

Внедрение YDB CDC на примере Yandex Monitoring

Егор Литвиненко

Yandex Infrastructure



Познакомимся

- Егор Литвиненко
- Работаю в Yandex Infrastructure
- Старший разработчик
Yandex Observability Platform
- Помогаю людям мониторить их сервисы
- В коммерческой разработке 10 лет

Lifestreet, Marketshare, Runawfe, Мультикарта,
АйТи Консалтинг, Магнит...

Цели доклада

1

Разобраться, как
обновлять данные
в сервисе

Цели доклада

1

Разобраться, как
обновлять данные
в сервисе

2

Разобраться,
как правильно
приготовить
YDB CDC

План

1. Теория
2. Практика использования CDC
3. Что пошло не так
4. Итоги

План

1. Теория
2. Практика использования CDC
3. Что пошло не так
4. Итоги

Yandex Monitoring в цифрах



Yandex Monitoring

8

700 млн

метрик в секунду записываем

Yandex Monitoring

9

700 млн

метрик в секунду записываем

это



8,5 GB

в секунду

Yandex Monitoring

10

700 млн

метрик в секунду записываем

это



8,5 GB

в секунду

Храним

6,4 PB

данных

Yandex Monitoring

11

700 млн

метрик в секунду записываем

это



8,5 GB

в секунду

Храним

6,4 PB

данных

Храним

8 млн

алертов пользователей

Yandex Monitoring

12

700 млн

метрик в секунду записываем

это



8,5 GB

в секунду

Храним

6,4 PB

данных

Храним

8 млн

алертов пользователей

Считаем

150 тысяч

алертов в секунду

Постановка задачи



Постановка задачи

- Сервис Yandex Monitoring запущен на тысячах нод

Постановка задачи

- Сервис Yandex Monitoring запущен на тысячах нод
- Сервису для работы нужны конфиги из таблиц базы

Постановка задачи

- Сервис Yandex Monitoring запущен на тысячах нод
- Сервису для работы нужны конфиги из таблиц базы
- Таблицы логически связаны и образуют иерархию

Постановка задачи

- Сервис Yandex Monitoring запущен на тысячах нод
- Сервису для работы нужны конфиги из таблиц базы
- Таблицы логически связаны и образуют иерархию
- Меняются редко

Постановка задачи

- Сервис Yandex Monitoring запущен на тысячах нод
- Сервису для работы нужны конфиги из таблиц базы
- Таблицы логически связаны и образуют иерархию
- Меняются редко
- Наша цель: чтобы ноды быстро забирали **только изменения**

Как мы решали задачу?



Как мы решали задачу?

- Таблицы небольшие

Как мы решали задачу?

- Таблицы небольшие
- Скачиваем таблицы целиком по таймеру

Как мы решали задачу?

- Таблицы небольшие
- Скачиваем таблицы целиком по таймеру
- Делаем это на каждой ноде

Как мы решали задачу?

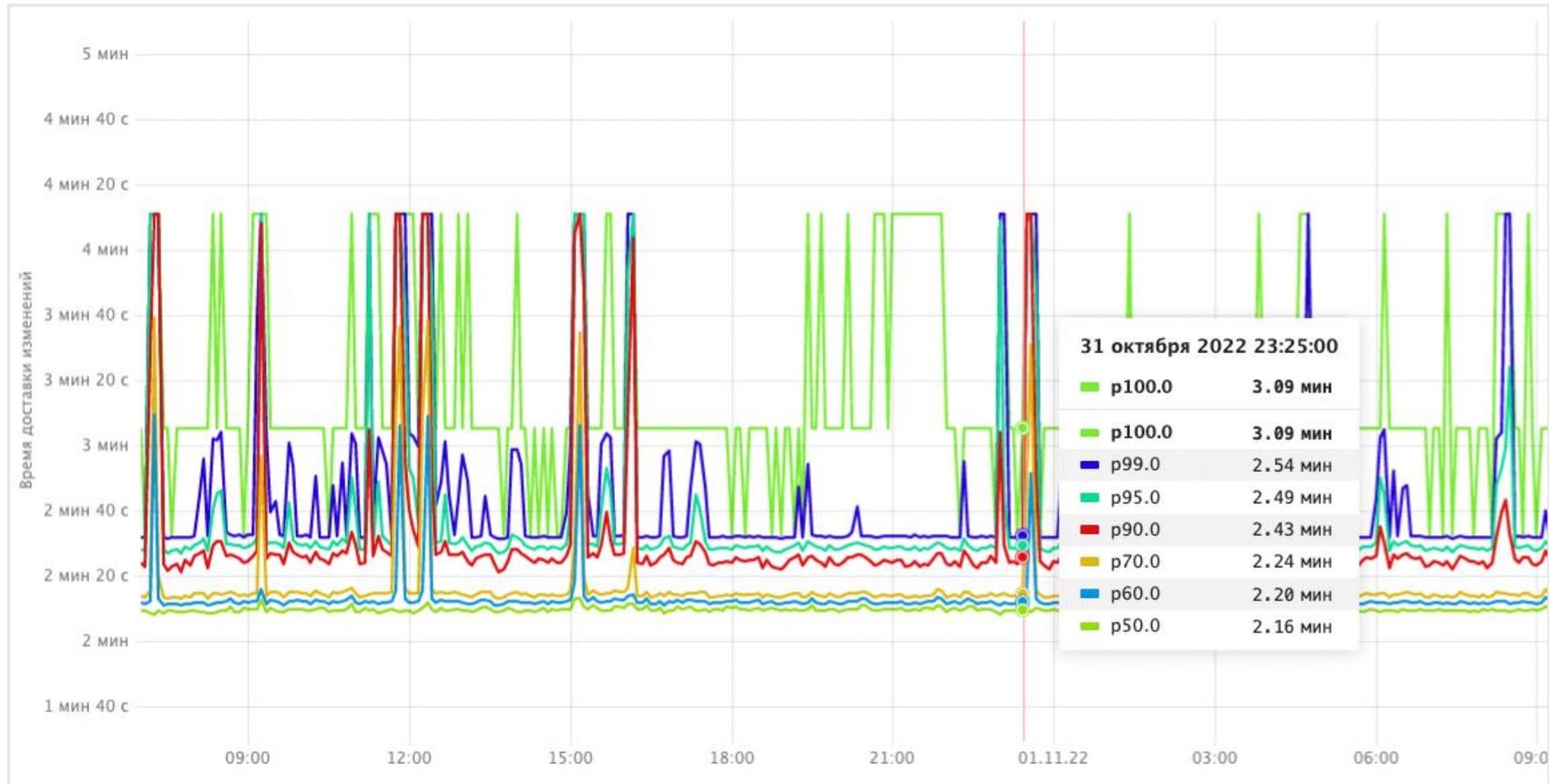
- Таблицы небольшие
- Скачиваем таблицы целиком по таймеру
- Делаем это на каждой ноде
- Используем механизм **YDB ReadTable (gRPC stream)**



YDB ReadTable
click.ru/34ccJd

Какой был лаг доставки?

24



Какие данные мы получаем?

25

- Таблицы Project, Cluster, ...

Shard

```
id : String
projectId : String
clusterId : String
serviceId : String
```

(*..1) (*..1)

Cluster

```
id : String
projectId : String
```

(*..1)

Project

```
id : String
```

Service

```
id : String
projectId : String
serviceProviderId : String
```

(*..1) (*..1)

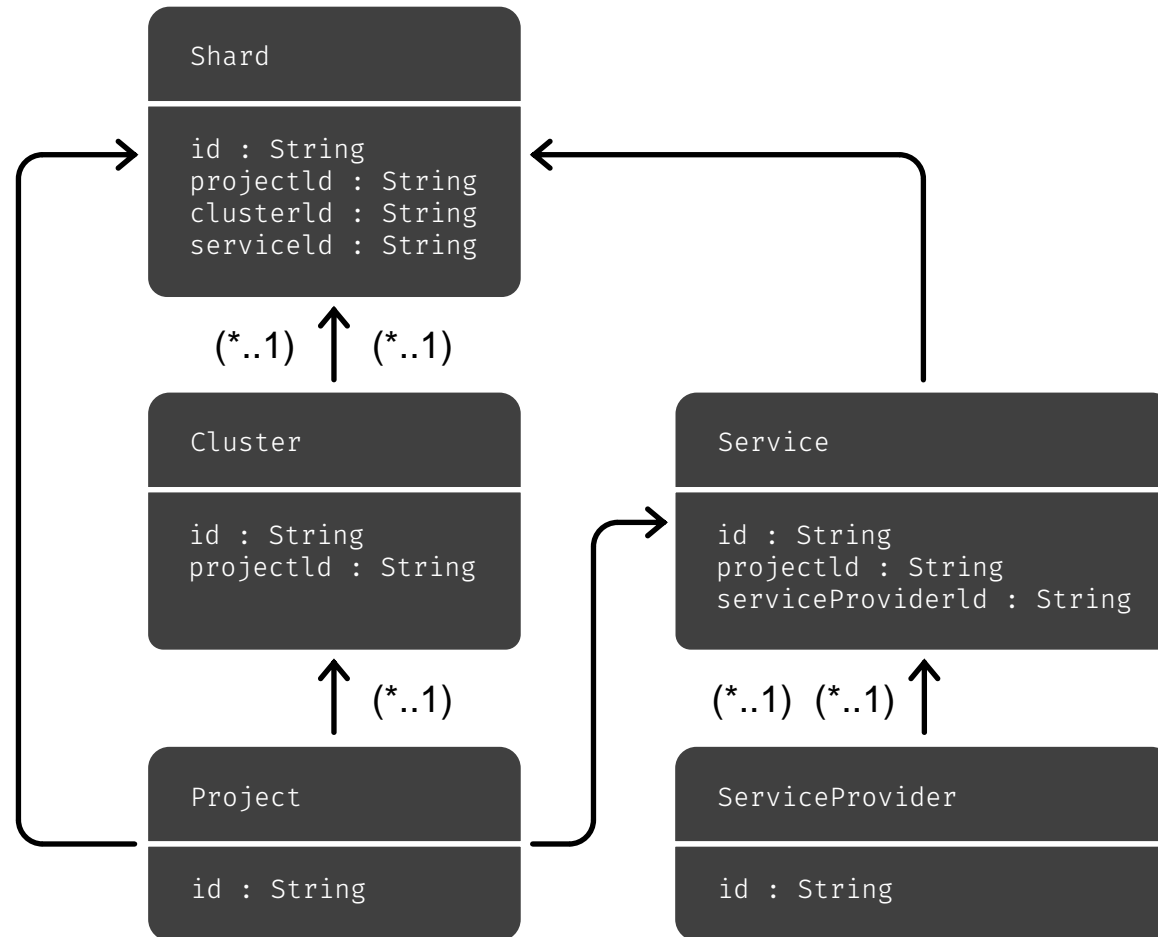
ServiceProvider

```
id : String
```

Какие данные мы получаем?

26

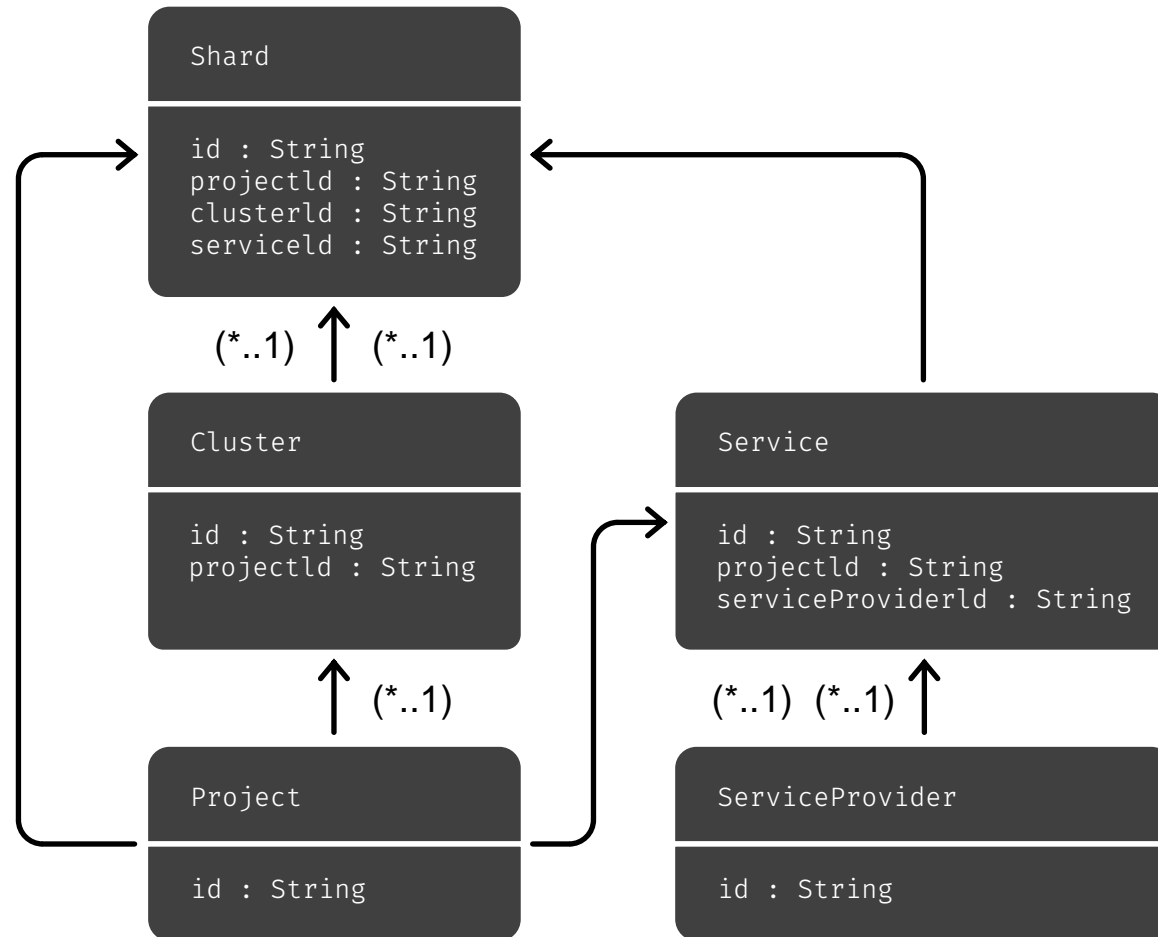
- Таблицы Project, Cluster, ...
- Сущности связаны между собой



Какие данные мы получаем?

