

Просто о сложном: как работает драйвер распределенной базы данных YDB

Мясников Алексей,
старший разработчик в команде YDB



HighLoad⁺⁺
2022

Яндекс

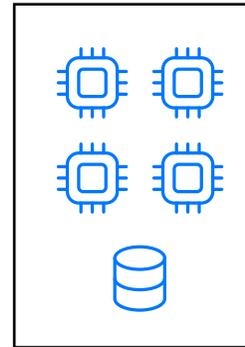
Содержание

- 01 Традиционные базы данных
- 02 YDB way
- 03 Требования к драйверу
- 04 YDB API
- 05 Жизненный цикл драйвера YDB

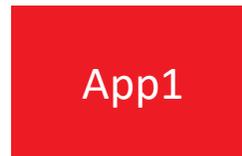
01

Традиционные базы данных

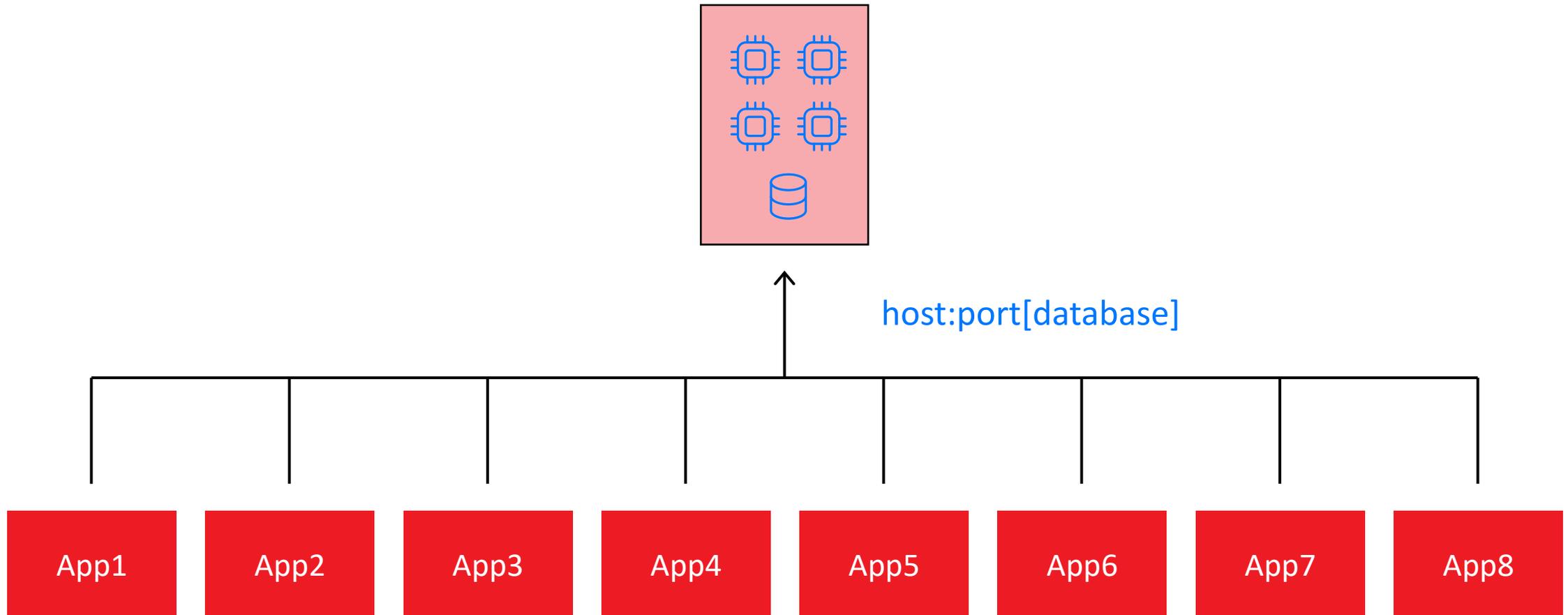
Маленькая база данных



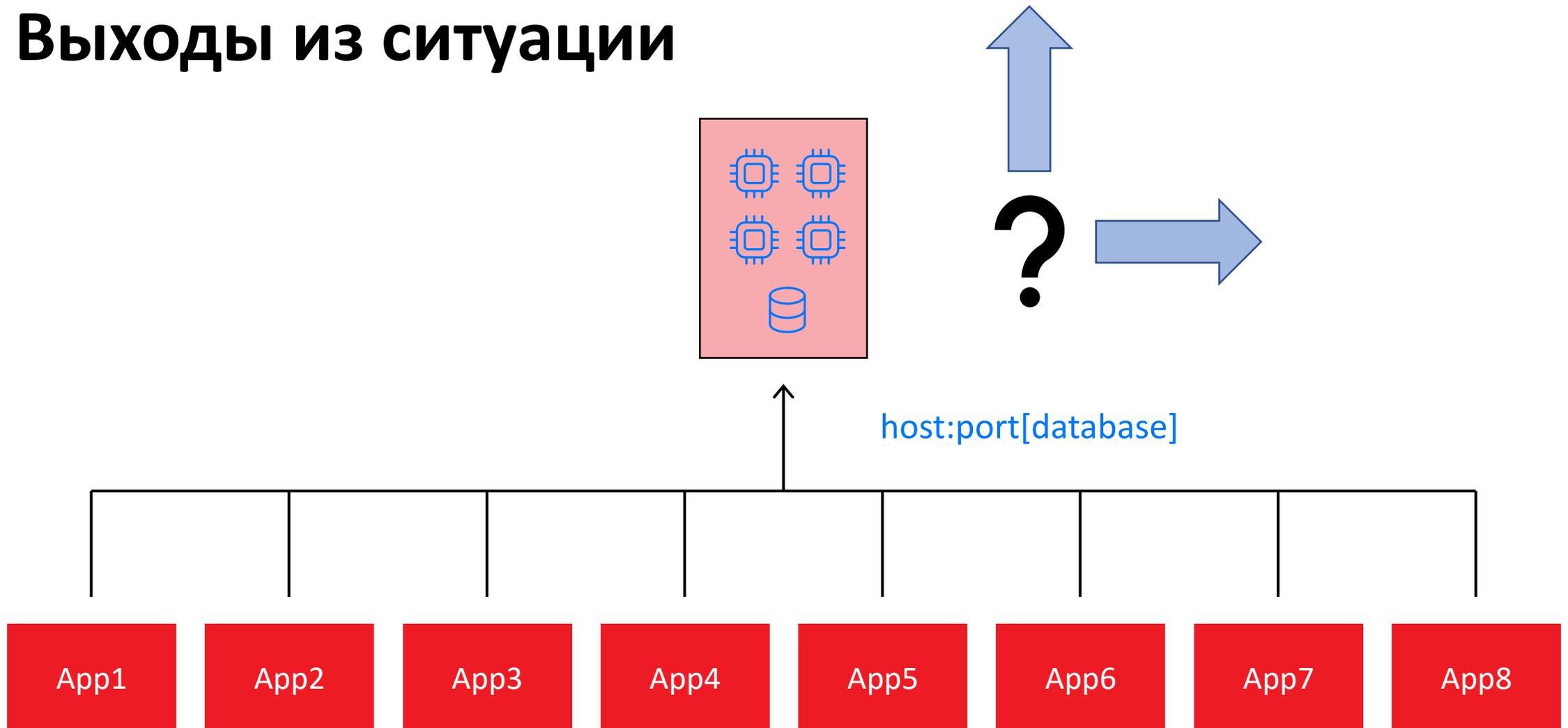
host:port[database]



