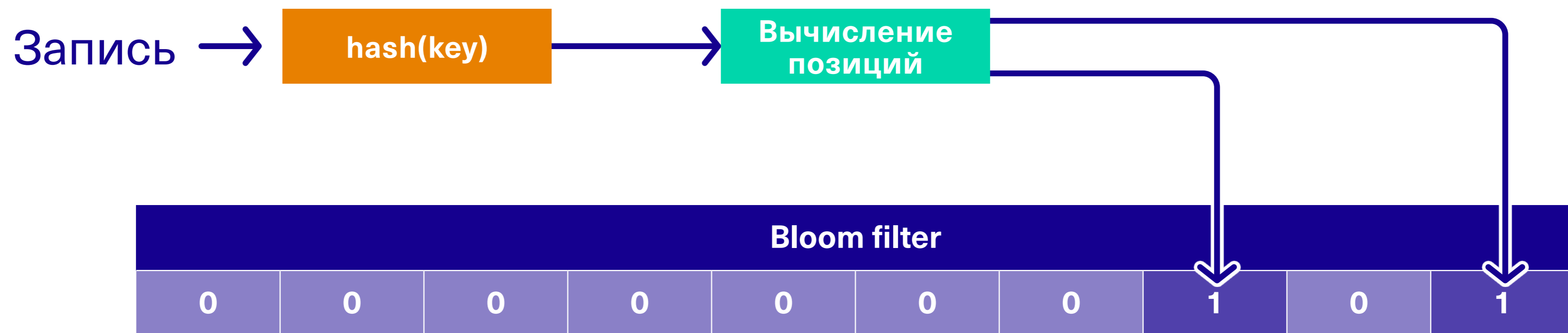


# Bloom filter



Может правда есть,  
а может это коллизия

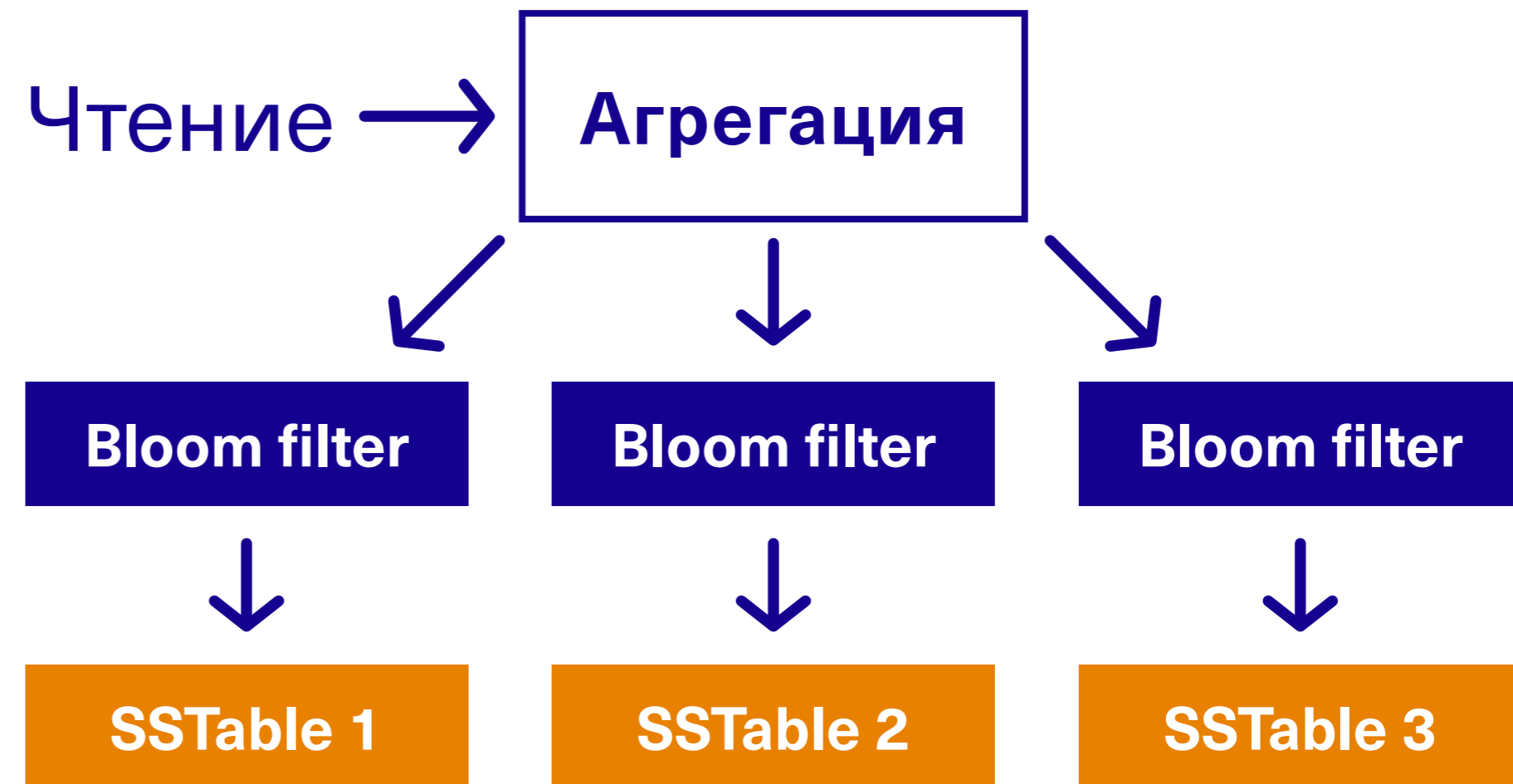
# Bloom filter



# Bloom filter



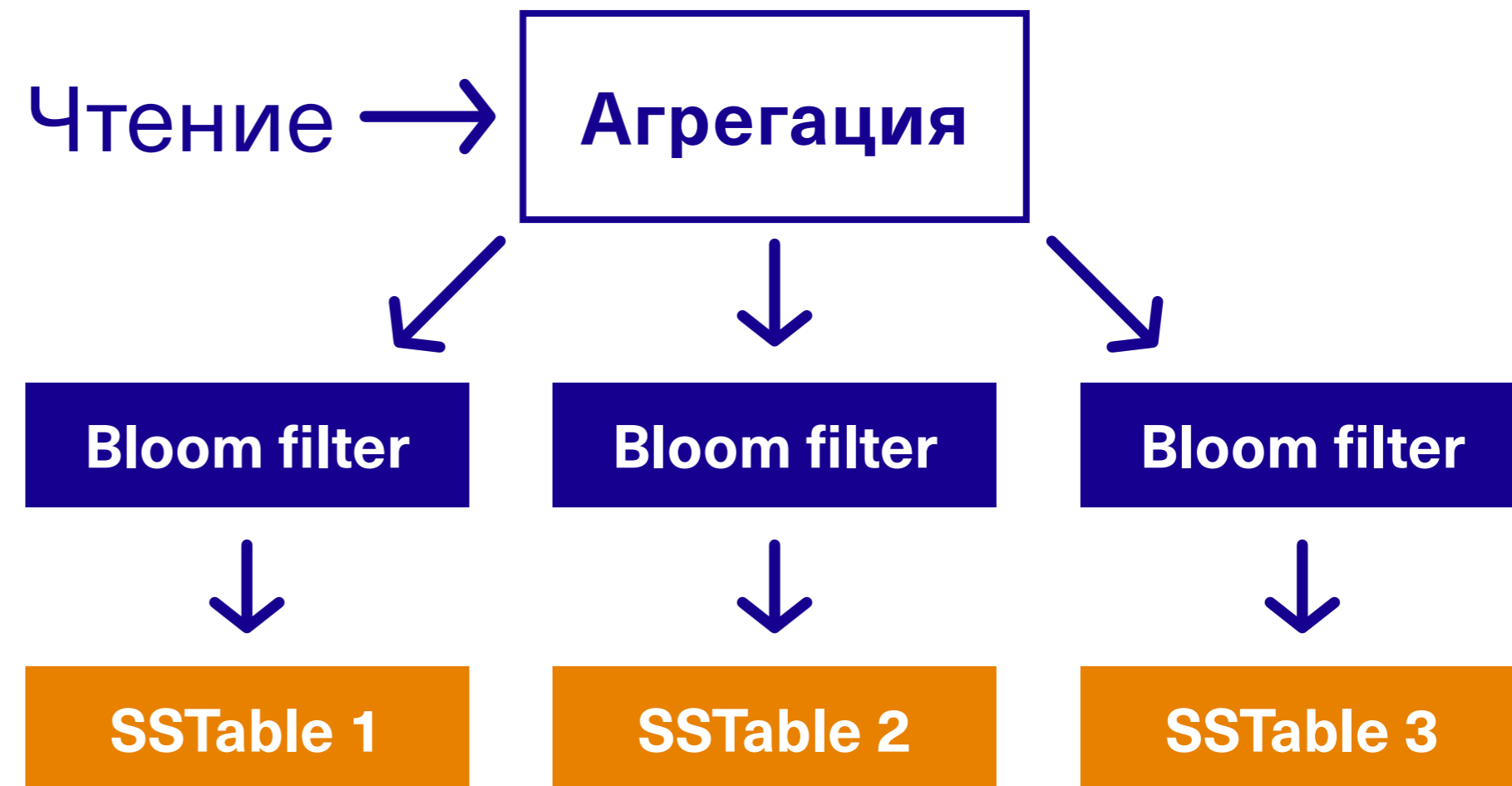
**bloom\_filter\_fp\_chance** —  
вероятность ошибки  
(от 0 до 1, рекомендуется  
от 0.01 до 0.1)



# Bloom filter



**bloom\_filter\_fp\_chance** —  
вероятность ошибки  
(от 0 до 1, рекомендуется  
от 0.01 до 0.1)



```
nodetool tablestats conf.timeline
...
Bloom filter false positives: 16721
Bloom filter false ratio: 0.0390
Bloom filter space used: 9392731
```

# Minor compaction



Много обновляем и читаем, значит compaction должен помочь