27

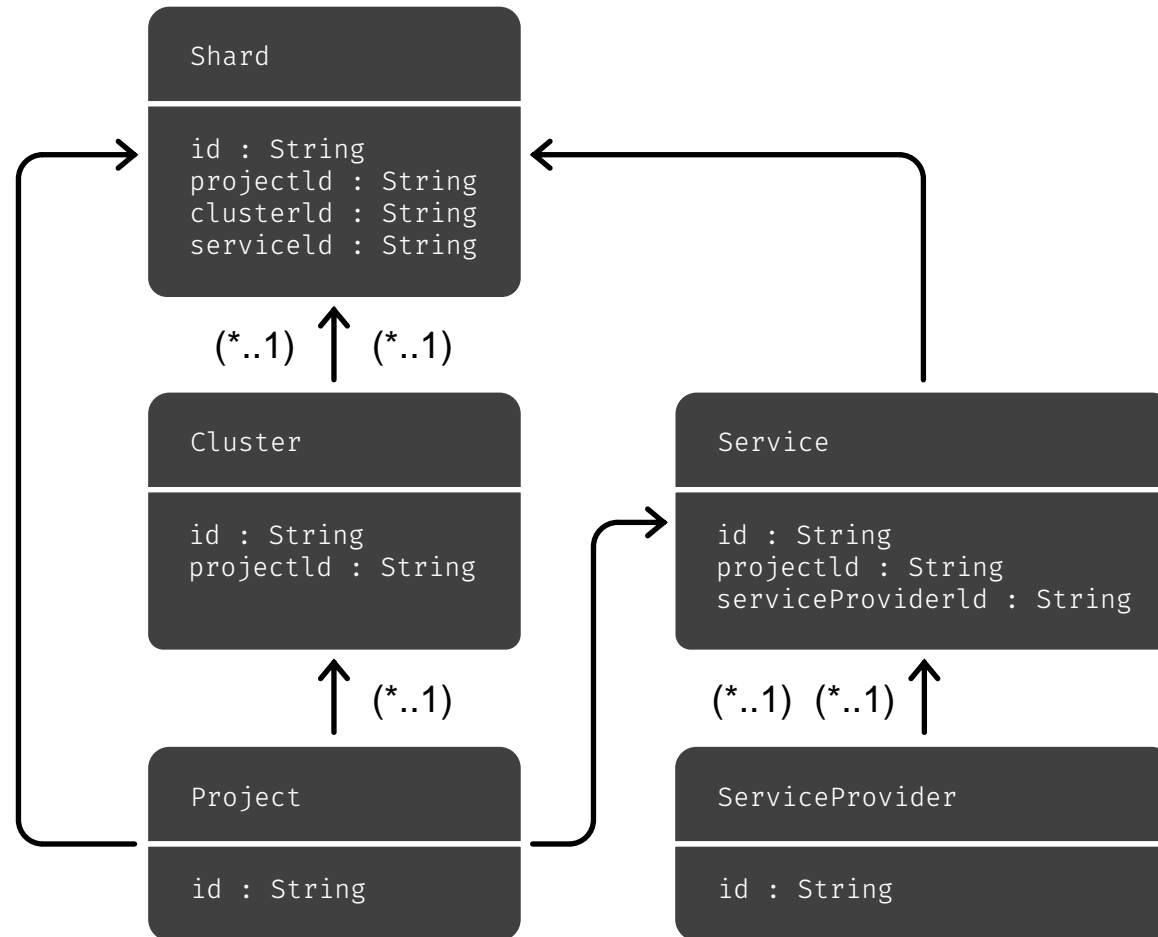
- Таблицы Project, Cluster, ...
- Сущности связаны между собой
- В базе foreign keys нет



Какие данные мы получаем?

28

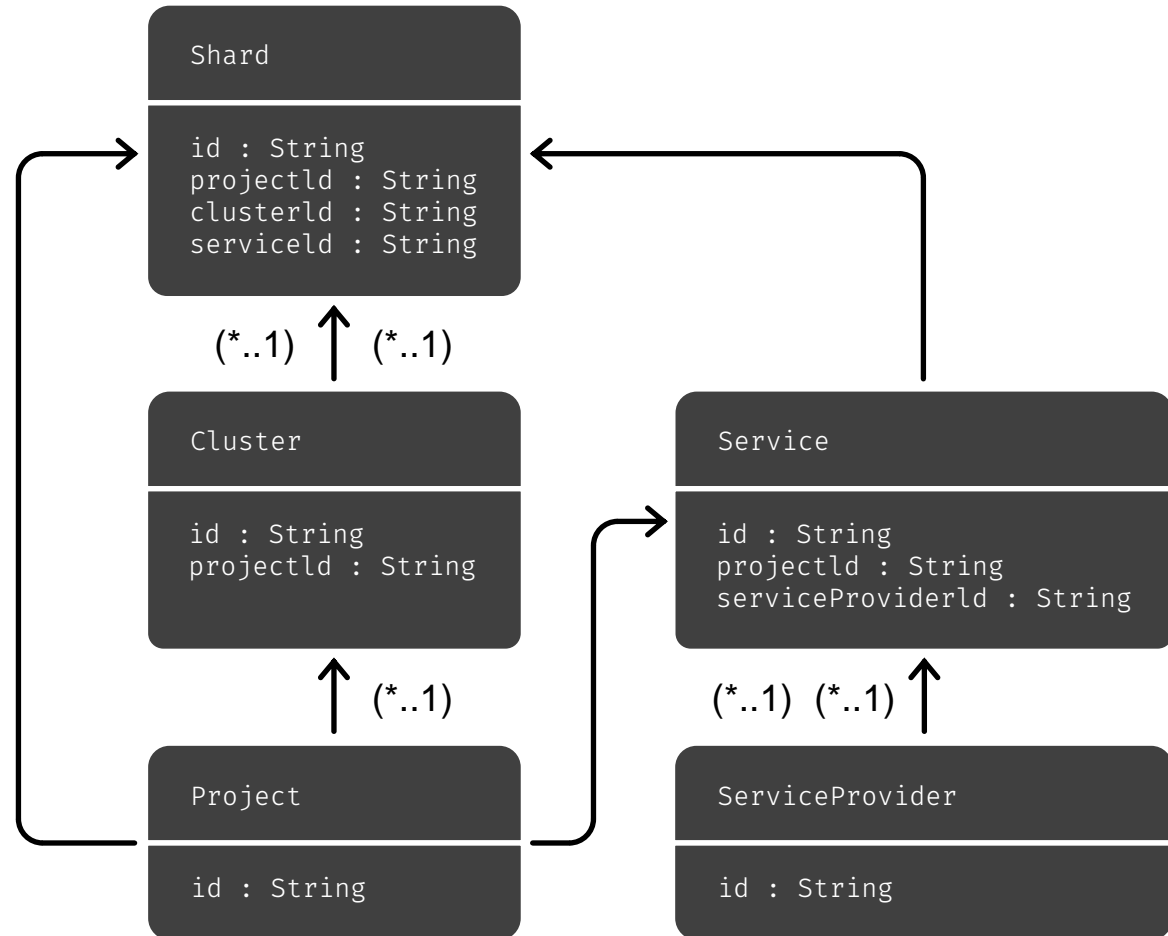
- Таблицы Project, Cluster, ...
- Сущности связаны между собой
- В базе foreign keys нет
- Меняются в разных транзакциях



Какие данные мы получаем?

29

- Таблицы Project, Cluster, ...
- Сущности связаны между собой
- В базе foreign keys нет
- Меняются в разных транзакциях
- Становятся только больше



Как брать из таблицы изменения



Как брать из таблицы изменения

31

- Slowly changing dimensions (SCD)

Любая база данных

- Triggers on tables

Любая база данных с триггерами



SCD

clck.ru/34cciJ

SCD Type4

```
SELECT * FROM SHARDS_HISTORY  
WHERE SHARD_ID = "123"  
AND UPDATED_AT > previous_read_time  
ORDER BY UPDATED_AT DESC  
LIMIT 1;
```

```
CREATE TRIGGER ON SHARDS  
INSERT, UPDATE, DELETE ...
```


SCD Type4

```
SELECT * FROM SHARDS_HISTORY  
WHERE SHARD_ID = "123"  
AND UPDATED_AT > previous_read_time  
ORDER BY UPDATED_AT DESC  
LIMIT 1;
```

// В YDB нет триггеров

```
CREATE TRIGGER ON SHARDS  
INSERT, UPDATE, DELETE ...
```

Как брать из таблицы изменения

- Slowly changing dimensions (SCD)

Любая база данных

- Triggers on tables

Любая база данных с триггерами

- Transaction log scanners

Binlog в MySQL, Replication log в Postgres, ...

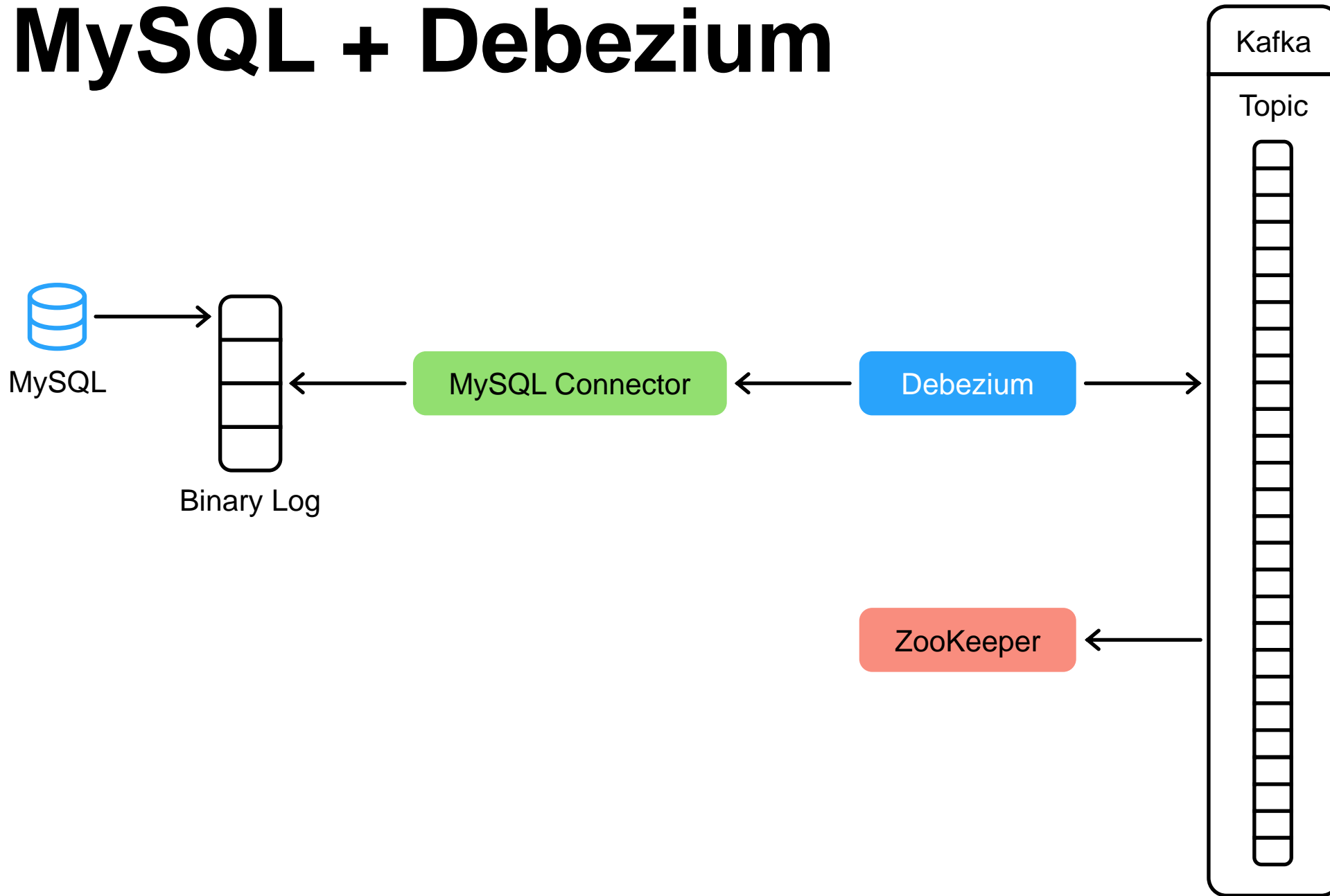
MySQL Binlog

```
mysql> SHOW BINARY LOGS;
```

Log_name	File_size
binlog.000130	27459
binlog.000131	13719
binlog.000132	43268

MySQL + Debezium

36



Как брать из таблицы изменения

37

- Slowly changing dimensions (SCD)