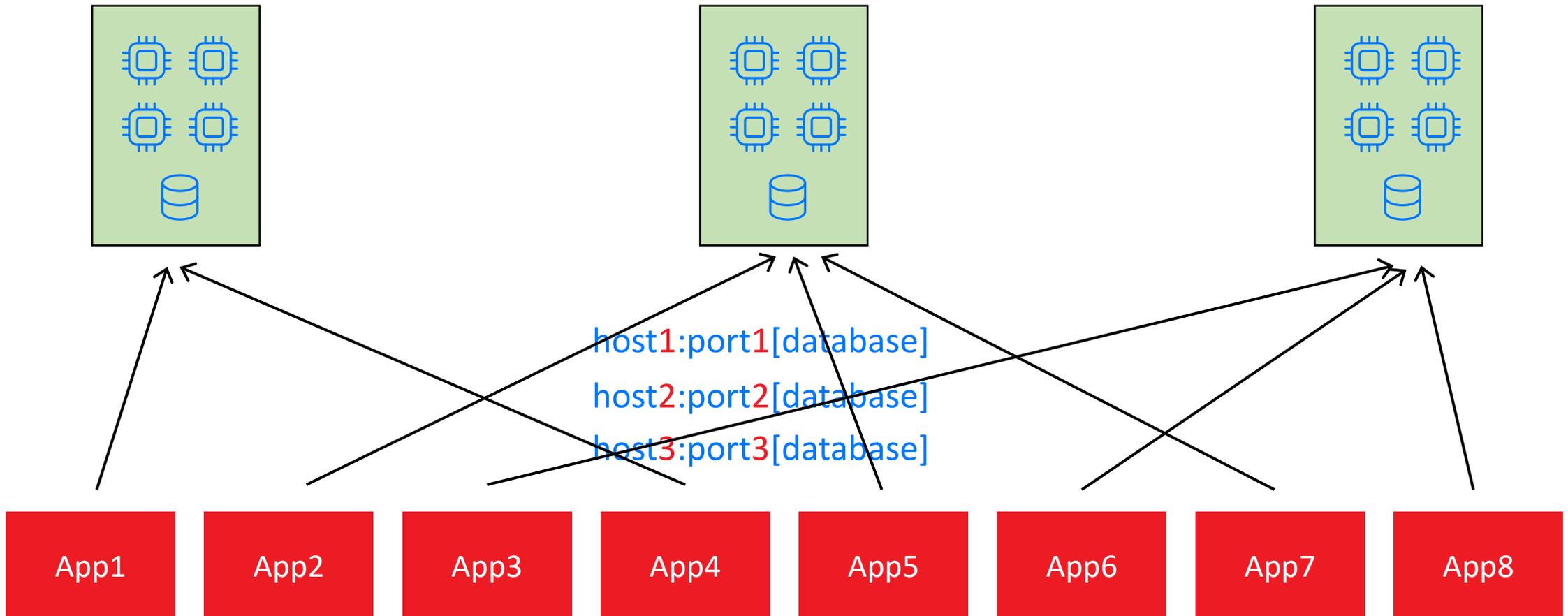
Растет нагрузка



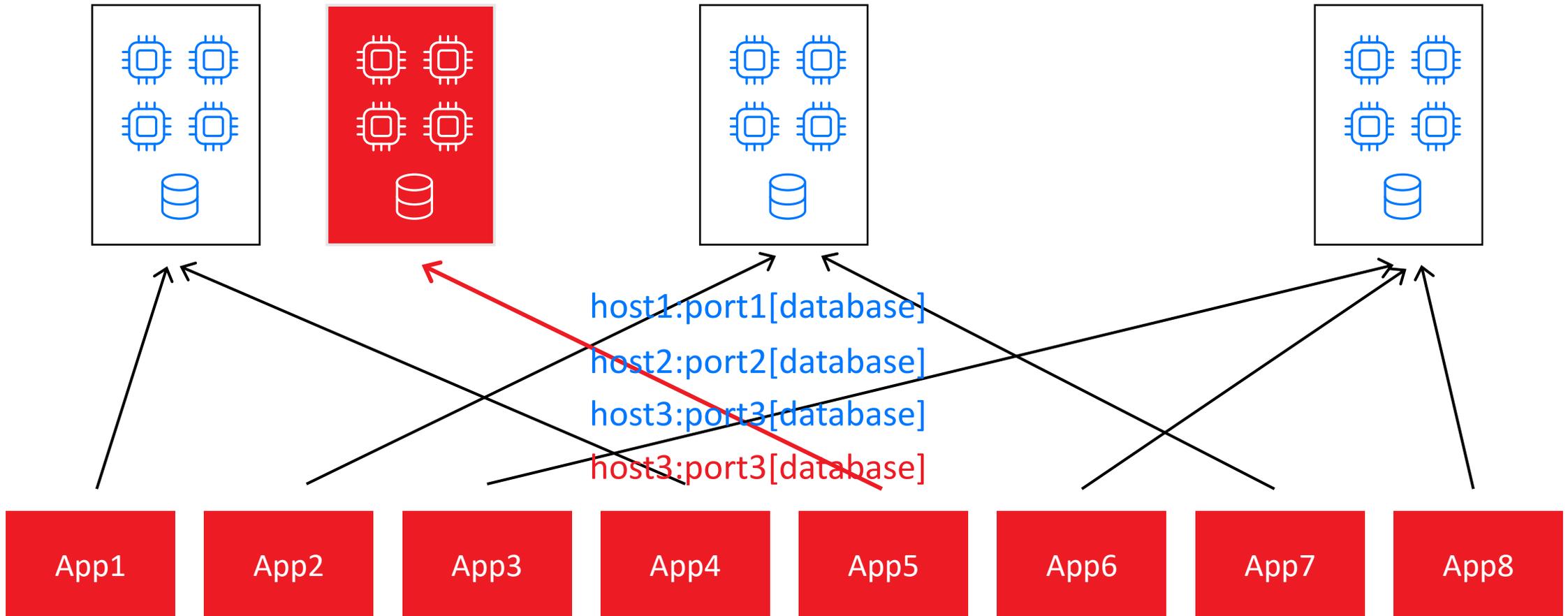