# Minor compaction



Много обновляем и читаем, значит compaction должен помочь



# Minor compaction



Много обновляем и читаем, значит compaction должен помочь

SSTable 1

SSTable 3

# Minor compaction



Много обновляем и читаем, значит compaction должен помочь





# Minor compaction



Много обновляем и читаем, значит compaction должен помочь

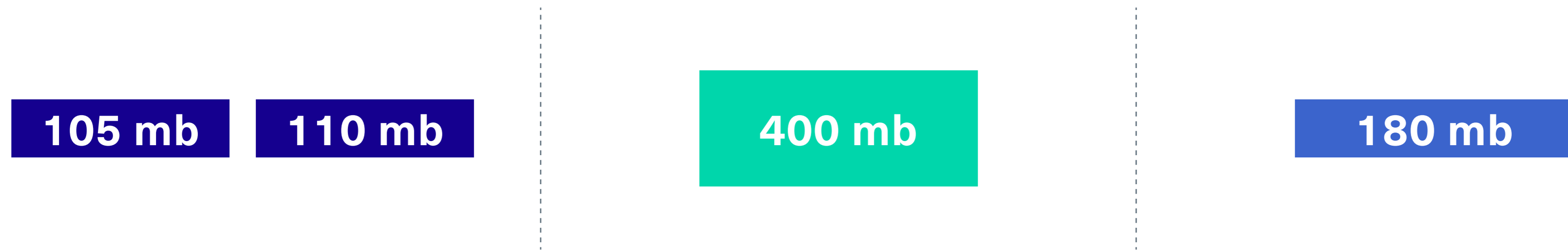


**Size-Tiered Compaction** <https://youtu.be/TyTXOjFMi7k?si=Hsid0evoaOzkEy7o>

**Leveled Compaction** [https://youtu.be/6yJEwqseMY4?si=VaOEeNLFm\\_IbFcoR](https://youtu.be/6yJEwqseMY4?si=VaOEeNLFm_IbFcoR)