Любая база данных

- Event source

Не зависит от базы данных

- Triggers on tables

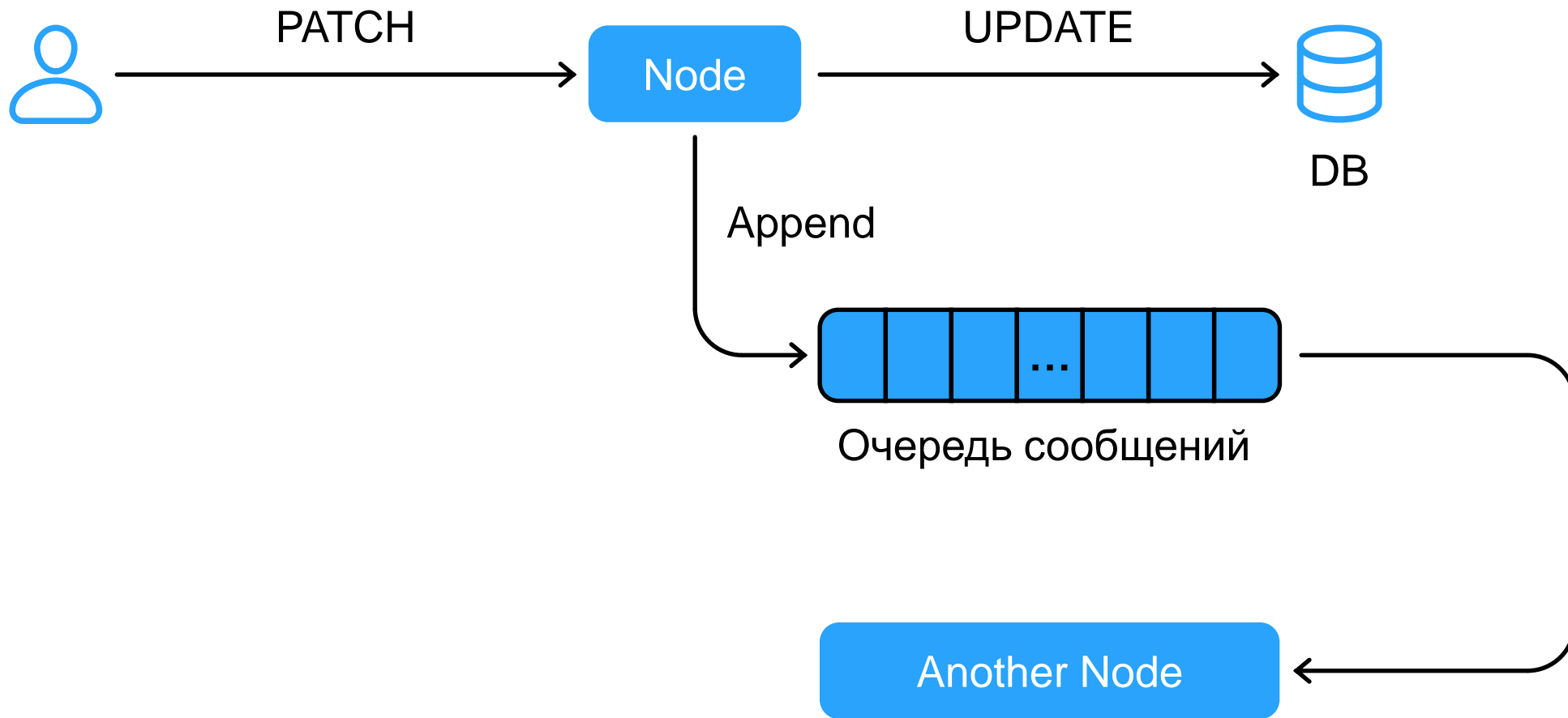
Любая база данных с триггерами

- Transaction log scanners

Binlog в MySQL, Replication log в Postgres, ...

Event source

38



Как брать из таблицы изменения

39

- Slowly changing dimensions (SCD)

Любая база данных

- Triggers on tables

Любая база данных с триггерами

- Transaction log scanners

Binlog в MySQL, Replication log в Postgres, ...

- Event source

Не зависит от базы данных

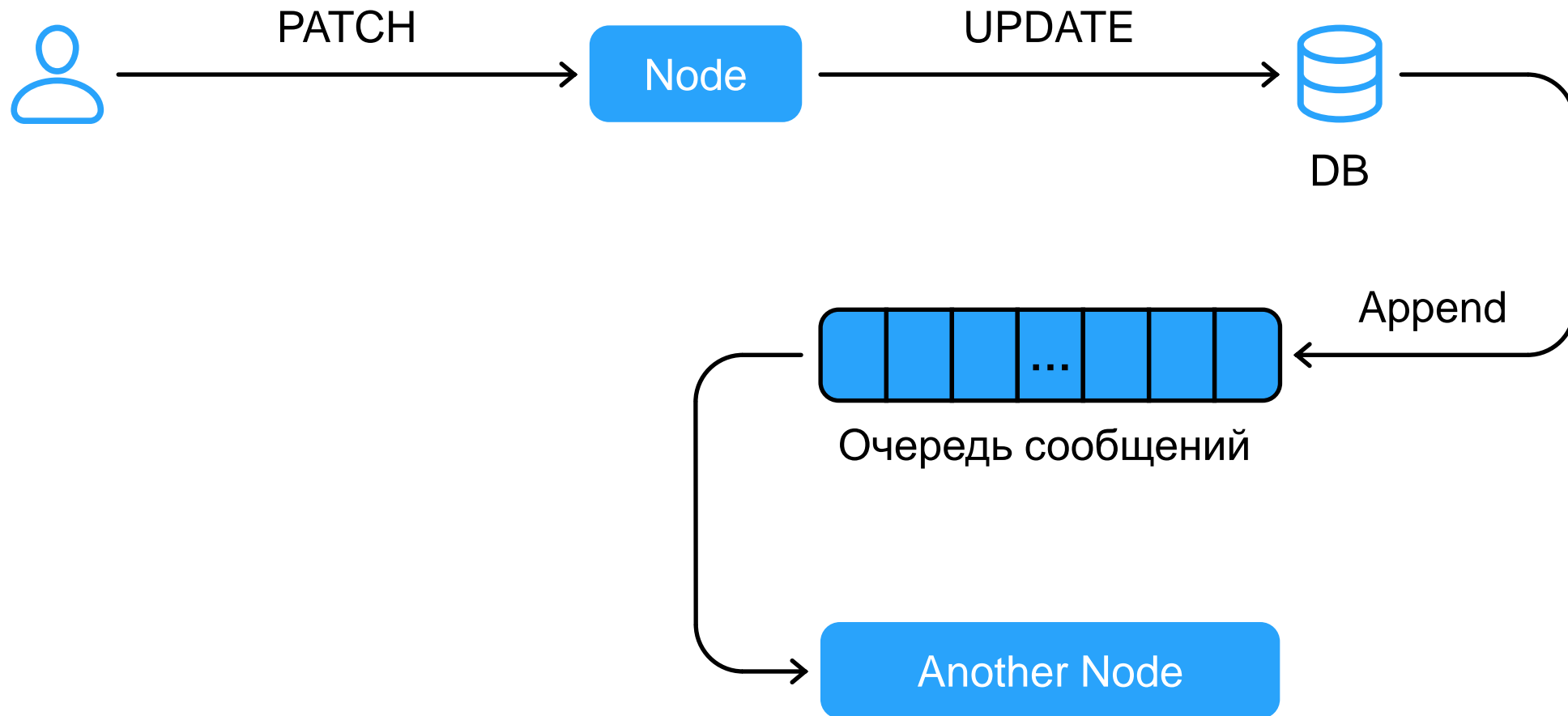
- Change Data Capture (CDC)

YDB, CockroachDB, RethinkDB, Tarantool...

Debezium, Databricks Live Table, Zendesk Daemon...

CDC

40



Сравнение подходов

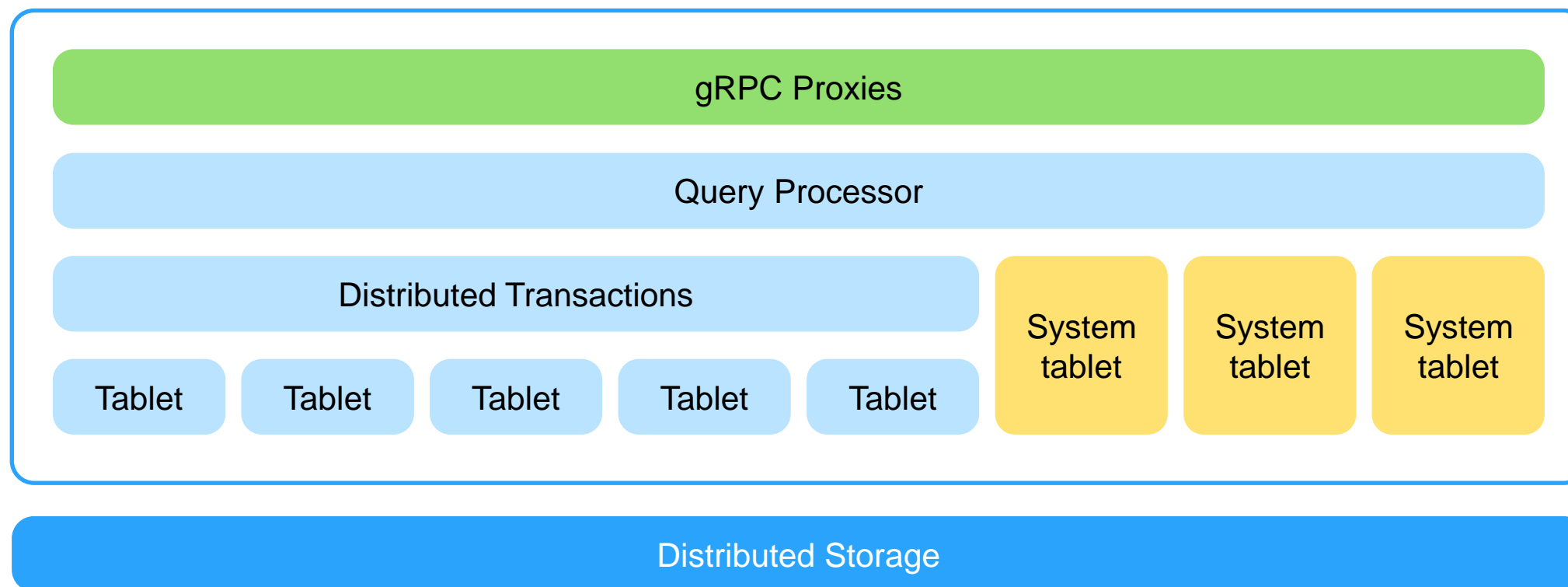
41

	SCD	Event sourcing	CDC
Писать код	Много	Много	Мало
Избыточность таблиц и индексов	Много	Нет	Нет

- Open-source распределённая база данных



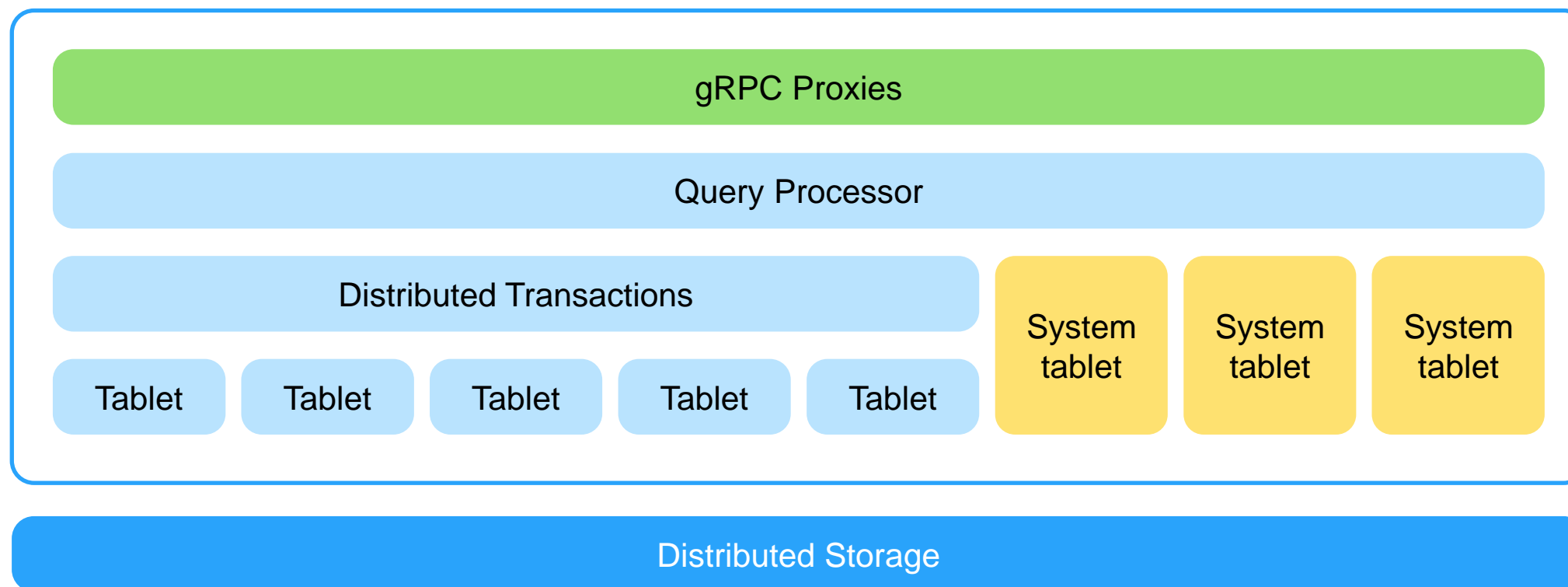
clck.ru/sG59F



- Open-source распределённая база данных
- Поддерживает очереди



clck.ru/sG59F



YDB CDC



YDB CDC

- Очереди партиционированы так же, как таблицы

YDB CDC

- Очереди партиционированы так же, как таблицы
- Все записи в очередь асинхронно

YDB CDC

- Очереди партиционированы так же, как таблицы
- Все записи в очередь асинхронно
- Упорядоченность изменений для primary key