Выходы из ситуации



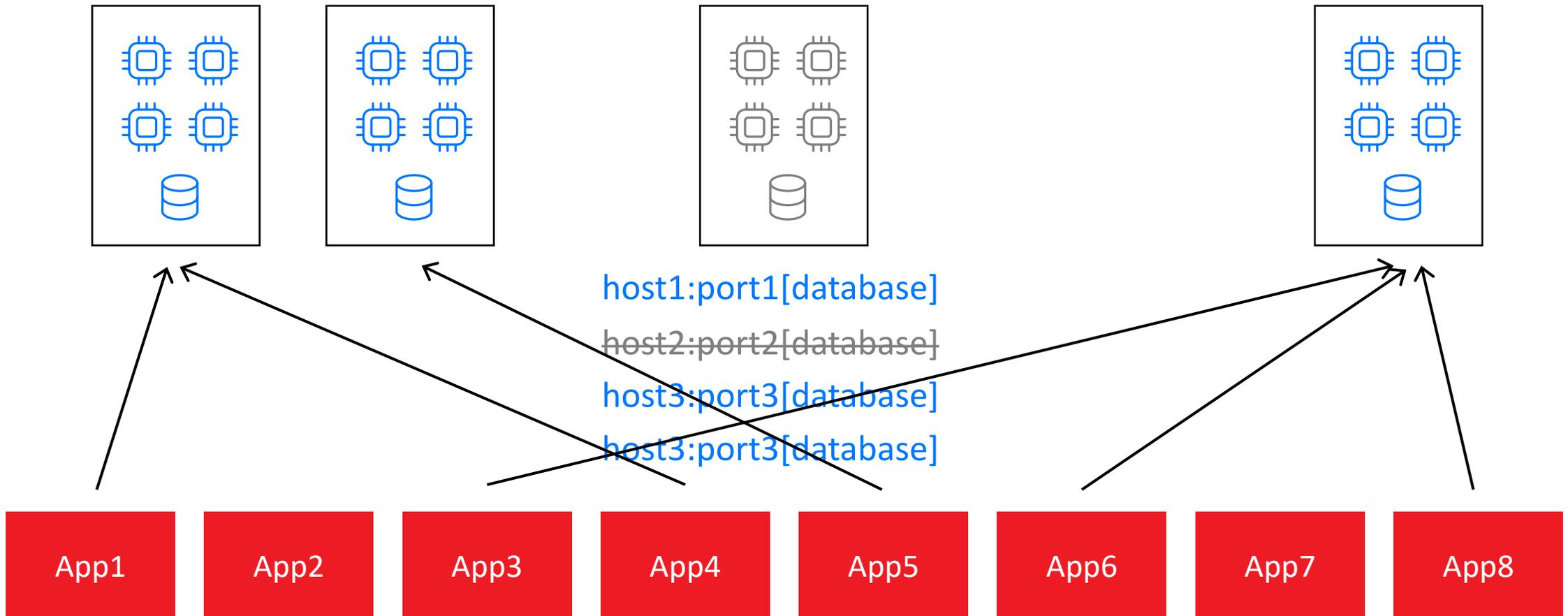
Хосты БД в конфиге