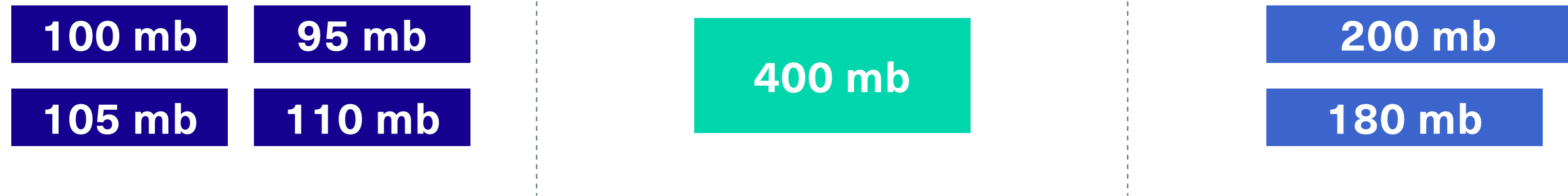
**Time Window Compaction** [https://youtu.be/P8jPlppa\\_8l?si=bVw8EVt3puEVn-MP](https://youtu.be/P8jPlppa_8l?si=bVw8EVt3puEVn-MP)

# Size Tiered Compaction Strategy (STCS)



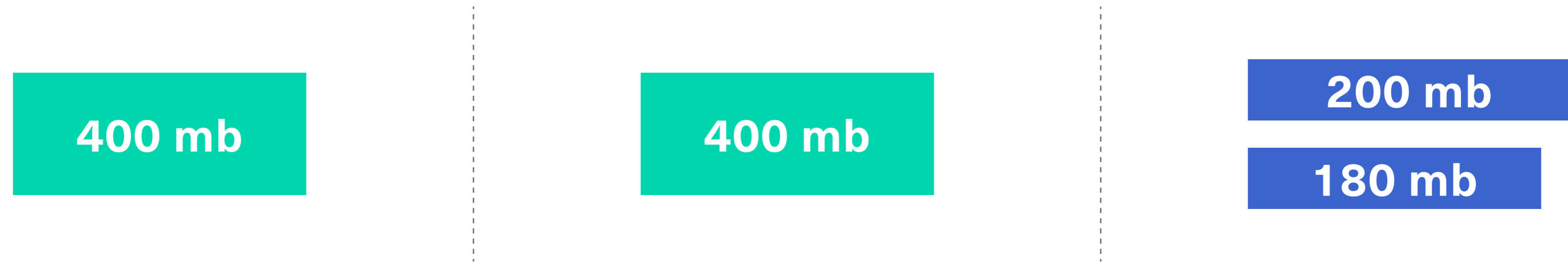
- Стратегия по умолчанию
- Объединяет близкие по размеру SSTable в одну
- Объединяет когда набралось нужное количество SSTable
- Оптимальна при активной записи (и если другие стратегии не подходят)

# Size Tiered Compaction Strategy (STCS)



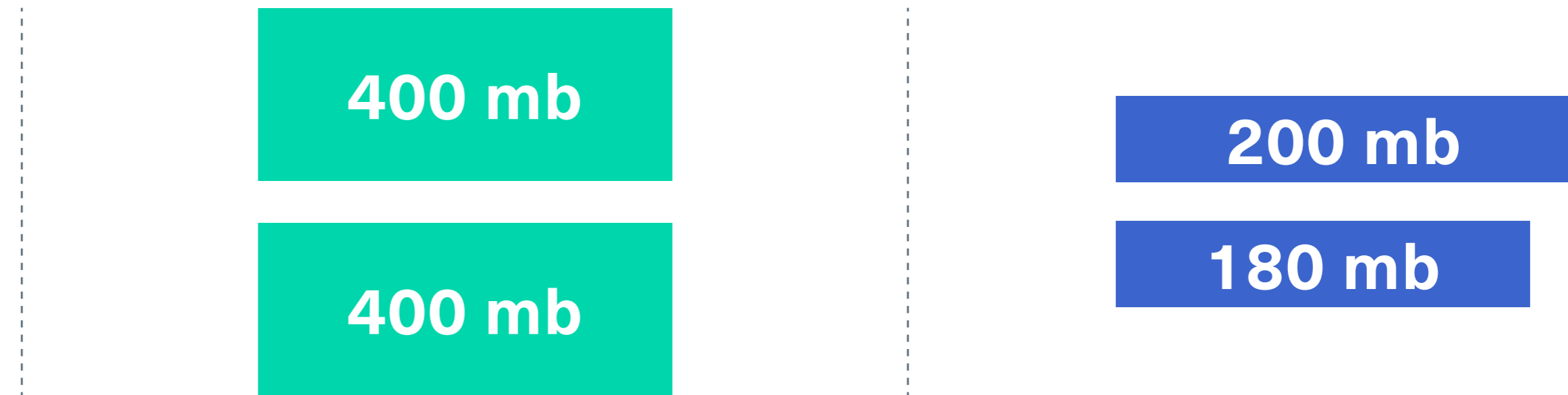
- Стратегия по умолчанию
- Объединяет близкие по размеру SStable в одну
- Объединяет когда набралось нужное количество SStable
- Оптимальна при активной записи (и если другие стратегии не подходят)