YDB CDC

- Очереди партиционированы так же, как таблицы
- Все записи в очередь асинхронно
- Упорядоченность изменений для primary key
- Exactly Once

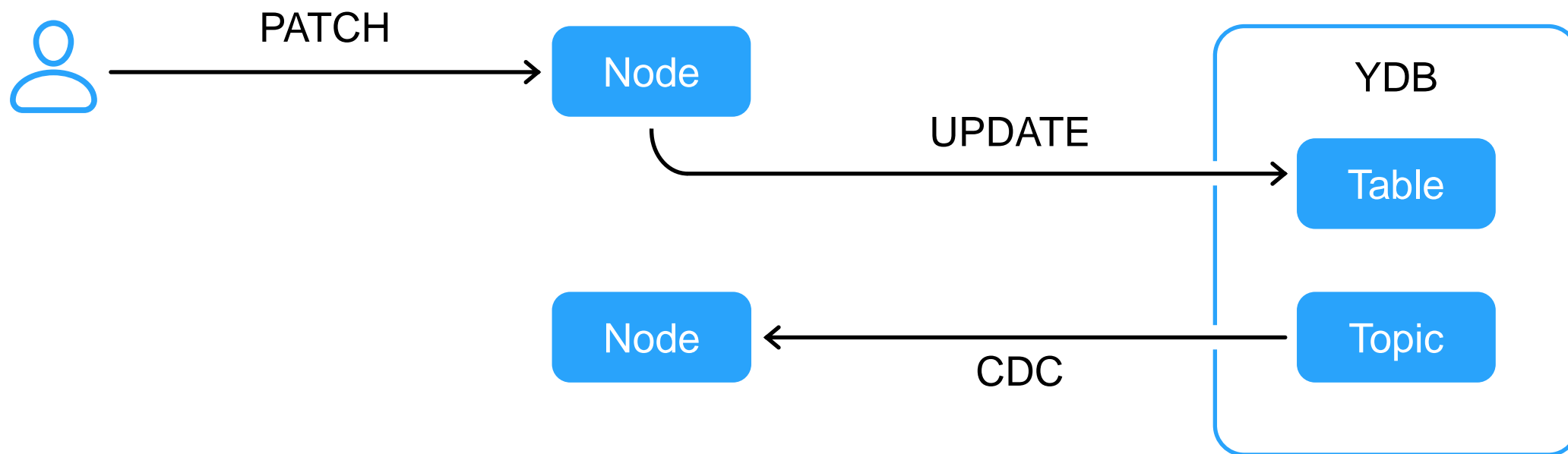
- Очереди партиционированы так же, как таблицы
- Все записи в очередь асинхронно
- Упорядоченность изменений для primary key
- Exactly Once
- TTL по умолчанию 24 часа, максимум 30 дней

Пример сообщения YDB CDC

```
{  
  "update": {},  
  "newImage": {  
    "createdAt": 1654947225226,  
    "updatedAt": 1654947225226,  
    "maxResponseSizeBytes": 10485760,  
    "state": "RW",  
    "maxFileSensors": 1000000,  
    "protoQuotas": "KGQwsQI=",  
    ...,  
    "maxMemSensors": 500000,  
  },  
  "key": [  
    "composite_key_part1",  
    "composite_key_part2"  
  ]  
}
```

Схема YDB CDC

51

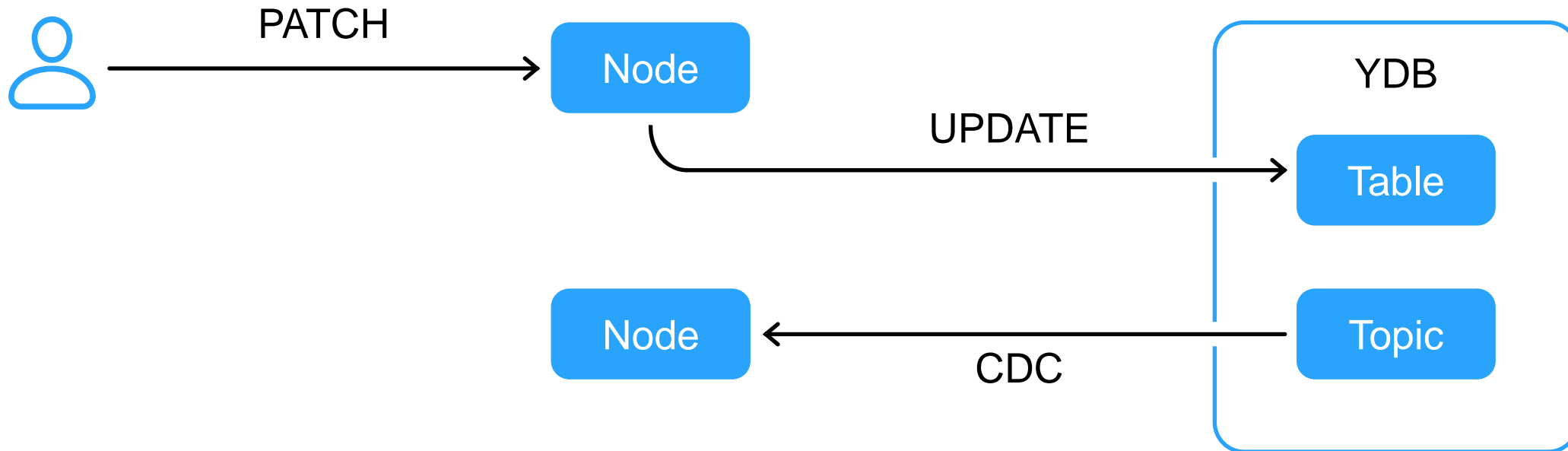


План

1. Теория
2. Практика использования CDC
3. Что пошло не так
4. Итоги

Схема YDB CDC

53



Первая загрузка таблицы



Первая загрузка таблицы

55

- Используем YDB ReadTable

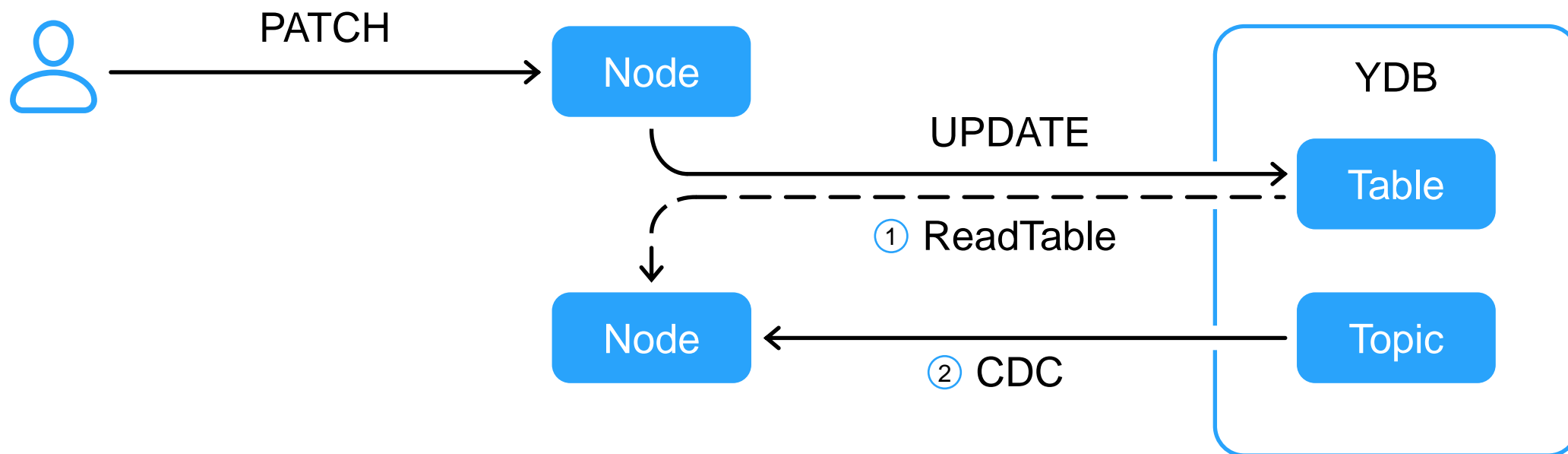
gRPC stream



YDB ReadTable
click.ru/34ccJd

Схема YDB CDC

56



Первая загрузка таблицы

57

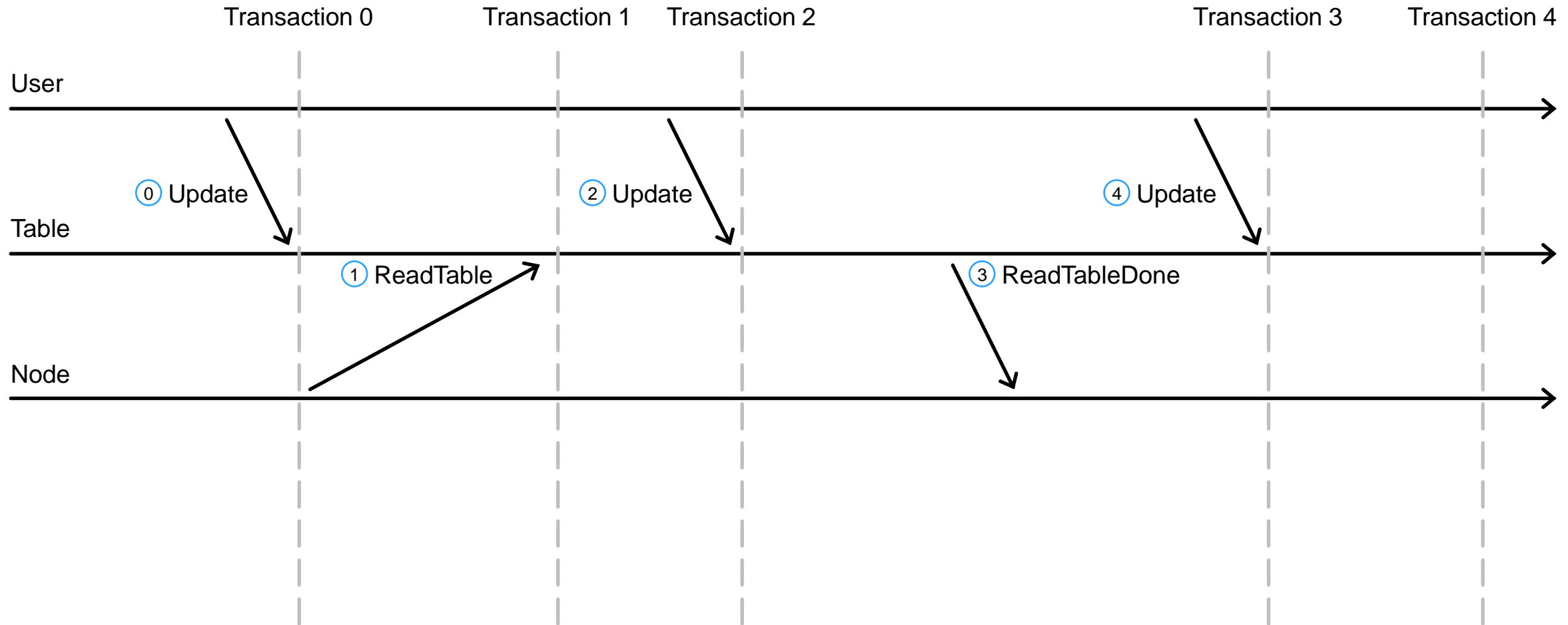
- Используем YDB ReadTable
gRPC stream
- Скачиваем только на старте



YDB ReadTable
click.ru/34ccJd

Первая загрузка таблицы

58



Первая загрузка таблицы

59

- ReadTable читает из снимка,
у которого есть virtual timestamp



virtual timestamp
clck.ru/34d94Q

Первая загрузка таблицы

60

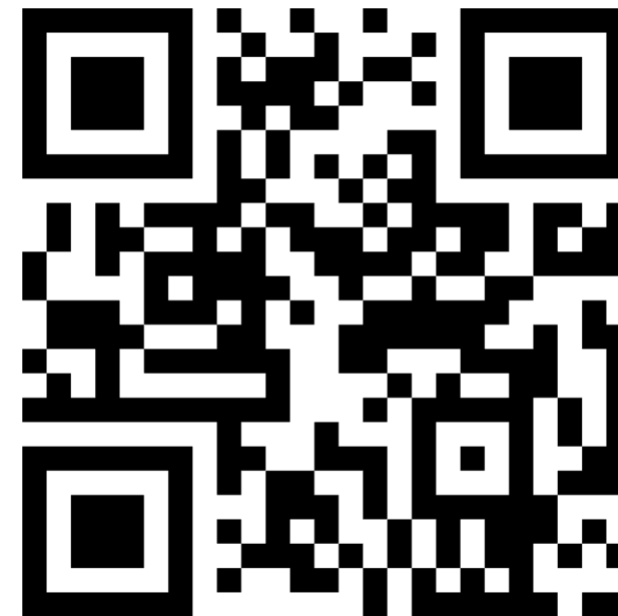
- ReadTable читает из снимка,
у которого есть virtual timestamp
- Сообщения тоже содержат
virtual timestamp



virtual timestamp
clck.ru/34d94Q

Первая загрузка таблицы

- ReadTable читает из снэпшота, у которого есть virtual timestamp
- Сообщения тоже содержат virtual timestamp
- virtual timestamp:
Coordinator time, transaction id



virtual timestamp
clck.ru/34d94Q

Первая загрузка таблицы

- ReadTable читает из снэпшота, у которого есть virtual timestamp
- Сообщения тоже содержат virtual timestamp
- virtual timestamp:
Coordinator time, transaction id
- Можно мержить, фильтровать, сортировать строки



virtual timestamp
clck.ru/34d94Q

Пример virtual timestamp

```
{
  "update": {},
  "newImage": {
    "createdAt": 1654947225226,
    "updatedAt": 1654947225226,
    "maxResponseSizeBytes": 10485760,
    "state": "RW",
    "maxFileSensors": 1000000,
    "protoQuotas": "KGQwsQI=",
    ...,
    "maxMemSensors": 500000,
  },
  "key": [
    "composite_key_part1",
    "composite_key_part2"
  ],
  "ts": [1670792400, 562949953607163]
}
```