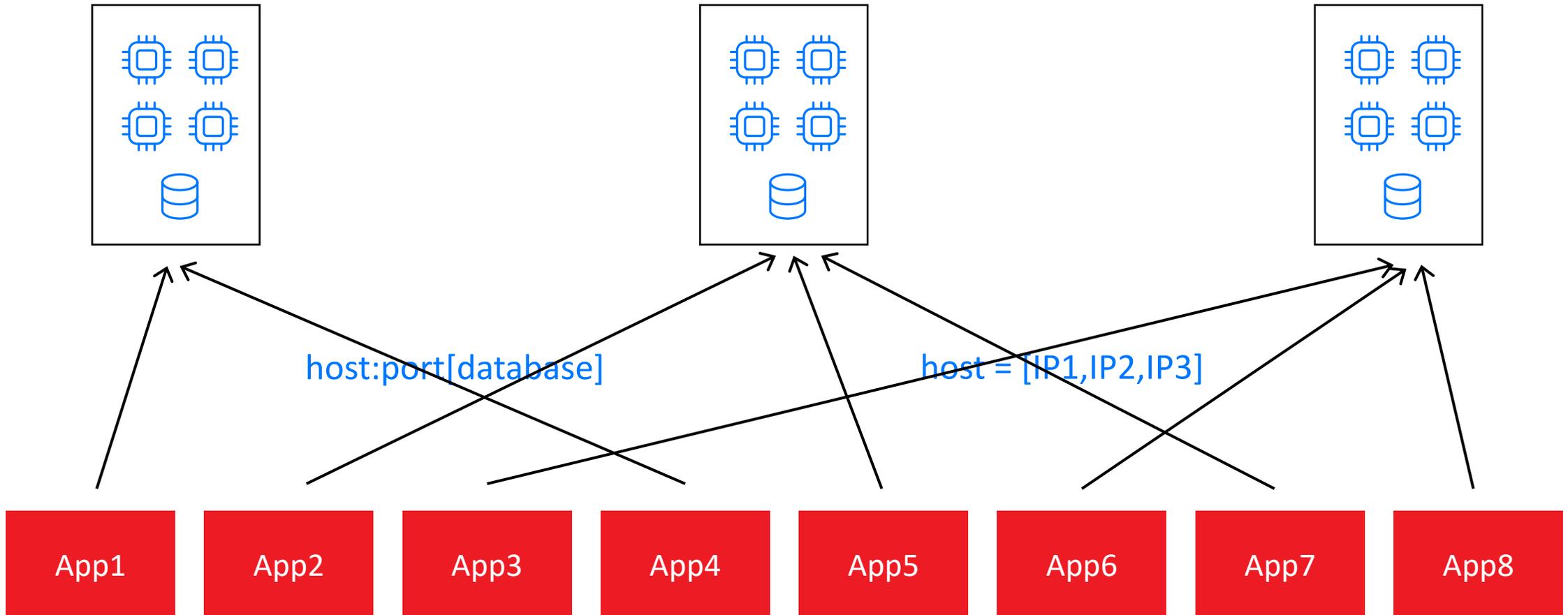
Хосты БД в конфиге: добавляем хост



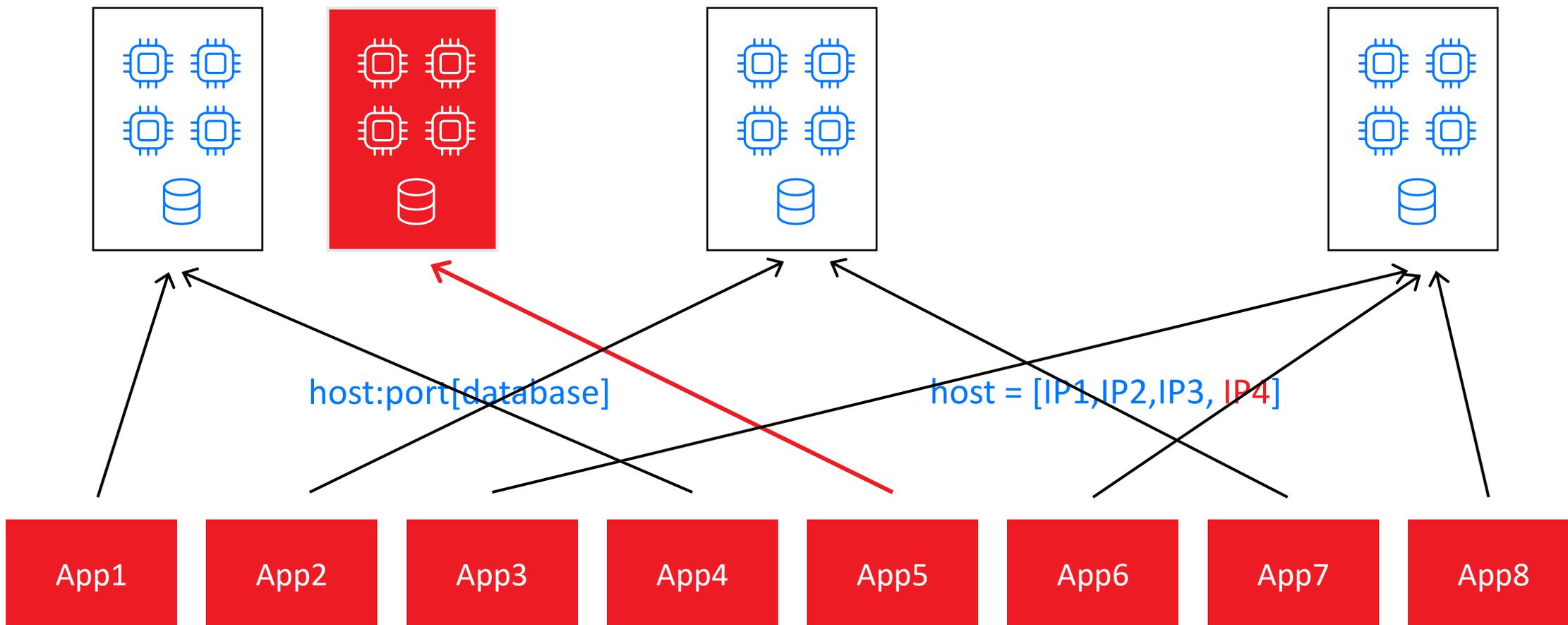