# Size Tiered Compaction Strategy (STCS)



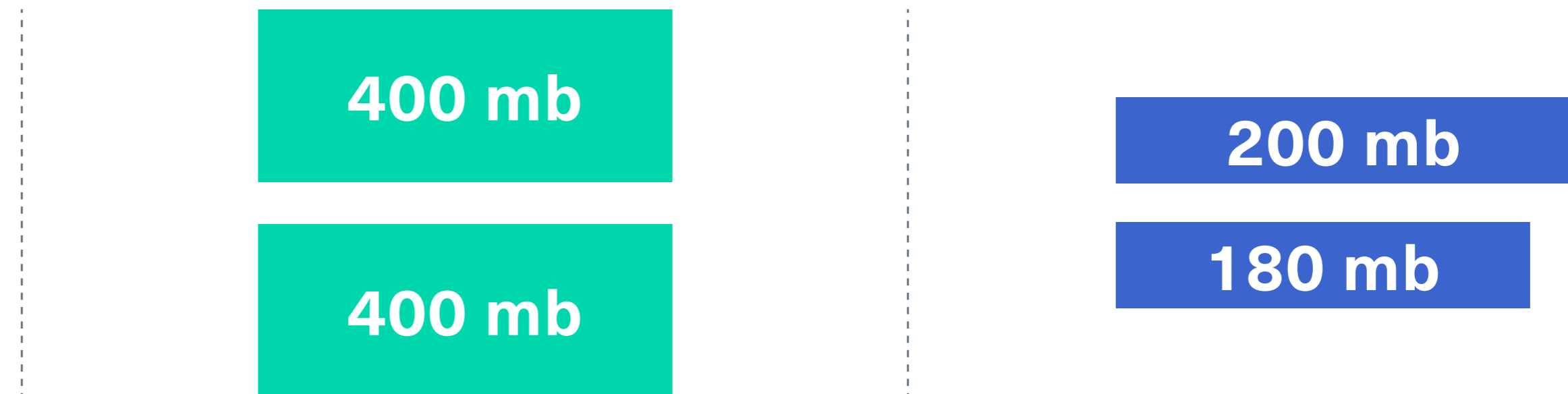
- Стратегия по умолчанию
- Объединяет близкие по размеру SSTable в одну
- Объединяет когда набралось нужное количество SSTable
- Оптимальна при активной записи (и если другие стратегии не подходят)

# Size Tiered Compaction Strategy (STCS)



- Стратегия по умолчанию
- Объединяет близкие по размеру SSTable в одну
- Объединяет когда набралось нужное количество SSTable
- Оптимальна при активной записи (и если другие стратегии не подходят)

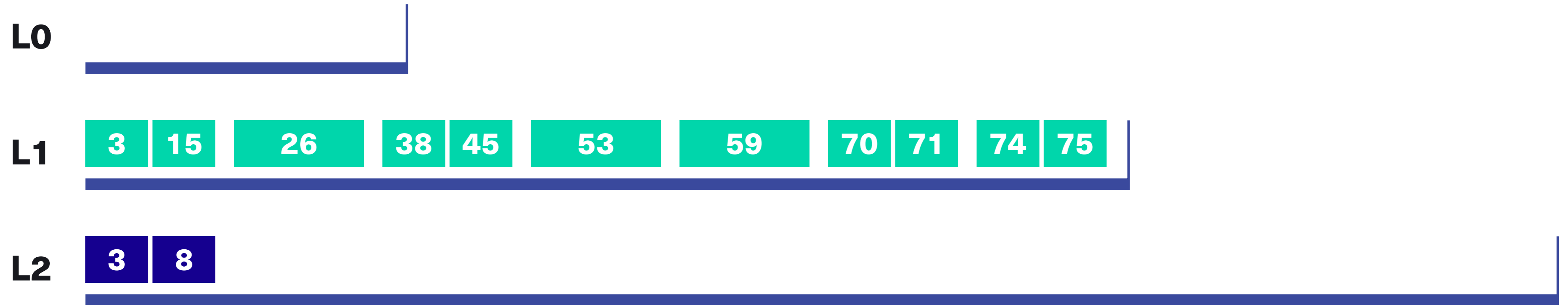
# Size Tiered Compaction Strategy (STCS)



- Стратегия по умолчанию
- Объединяет близкие по размеру SSTable в одну
- Объединяет когда набралось нужное количество SSTable
- Оптимальна при активной записи (и если другие стратегии не подходят)

```
ALTER TABLE tablename WITH compaction =  
  {'class' : 'SizeTieredCompactionStrategy'}
```

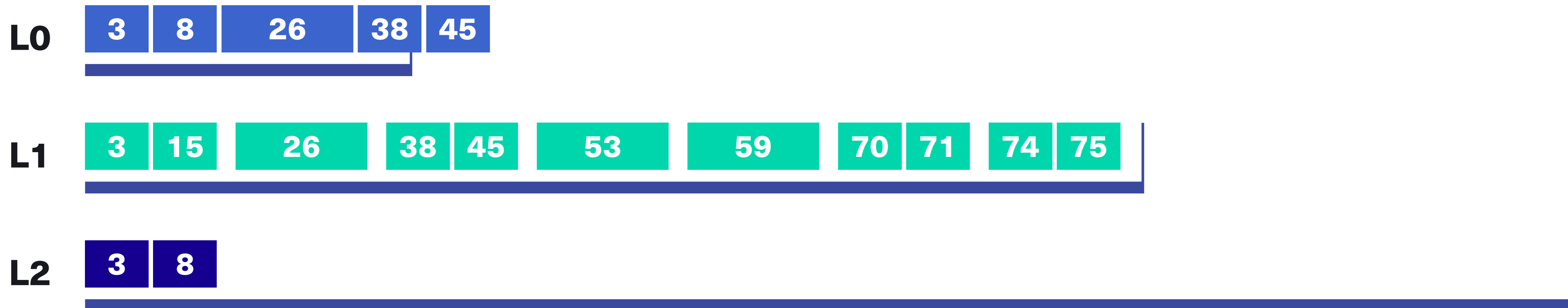
# Levelled Compaction Strategy (LCS)