Получение CDC после ReadTable



Получение CDC после ReadTable

65

- У очереди нельзя запросить сообщения по virtual timestamp 😬

Получение CDC после ReadTable

66

- У очереди нельзя запросить сообщения по virtual timestamp 😬
- Поэтому запрашиваем по времени снэпшота минус лаг

Получение CDC после ReadTable

67

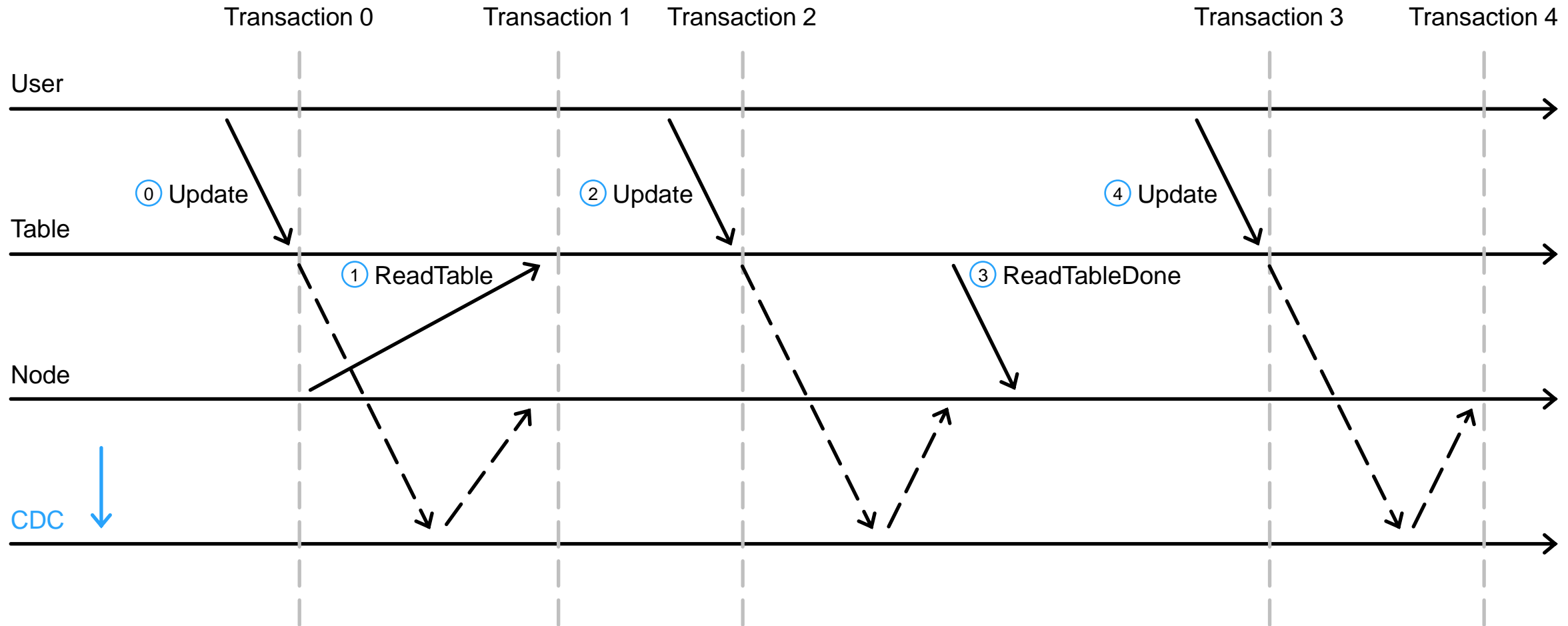


Схема таблиц

68

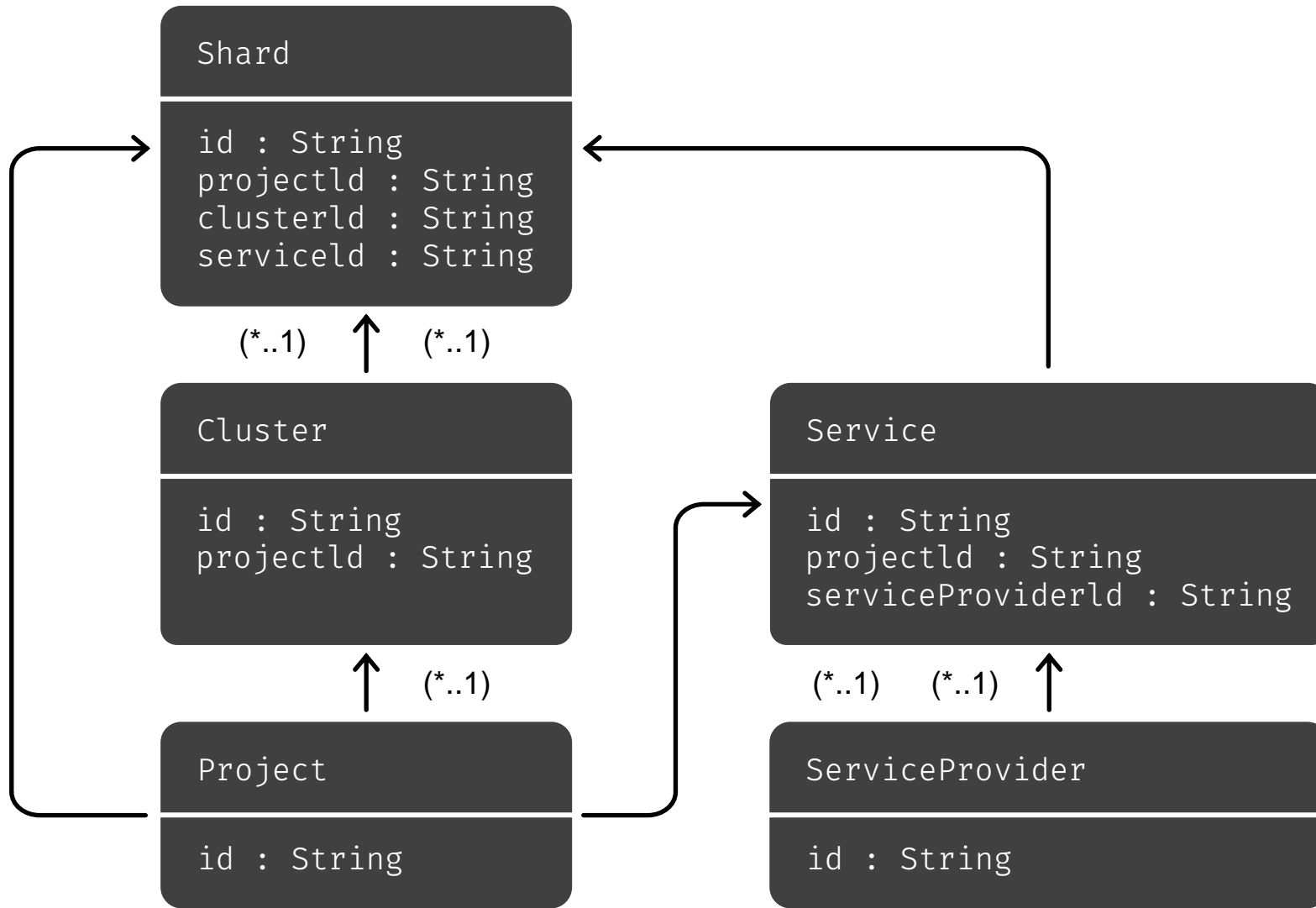
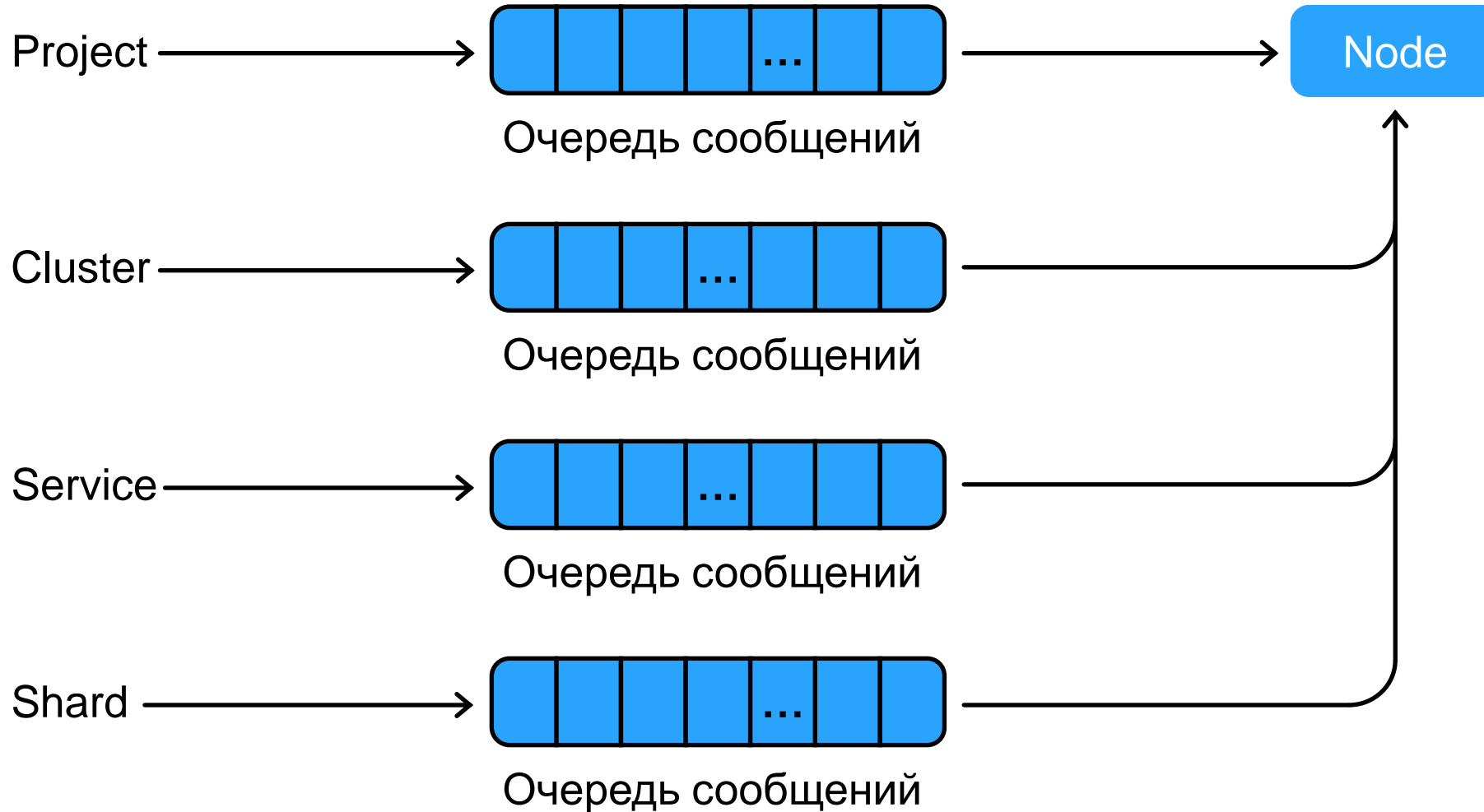


Схема очередей



Целостность без транзакций



Целостность без транзакций

71

- Считываем изменения из CDC

Целостность без транзакций

- Считываем изменения из CDC
- На лету валидируем все связи

Целостность без транзакций

- Считываем изменения из CDC
- На лету валидируем все связи
- Если состояние невалидно — не используем и ждём

Целостность без транзакций

- Считываем изменения из CDC
- На лету валидируем все связи
- Если состояние невалидно — не используем и ждём
- Если не дождались — перечитываем запросом