Хосты БД в конфиге: удаляем хост



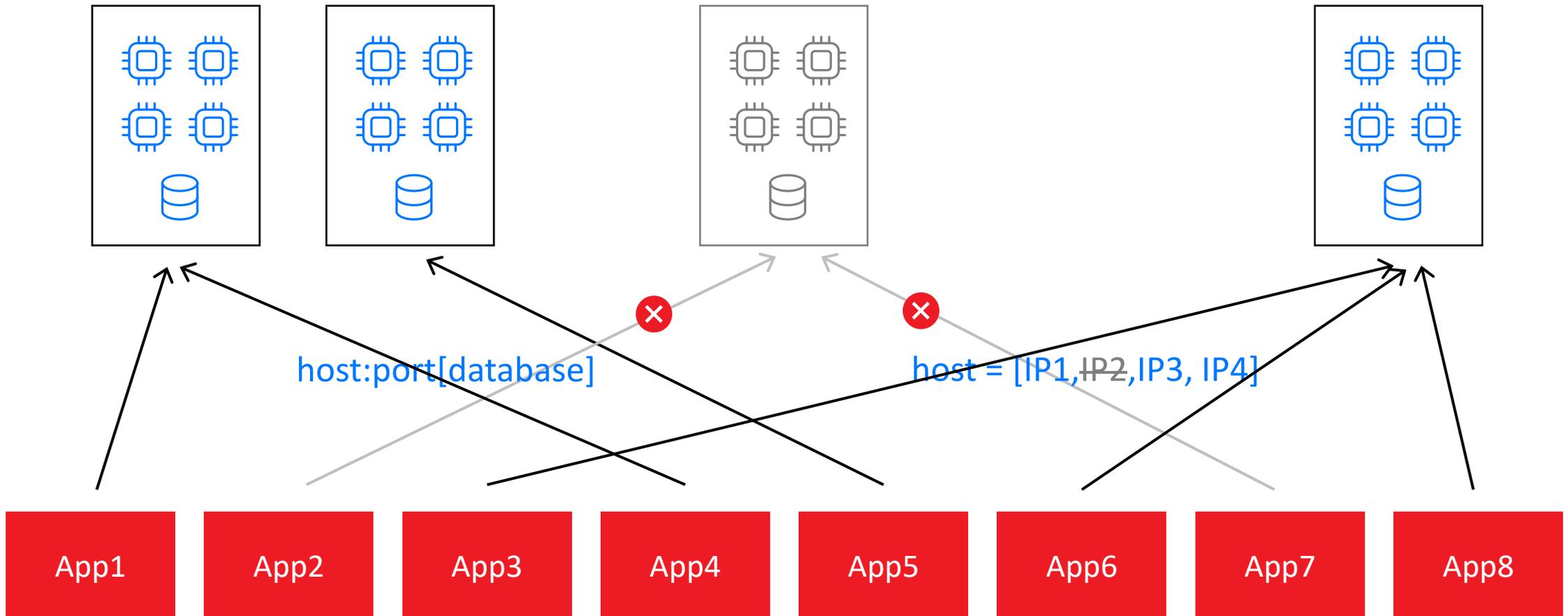
Хосты БД за DNS