- Каждый уровень в 10 раз больше предыдущего
- Дублирования ключей внутри уровня нет

```
ALTER TABLE tablename WITH  
  compaction = { 'class' : 'LevelledCompactionStrategy' }
```

# Levelled Compaction Strategy (LCS)

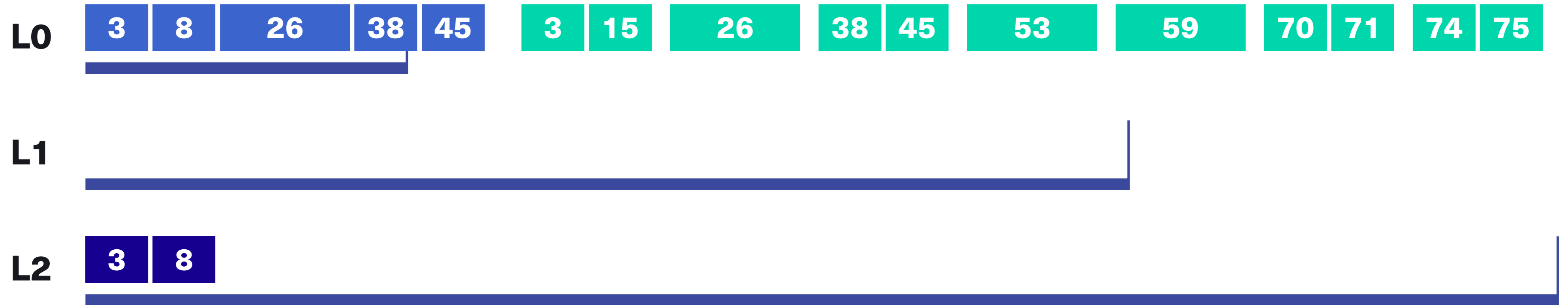


- Каждый уровень в 10 раз больше предыдущего
- Дублирования ключей внутри уровня нет

```
ALTER TABLE tablename WITH  
  compaction = { 'class' : 'LevelledCompactionStrategy' }
```



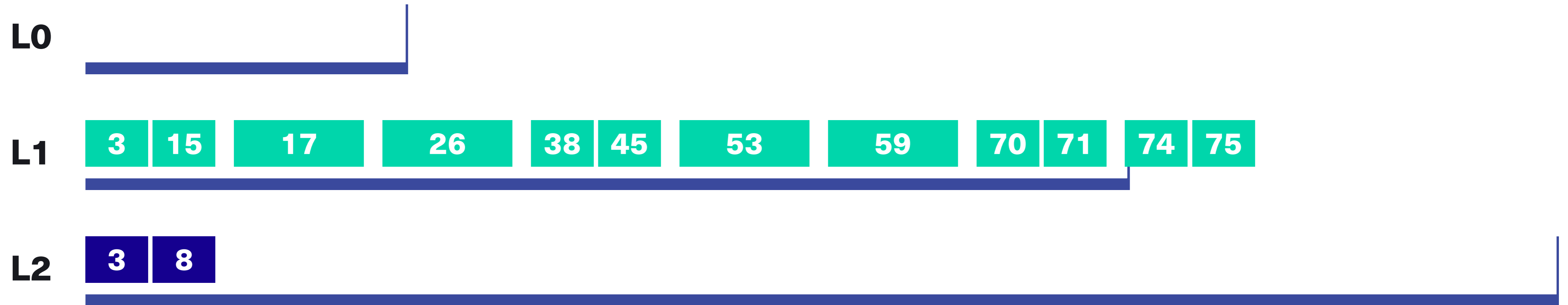
# Leveled Compaction Strategy (LCS)



- Каждый уровень в 10 раз больше предыдущего
- Дублирования ключей внутри уровня нет

```
ALTER TABLE tablename WITH  
  compaction = { 'class' : 'LeveledCompactionStrategy' }
```

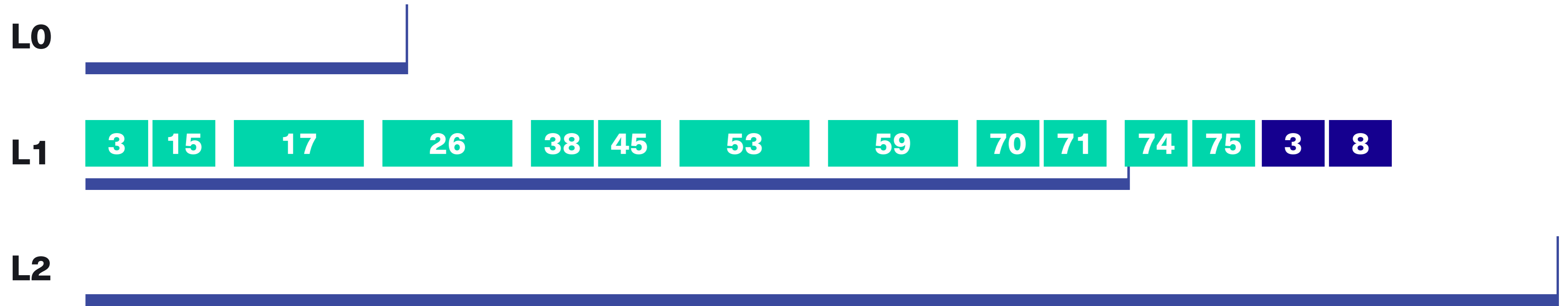
# Levelled Compaction Strategy (LCS)