Схема YDB CDC

75

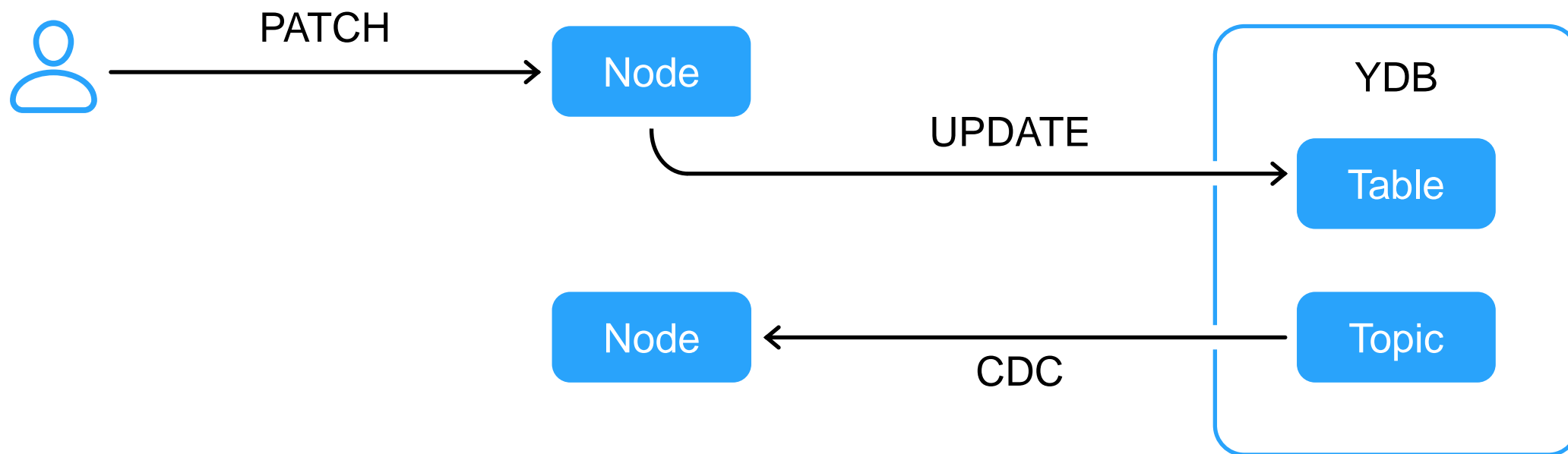
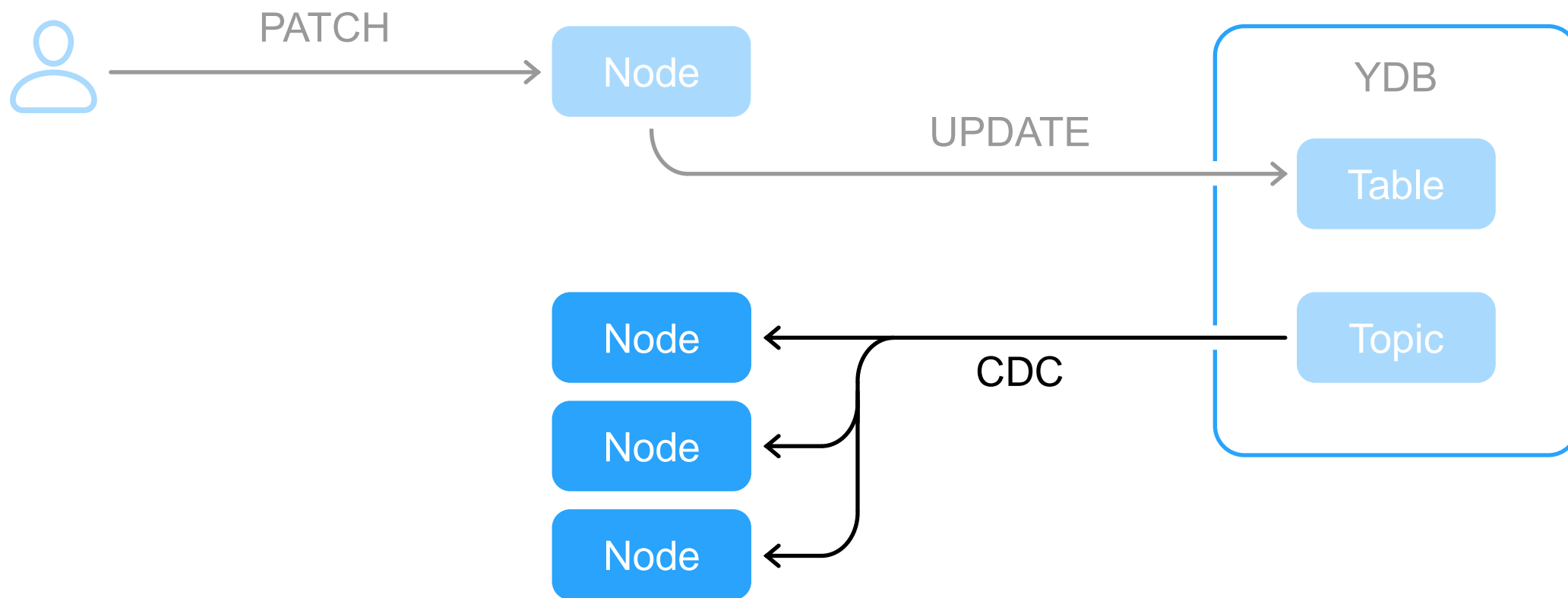


Схема YDB CDC

76



Лимиты YDB CDC по умолчанию



Лимиты YDB CDC по умолчанию

78

2 MB

max per sec

на подписчика

Лимиты YDB CDC по умолчанию

79

2 MB

max per sec

на подписчика

10

max

подписчиков на очередь

Лимиты YDB CDC по умолчанию

80

2 MB

max per sec

на подписчика

10

max

подписчиков на очередь

5

max

очередей на таблицу

Лимиты YDB CDC по умолчанию

81

2 MB

max per sec

на подписчика

10

max

подписчиков на очередь

5

max

очередей на таблицу

100 MB

max per sec

итого на таблицу

Как не надо было делать



Как не надо было делать

- Обновления редкие

Как не надо было делать

- Обновления редкие
- Хотим подписать на очередь 1500 нод

Как не надо было делать

85

- Обновления редкие
- Хотим подписать на очередь 1500 нод
- Готово 😎

Как не надо было делать

86

- Обновления редкие
- Хотим подписать на очередь 1500 нод
- Готово 😎
- Надо было: зеркалировать или повышать лимиты



Повышать лимиты
clck.ru/34d9y6

Как не надо было делать

87

- Обновления редкие
- Хотим подписать на очередь 1500 нод
- Готово 😎
- Надо было: зеркалировать или повышать лимиты
- На самом деле наш сценарий уже поддерживается



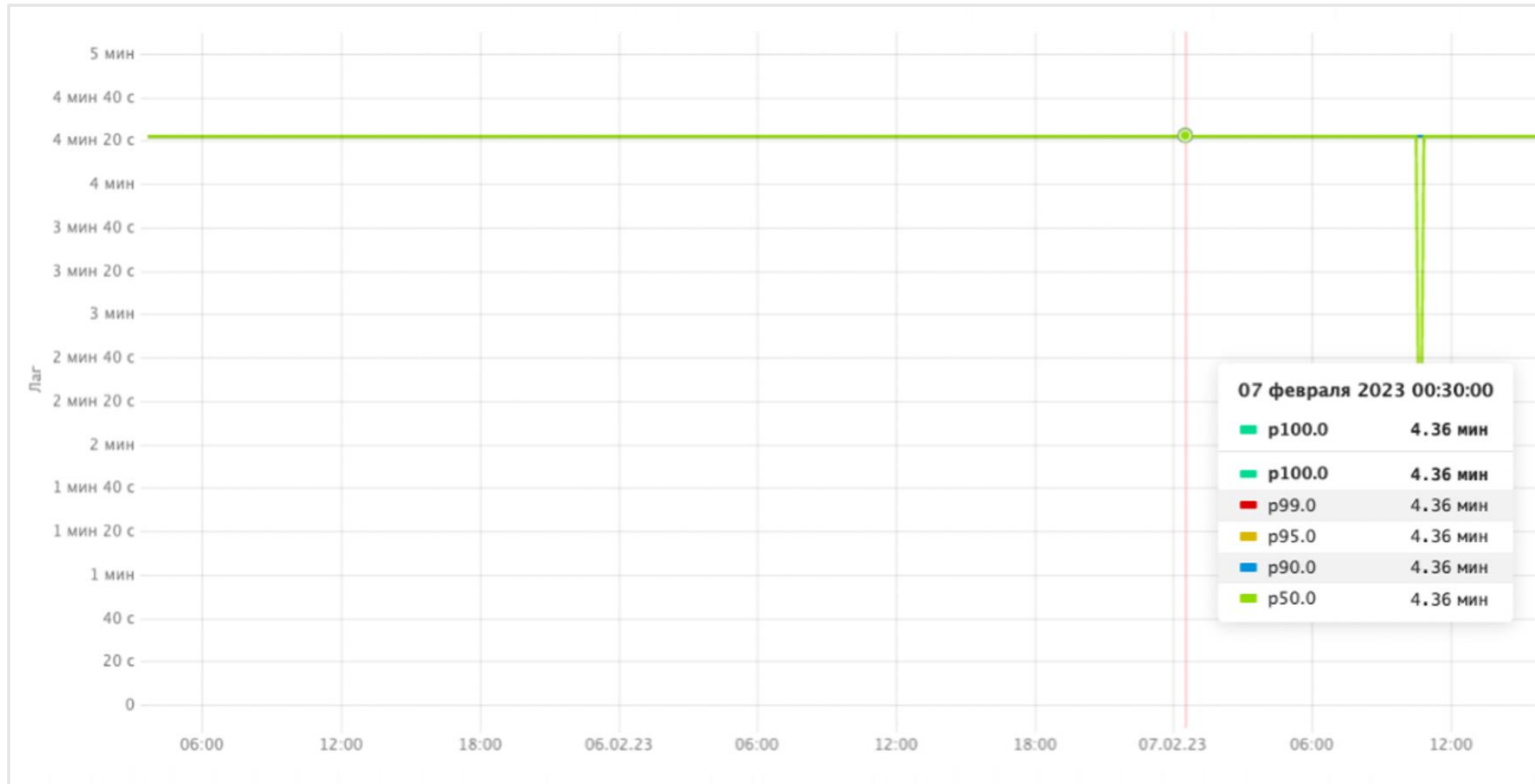
Повышать лимиты
clck.ru/34d9y6

План

1. Теория
2. Практика использования CDC
3. Что пошло не так
4. Итоги

Проблема: лаг доставки 4 минуты

89



Решение: лаг доставки 4 минуты

90

- Перезагружаем ноды

Решение: лаг доставки 4 минуты

- Перезагружаем ноды
- Они загружают первый раз таблицы

Решение: лаг доставки 4 минуты

- Перезагружаем ноды
- Они загружают первый раз таблицы
- Это долго, отсюда большие персентили вверху

**Решение: лаг доставки
4 минуты**



Решение: лаг доставки 4 минуты

94

Смотрим
на лаги доставки
в очереди YDB
и загрузку



Загрузка

clck.ru/34dAM7

Решение: лаг доставки 4 минуты

95

Смотрим
на лаги доставки
в очереди YDB
и загрузку



Загрузка
clck.ru/34dAM7

Тяжёлые метрики,
когда много
подписчиков

Исправлено в Upstream



Upstream
clck.ru/34dANP

Решение: лаг доставки 4 минуты

96

Смотрим
на лаги доставки
в очереди YDB
и загрузку



Загрузка
clck.ru/34dAM7

Тяжёлые метрики,
когда много
подписчиков
[Исправлено в Upstream](#)



Upstream
clck.ru/34dANP

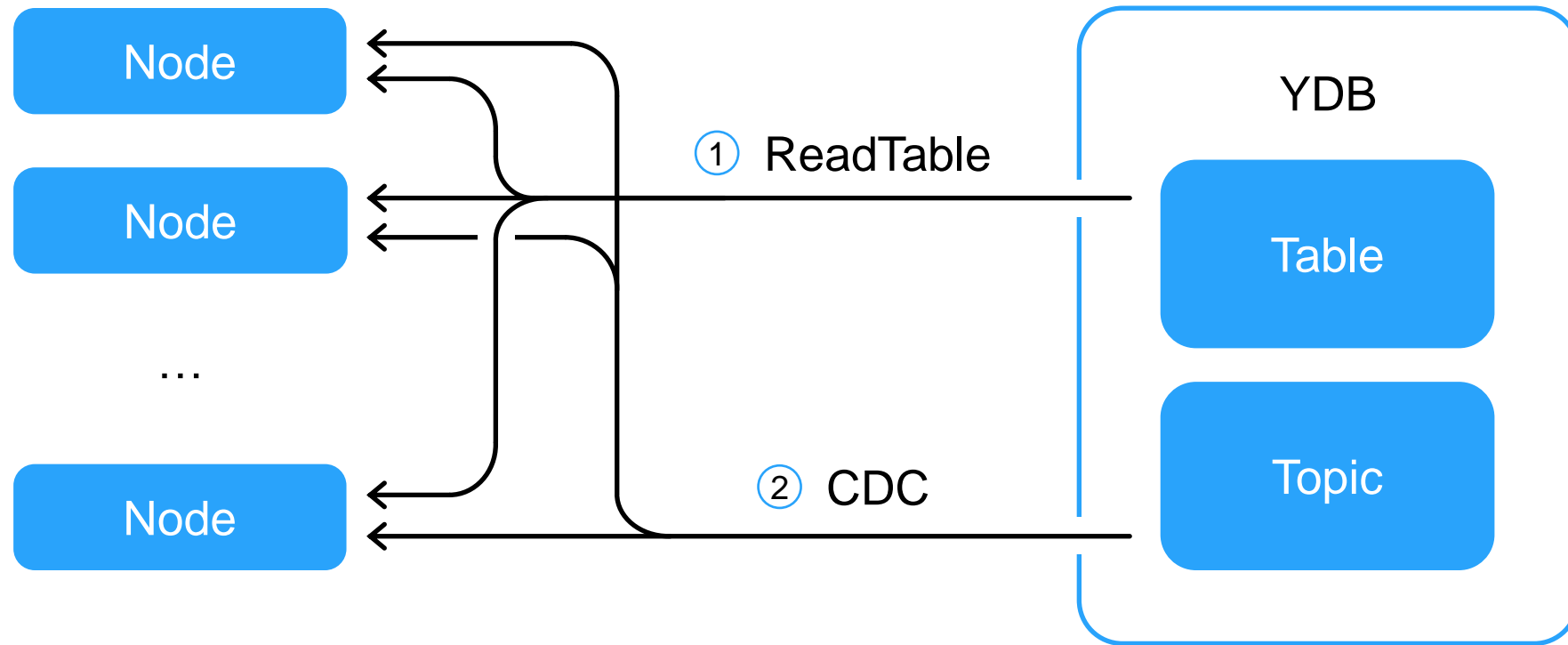
Тяжёлые
SQL-запросы
[Оптимизировать](#)
[Добавлять ресурсы](#)