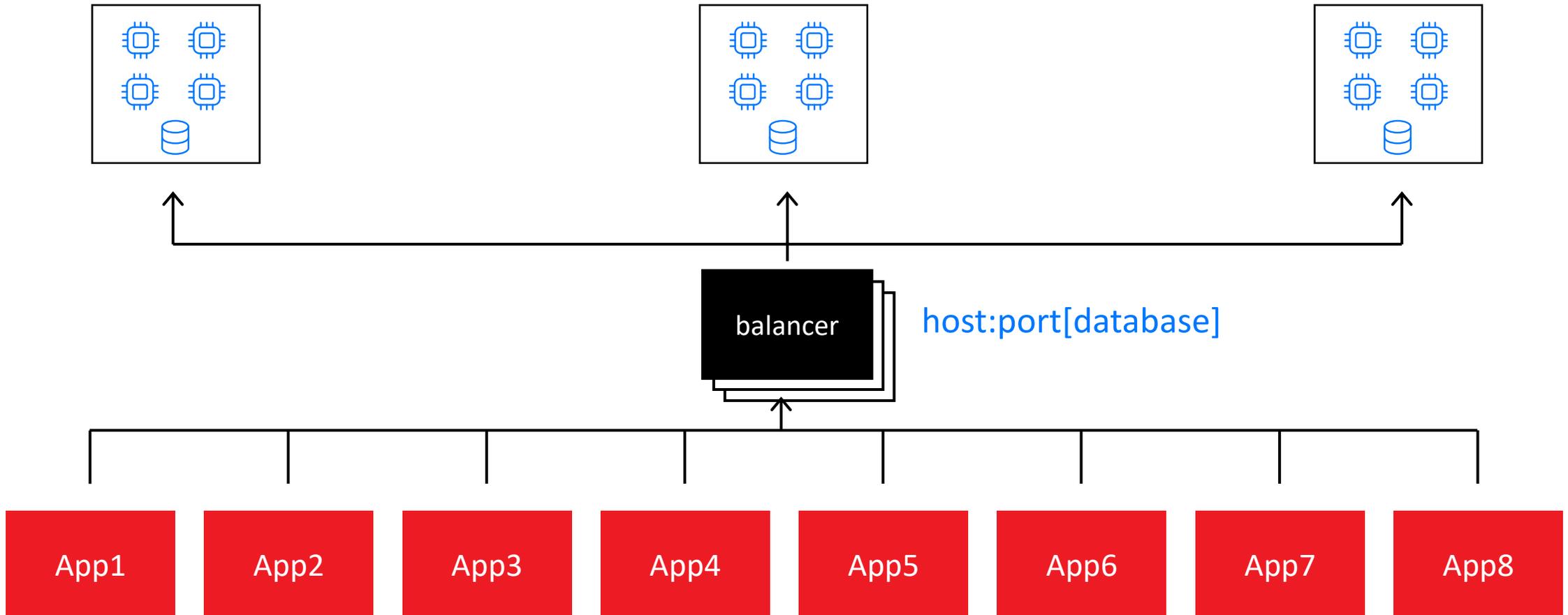
DNS: добавляем хост



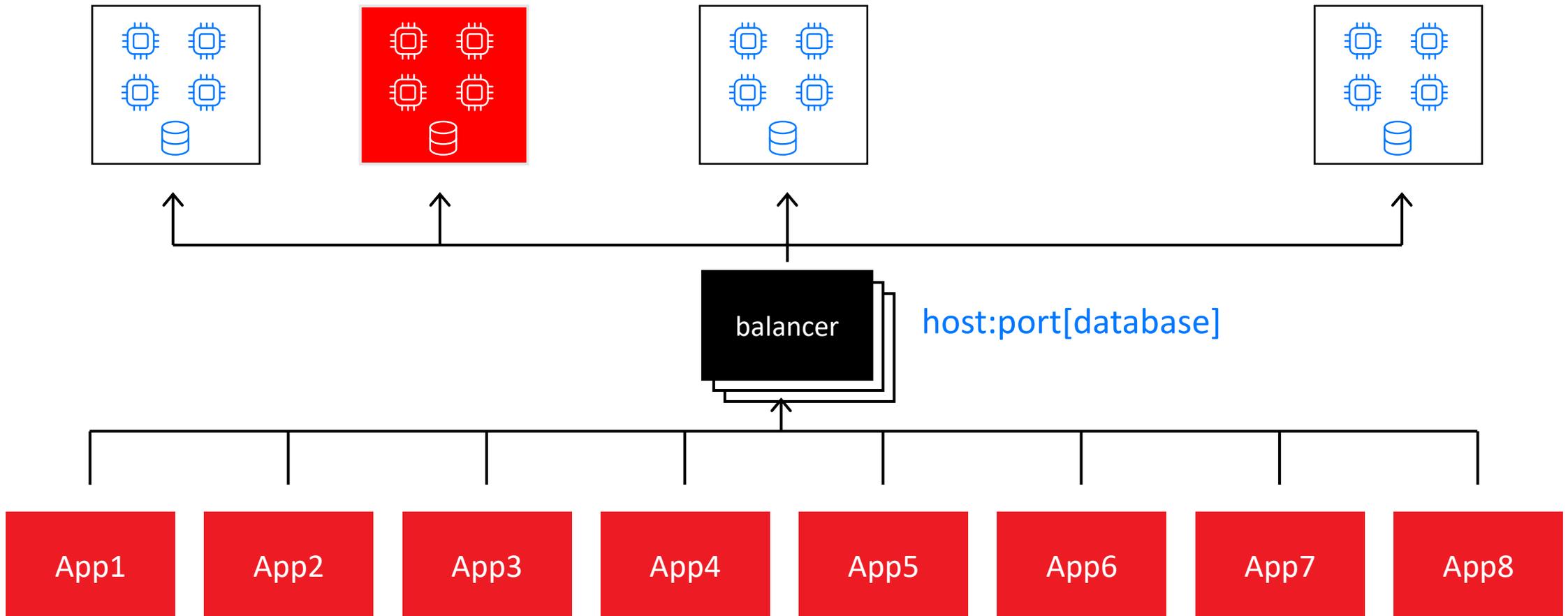