- Каждый уровень в 10 раз больше предыдущего
- Дублирования ключей внутри уровня нет

```
ALTER TABLE tablename WITH  
  compaction = { 'class' : 'LevelledCompactionStrategy' }
```

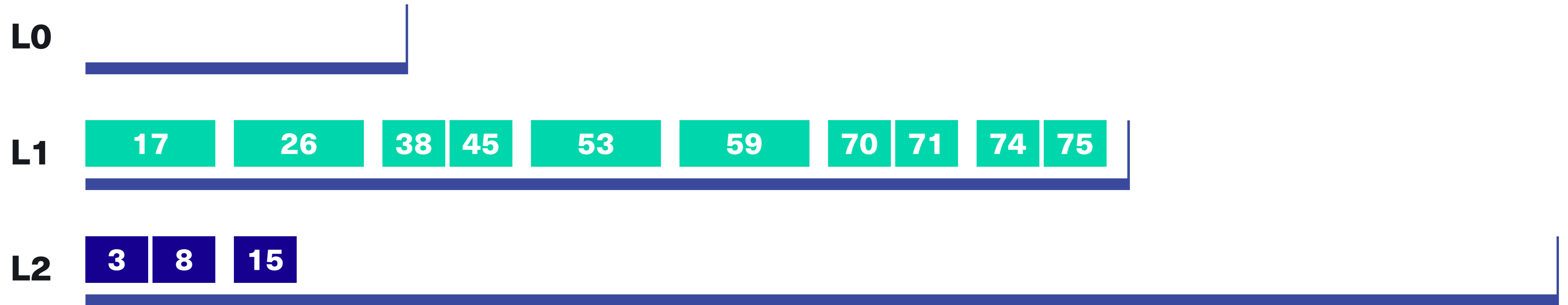
# Levelled Compaction Strategy (LCS)



- Каждый уровень в 10 раз больше предыдущего
- Дублирования ключей внутри уровня нет

```
ALTER TABLE tablename WITH  
  compaction = { 'class' : 'LevelledCompactionStrategy' }
```

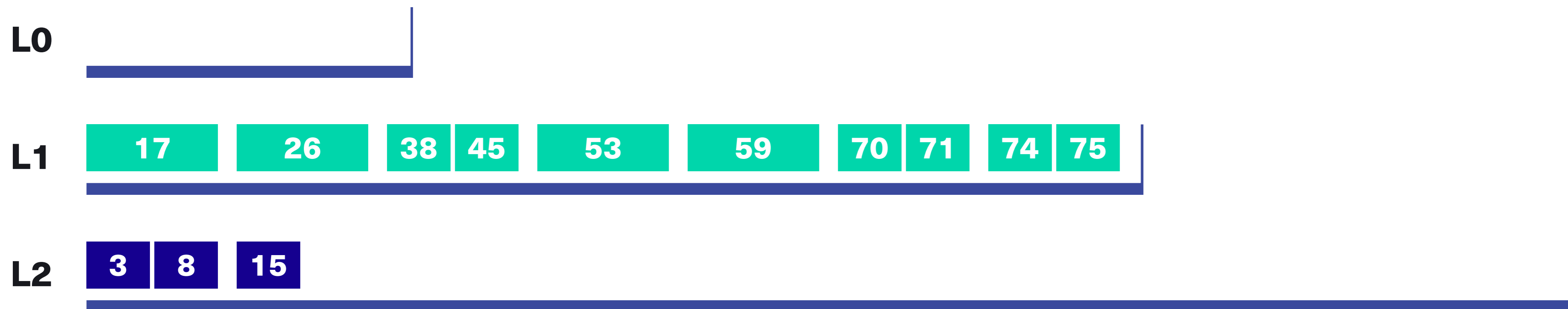
# Levelled Compaction Strategy (LCS)



- Каждый уровень в 10 раз больше предыдущего
- Дублирования ключей внутри уровня нет

```
ALTER TABLE tablename WITH  
  compaction = { 'class' : 'LevelledCompactionStrategy' }
```

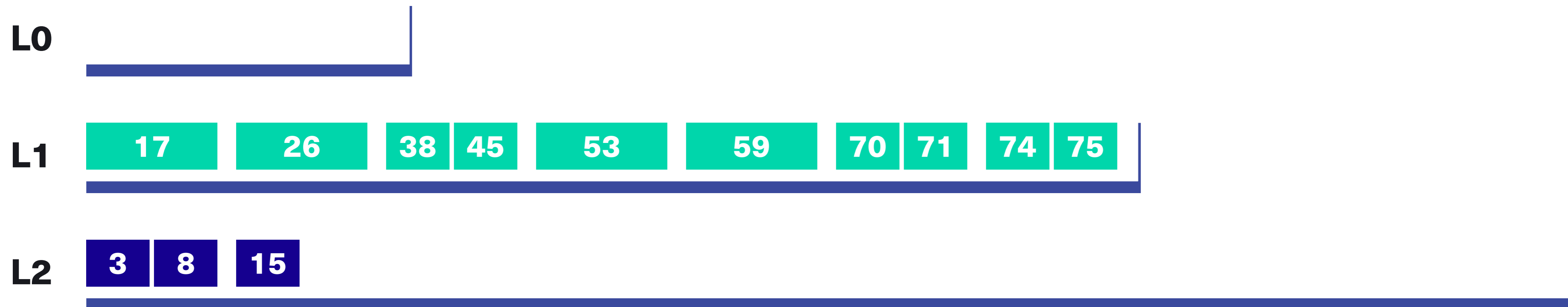
# Levelled Compaction Strategy (LCS)