Performance
clck.ru/34hXJg

Проблема: массовый рестарт

97



Решение: массовый рестарт

1

Повторяем упавшие
ReadTable на клиентах

Решение: массовый рестарт

1

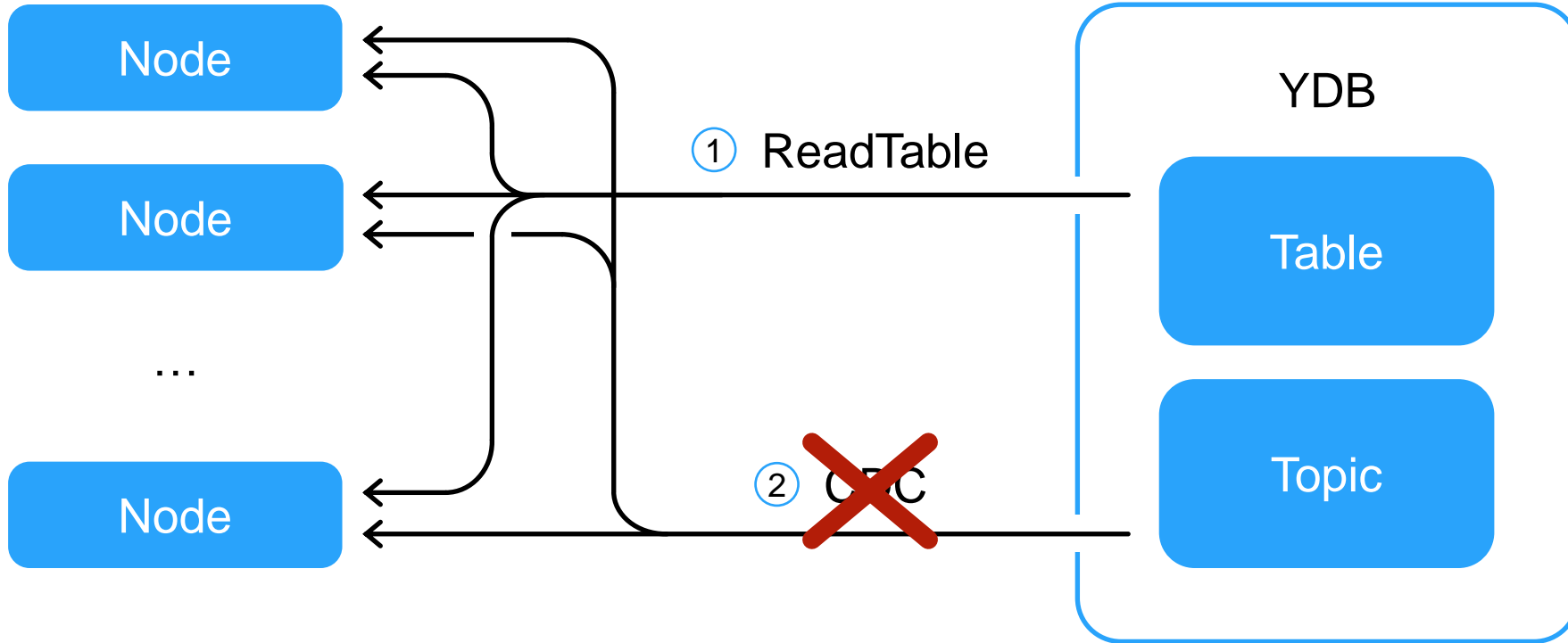
Повторяем упавшие
ReadTable на клиентах

2

Distributed semaphore
на клиентах

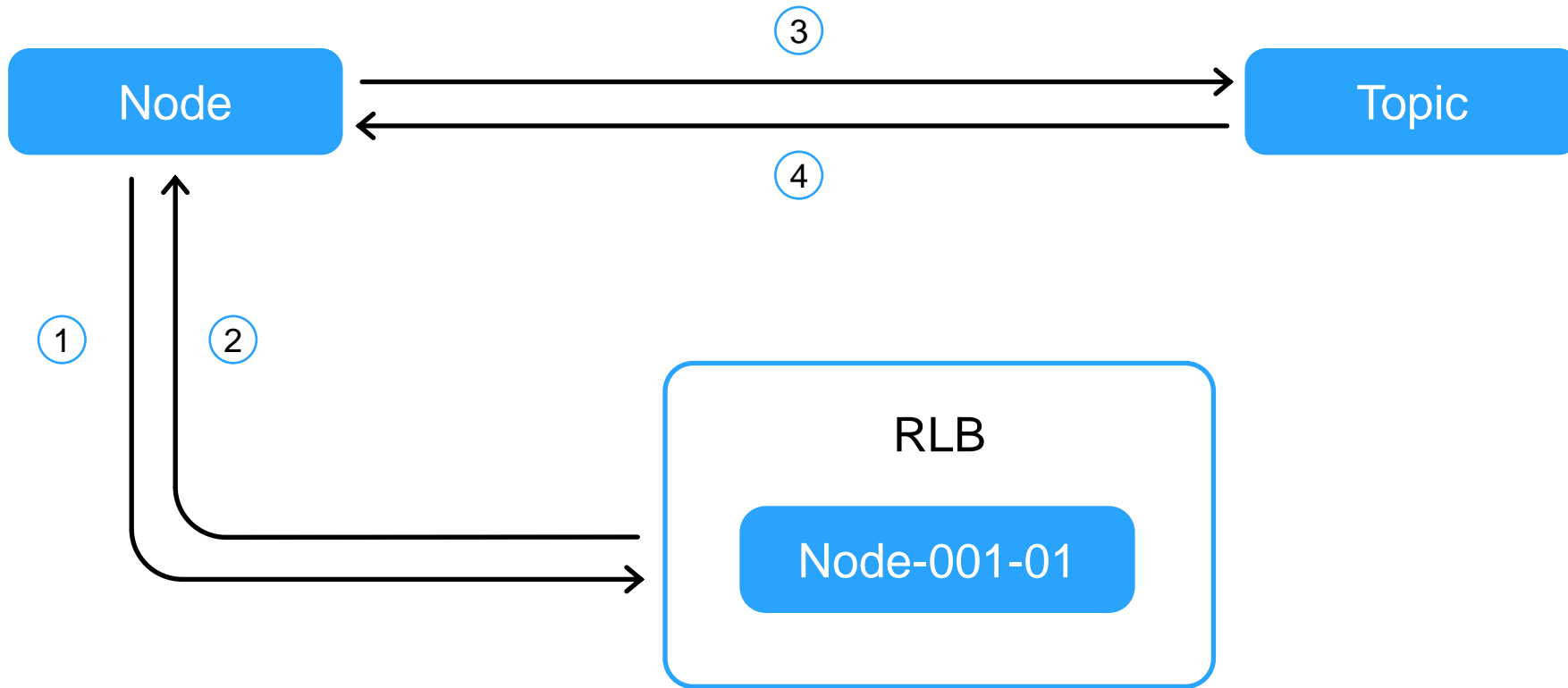
Проблема: рестарт сервиса

100



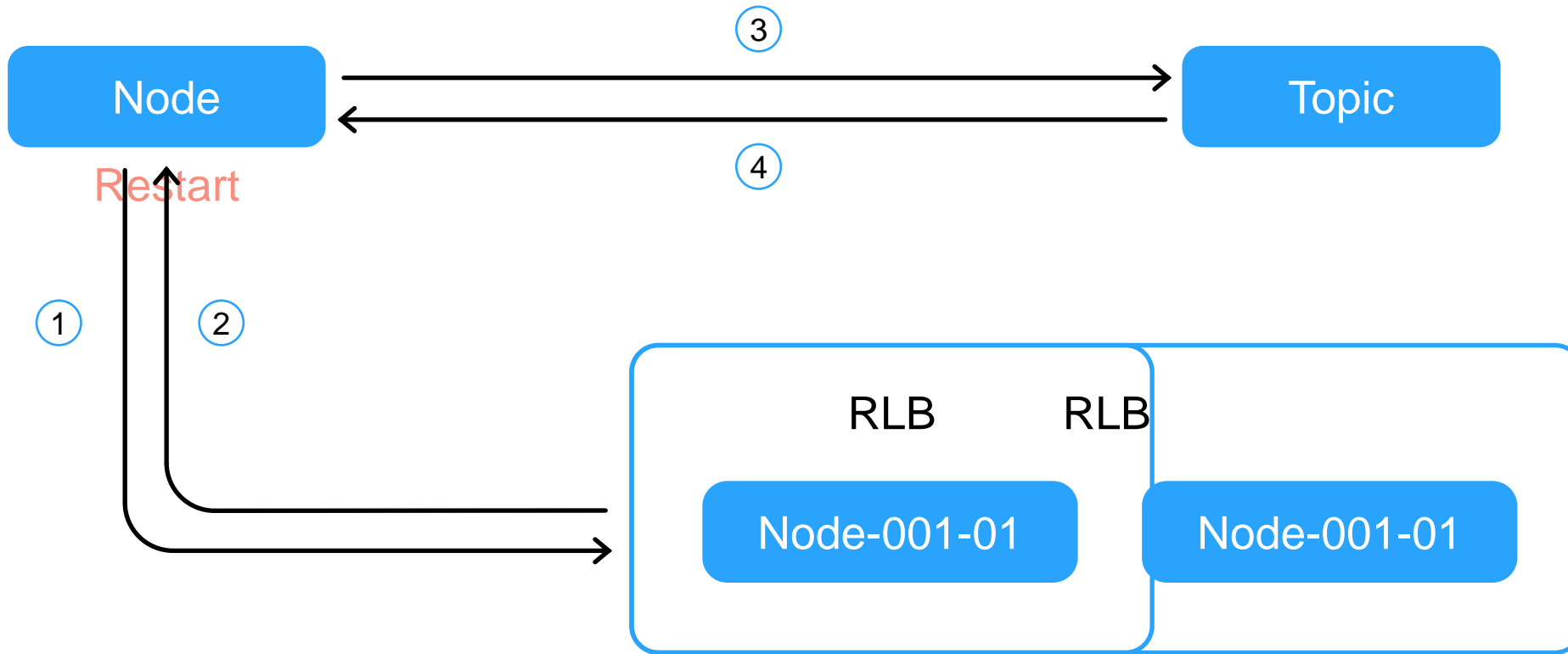
Проблема: рестарт сервиса

101



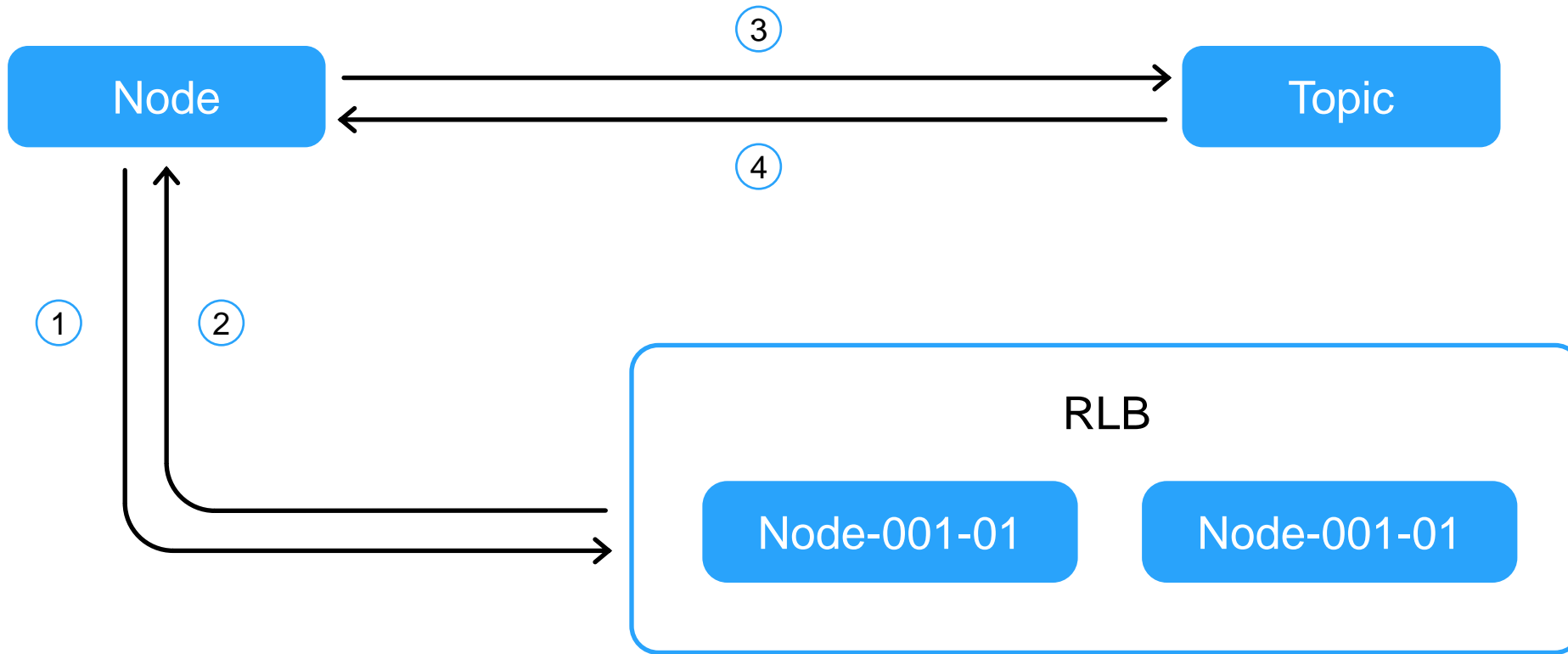
Проблема: рестарт сервиса

102



Проблема: рестарт сервиса

103

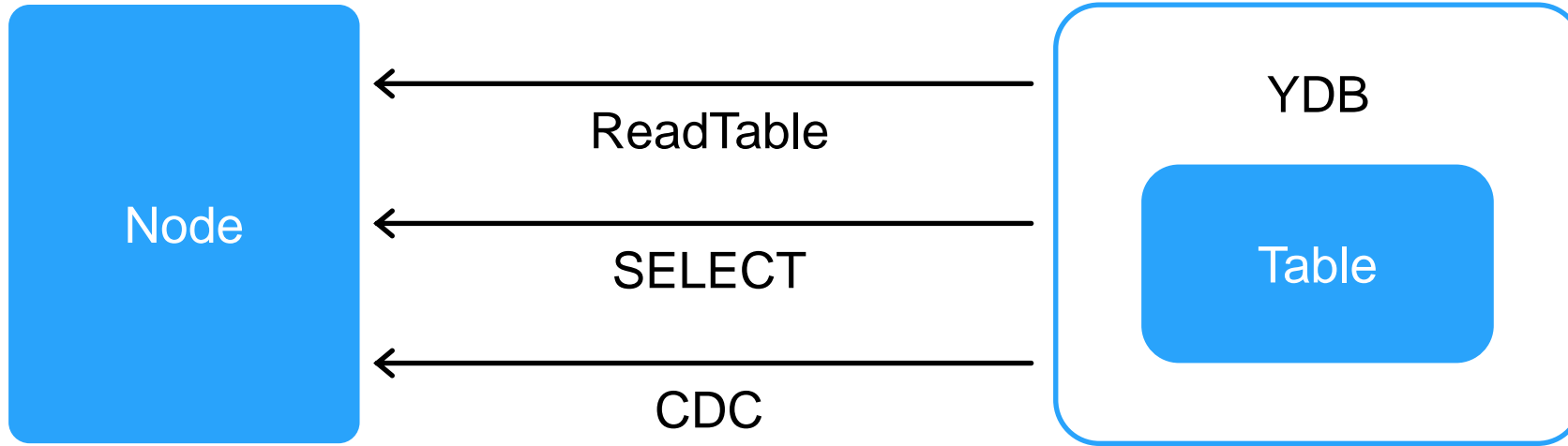


Решение: рестарт сервиса

Client GRPC keep alive/idle timeout —
помогают умирать старой сессии

Проблема: методы чтения

105



Решение: методы чтения



Решение: методы чтения

- Round-trip tests

Подняли локально YDB

Записали сущность

Прочитали через CDC

Проверили, что все поля заполнились корректно

Решение: методы чтения

- Round-trip tests

Подняли локально YDB

Записали сущность

Прочитали через CDC

Проверили, что все поля заполнились корректно

- Унификация кода, читающего данные

Решение: методы чтения

- Round-trip tests

Подняли локально YDB

Записали сущность

Прочитали через CDC

Проверили, что все поля заполнились корректно

- Унификация кода, читающего данные
- Флаг, включающий один метод

У нас ReadTable

Флаг на ReadTable



Флаг на ReadTable

- Когда всё хорошо, конфиги совпадают с небольшим лагом

Флаг на ReadTable

- Когда всё хорошо, конфиги совпадают с небольшим лагом
- Если найдено расхождение, то ставится флаг

Флаг на ReadTable

- Когда всё хорошо, конфиги совпадают с небольшим лагом
- Если найдено расхождение, то ставится флаг
- Расхождение может возникнуть

Баги в коде при выкатке новых версий

Флаг на ReadTable

- Когда всё хорошо, конфиги совпадают с небольшим лагом
- Если найдено расхождение, то ставится флаг
- Расхождение может возникнуть

Баги в коде при выкатке новых версий

Потеря сетевого соединения и т. п.

Флаг на ReadTable

115

- Когда всё хорошо, конфиги совпадают с небольшим лагом