DNS: удаляем хост



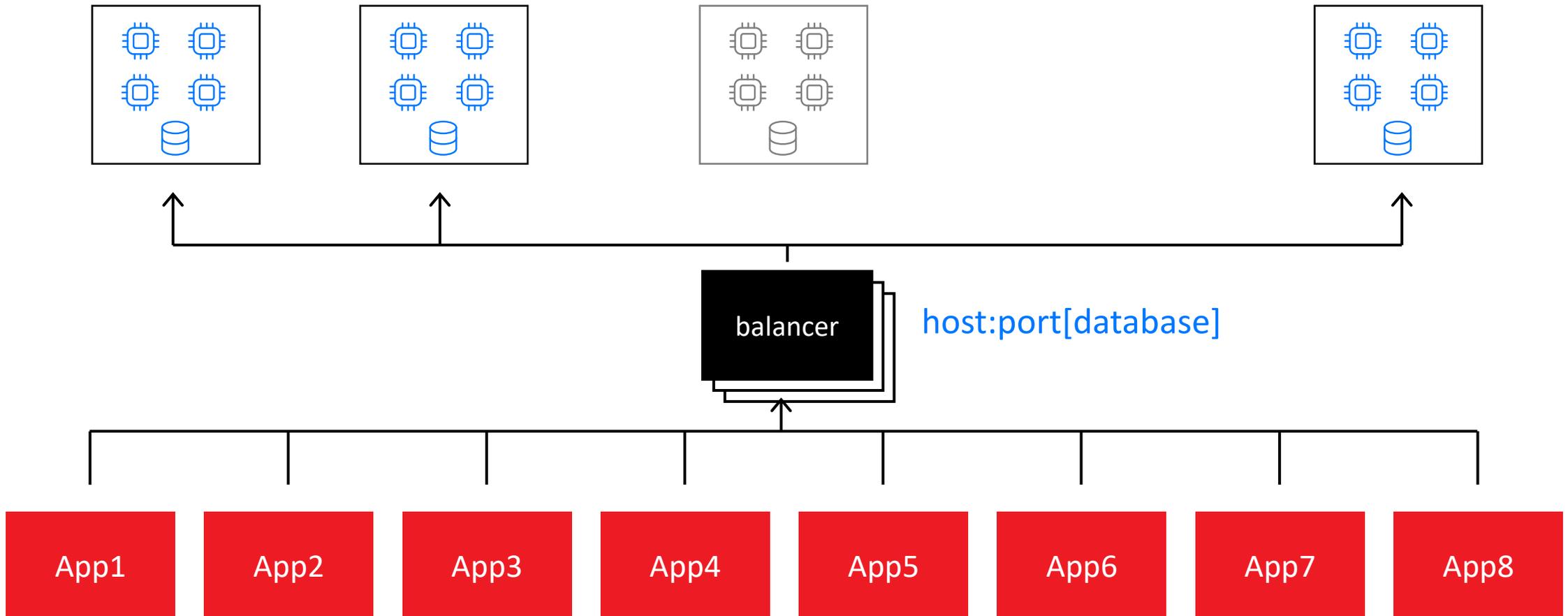
Хосты БД за балансером (L3, L7)