- Каждый уровень в 10 раз больше предыдущего
- Дублирования ключей внутри уровня нет
- Оптимальна при интенсивных чтениях обновляемых данных

```
ALTER TABLE tablename WITH  
  compaction = { 'class' : 'LevelledCompactionStrategy' }
```

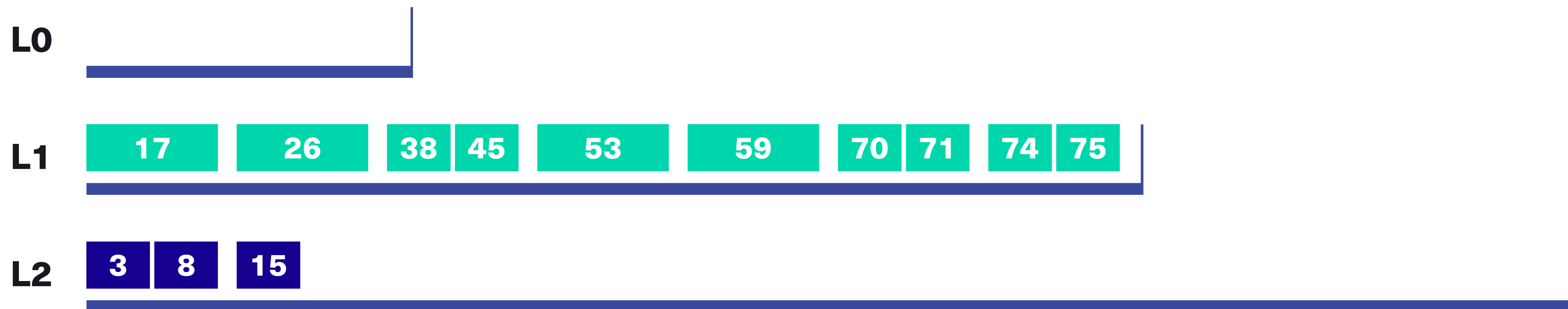
# Levelled Compaction Strategy (LCS)



- Каждый уровень в 10 раз больше предыдущего
- Дублирования ключей внутри уровня нет
- Оптимальна при интенсивных чтениях обновляемых данных
- Нагружает IO

```
ALTER TABLE tablename WITH  
  compaction = { 'class' : 'LevelledCompactionStrategy' }
```

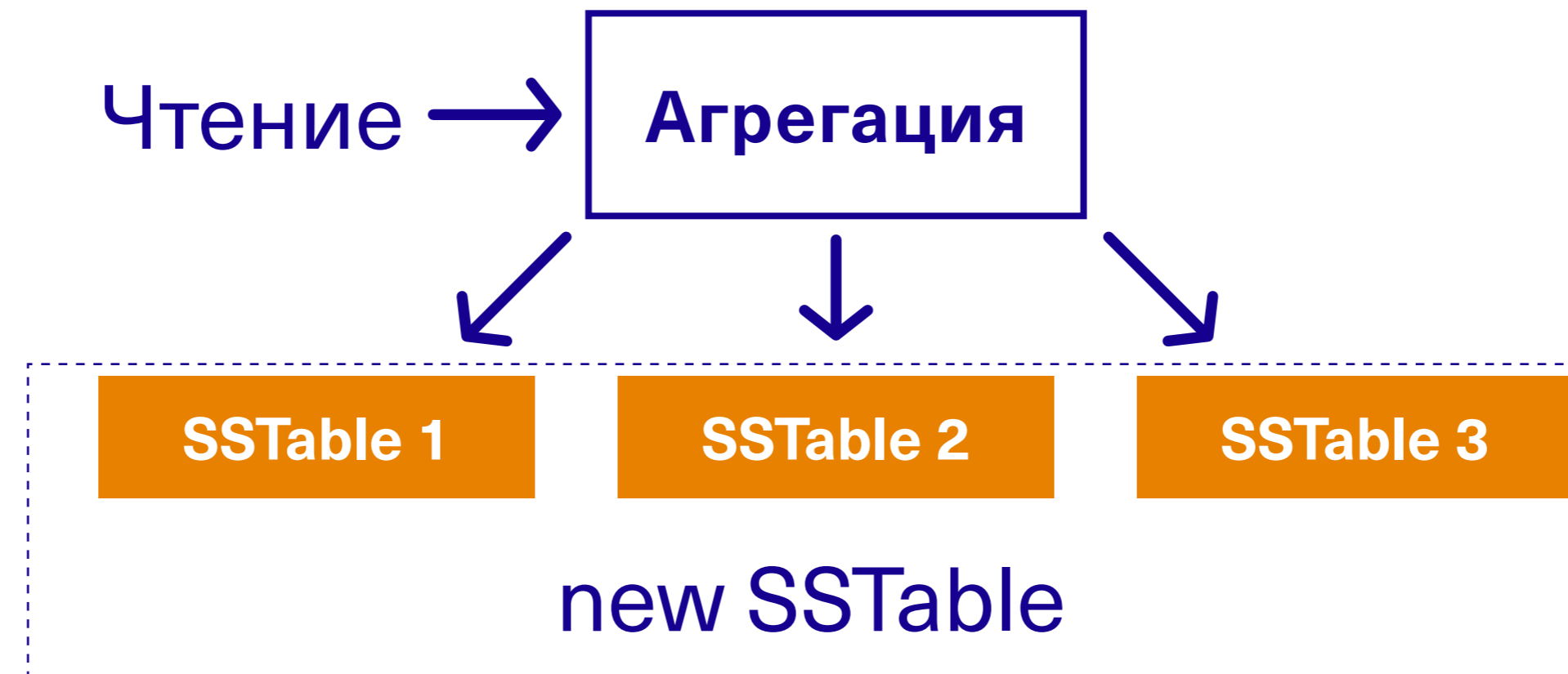
# Levelled Compaction Strategy (LCS)