- Если найдено расхождение, то ставится флаг
- Расхождение может возникнуть

Баги в коде при выкатке новых версий

Потеря сетевого соединения и т. п.

Космическое излучение меняет битик 🙄

Как всё это мониторить

116

- Метрики сервиса

Как всё это мониторить

117

- Метрики сервиса

Задержка доставки

Как всё это мониторить

- Метрики сервиса

Задержка доставки

Время последнего события
(cross dc, cross nodes)

Как всё это мониторить

- Метрики сервиса

Задержка доставки

Время последнего события
(cross dc, cross nodes)

Объёмы чтения и спайки

Как всё это мониторить

120

- Метрики сервиса

Задержка доставки

Время последнего события
(cross dc, cross nodes)

Объёмы чтения и спайки

Ошибки обработки

Как всё это мониторить

121

- Метрики сервиса

Задержка доставки

Время последнего события
(cross dc, cross nodes)

Объёмы чтения и спайки

Ошибки обработки

- Метрики YDB

Метрики топигов

Задержка чтения

Сообщения в очереди

Загрузку CPU



Метрики топигов

<https://clck.ru/34mn8M>



CPU

clck.ru/34dAM7

Как всё это мониторить

122

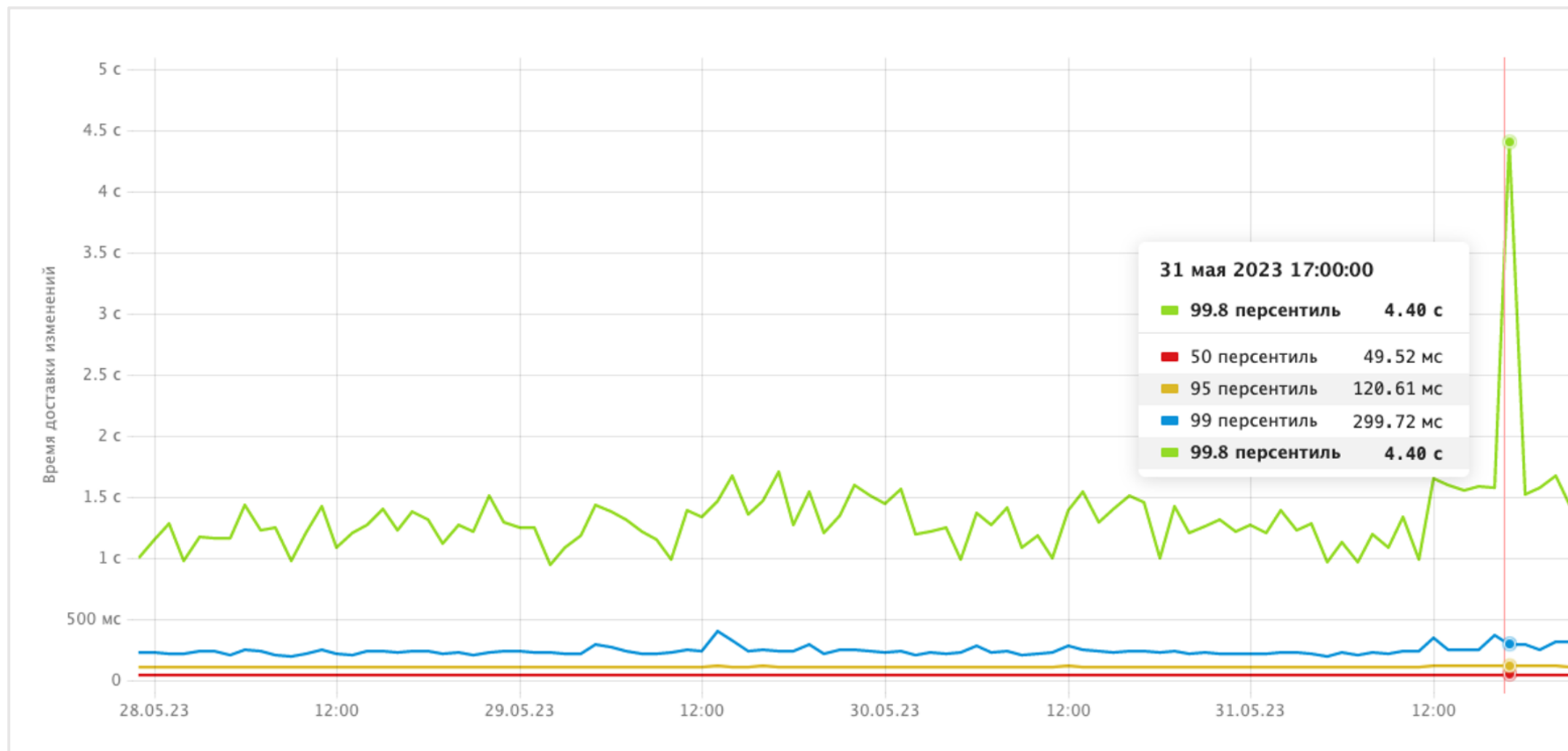


План

1. Теория
2. Практика использования CDC
3. Что пошло не так
4. Итоги

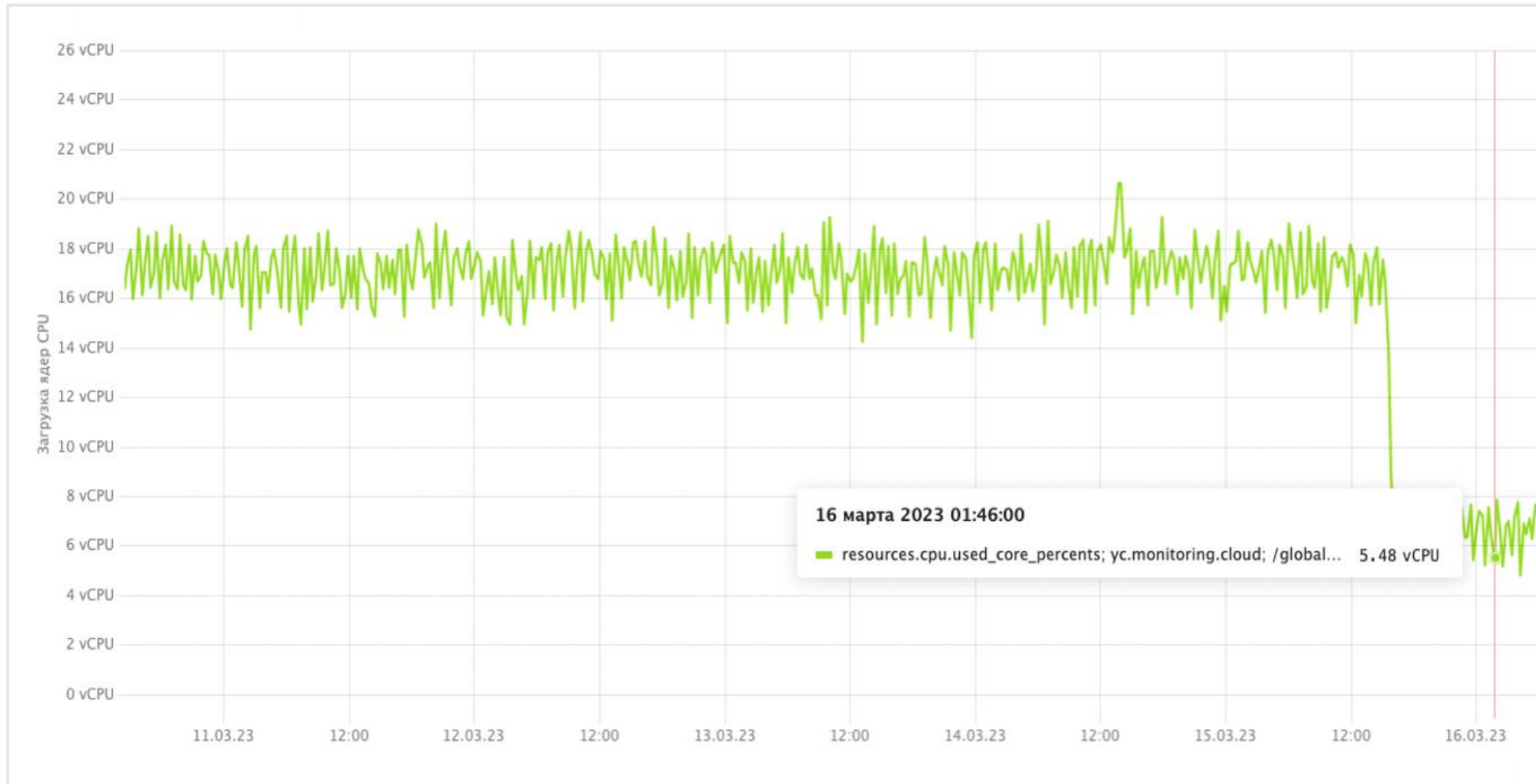
Какой стал лаг доставки?

124



Как упал CPU?

125



Наш CDC в цифрах

126

1530

подписчиков на очередь

~6 тыс.

сообщений в секунду

80 MB

в секунду в среднем

180 MB

в секунду в пике

<800

миллисекунд —
99 персентиль
лага доставки

В 2 раза

упала загрузка
ядер CPU

Дальнейшие планы

- Ещё больше таблиц переводить на CDC
- Протестировать CDC initial table scan
- Проверить возможности YDB SDK v2
- Планируется в CDC

Консистентность между очередями

Эфемерные подписчики



initial table scan
clck.ru/34dCo2



YDB SDK v2
clck.ru/34dCp7

Выводы

- Получили реализацию механизма Change Data Capture на YDB CDC, которую мы научились готовить
- Получили быструю доставку конфигураций на 1500 подписчиков
- Улучшили доступность мониторинга для наших пользователей

Голосуйте за мой доклад



Егор Литвиненко
Yandex Infrastructure