- Каждый уровень в 10 раз больше предыдущего
- Дублирования ключей внутри уровня нет
- Оптимальна при интенсивных чтениях обновляемых данных
- Нагружает IO — **в нашем кейсе нет запаса IO**

```
ALTER TABLE tablename WITH  
  compaction = { 'class' : 'LevelledCompactionStrategy' }
```

# Major compaction

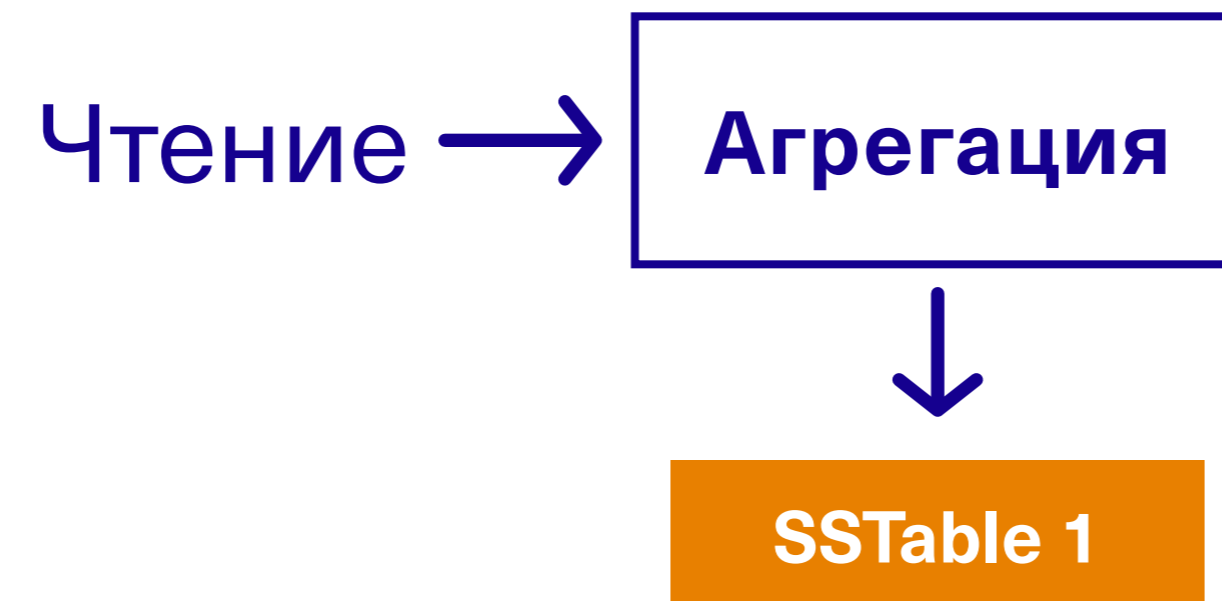


- Мёржит все SSTable в одну
- Дублирования ключа (устаревших версий данных) нет
- Запускается по команде

```
nodetool compact keyspace table
```



# Major compaction



- Мёржит все SSTable в одну
- Дублирования ключа (устаревших версий данных) нет
- Запускается по команде

```
nodetool compact keyspace table
```

# Major compaction

## Небезопасное решение:

- Большая нагрузка на CPU, IO
- Требуется много места на диске



**WARNING**

# Major compaction

## Небезопасное решение:

- Большая нагрузка на CPU, IO
- Требуется много места на диске
- Отрабатывает долго



**WARNING**

# Major compaction

## Небезопасное решение:

- Большая нагрузка на CPU, IO
- Требуется много места на диске
- Отрабатывает долго
- Большой SSTable может мешать Size Tiered Compaction Strategy



**WARNING**

# Major compaction

## Небезопасное решение:

- Большая нагрузка на CPU, IO
- Требуется много места на диске
- Отрабатывает долго
- Большой SSTable может мешать Size Tiered Compaction Strategy



**WARNING**

**1 TB**

# Major compaction

## Небезопасное решение:

- Большая нагрузка на CPU, IO
- Требуется много места на диске
- Отрабатывает долго
- Большой SSTable может мешать Size Tiered Compaction Strategy

105 mb

110 mb

180 mb



**WARNING**

1 TB

# Major compaction

## Небезопасное решение:

- Большая нагрузка на CPU, IO
- Требуется много места на диске
- Отрабатывает долго
- Большой SSTable может мешать Size Tiered Compaction Strategy



**WARNING**

100 mb

95 mb

105 mb

110 mb

200 mb

180 mb

1 TB

# Major compaction

## Небезопасное решение:

- Большая нагрузка на CPU, IO
- Требуется много места на диске
- Отрабатывает долго
- Большой SSTable может мешать Size Tiered Compaction Strategy

400 mb

200 mb

180 mb

**WARNING**

1 TB



## Apache Cassandra полезный инструмент:

- Настраиваемые согласованность и доступность
- Высокая скорость записи из коробки

## Apache Cassandra полезный инструмент:

- Настраиваемые согласованность и доступность
- Высокая скорость записи из коробки

## Необходимо понимать как она работает под капотом:

- У `using timestamp` есть нюансы
- Tombstone имеют свою цену

## Apache Cassandra полезный инструмент:

- Настраиваемые согласованность и доступность
- Высокая скорость записи из коробки

## Необходимо понимать как она работает под капотом:

- У using timestamp есть нюансы
- Tombstone имеют свою цену

## Нагрузку на IO можно снизить:

- Кэши
- Bloom filter
- Compaction strategies



**Спасибо за внимание!**