Хосты БД за балансером (L3, L7): добавляем хост



Хосты БД за балансером (L3, L7): удаляем хост



02

YDB way

YDB

Распределённая open-source SQL-база данных для операционных нагрузок

- + SQL для OLTP
- + Горизонтальное масштабирование
- + Транзакции с гарантиями ACID в нескольких AZ
- + Работоспособность и автоматическое восстановление при отказах
- + Масштабирование на миллионы транзакций в секунду и сотни терабайт данных
- + Поддерживаются OLAP-сценарии, координация распределённых систем (like ZooKeeper), доставка сообщений (like Kafka)

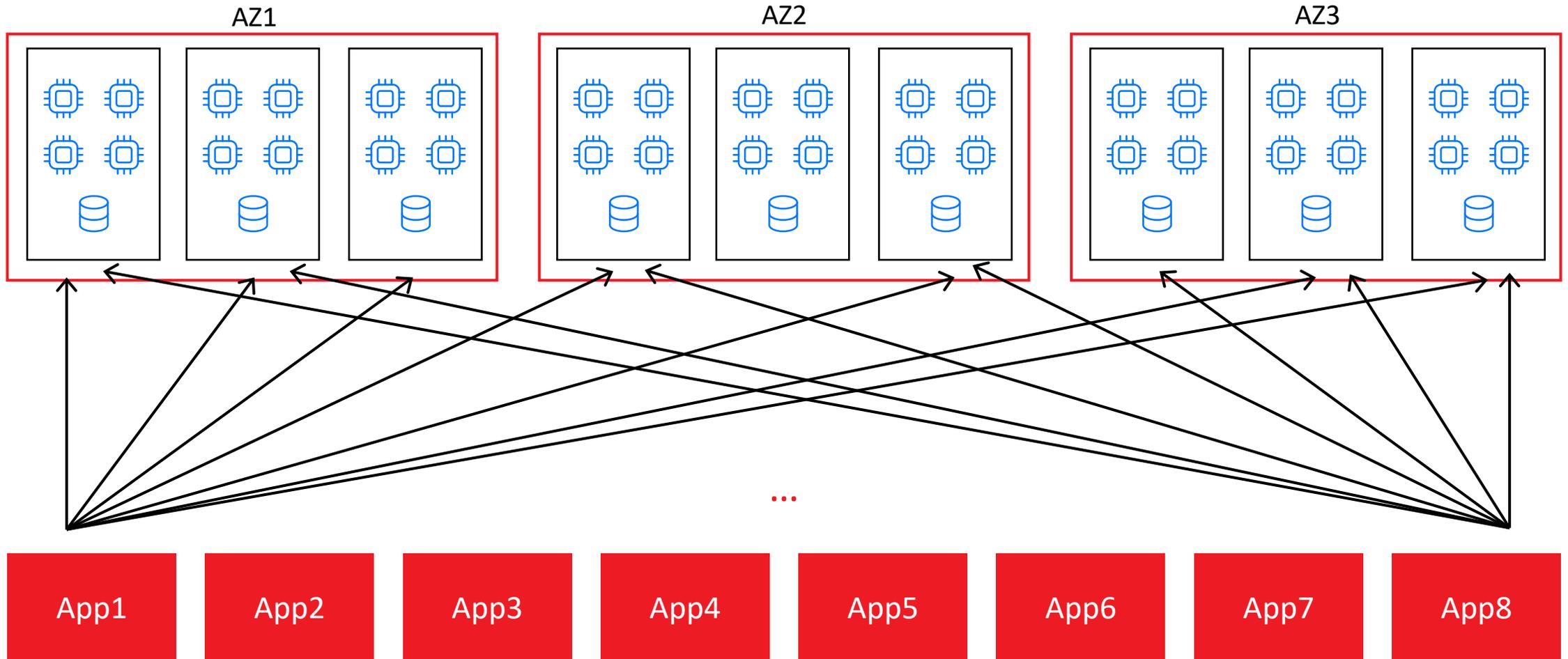
<https://ydb.tech>

<https://github.com/ydb-platform/ydb>

https://t.me/ydb_ru

https://t.me/ydb_en

Кластер YDB с точки зрения клиента



Что такое нода YDB?

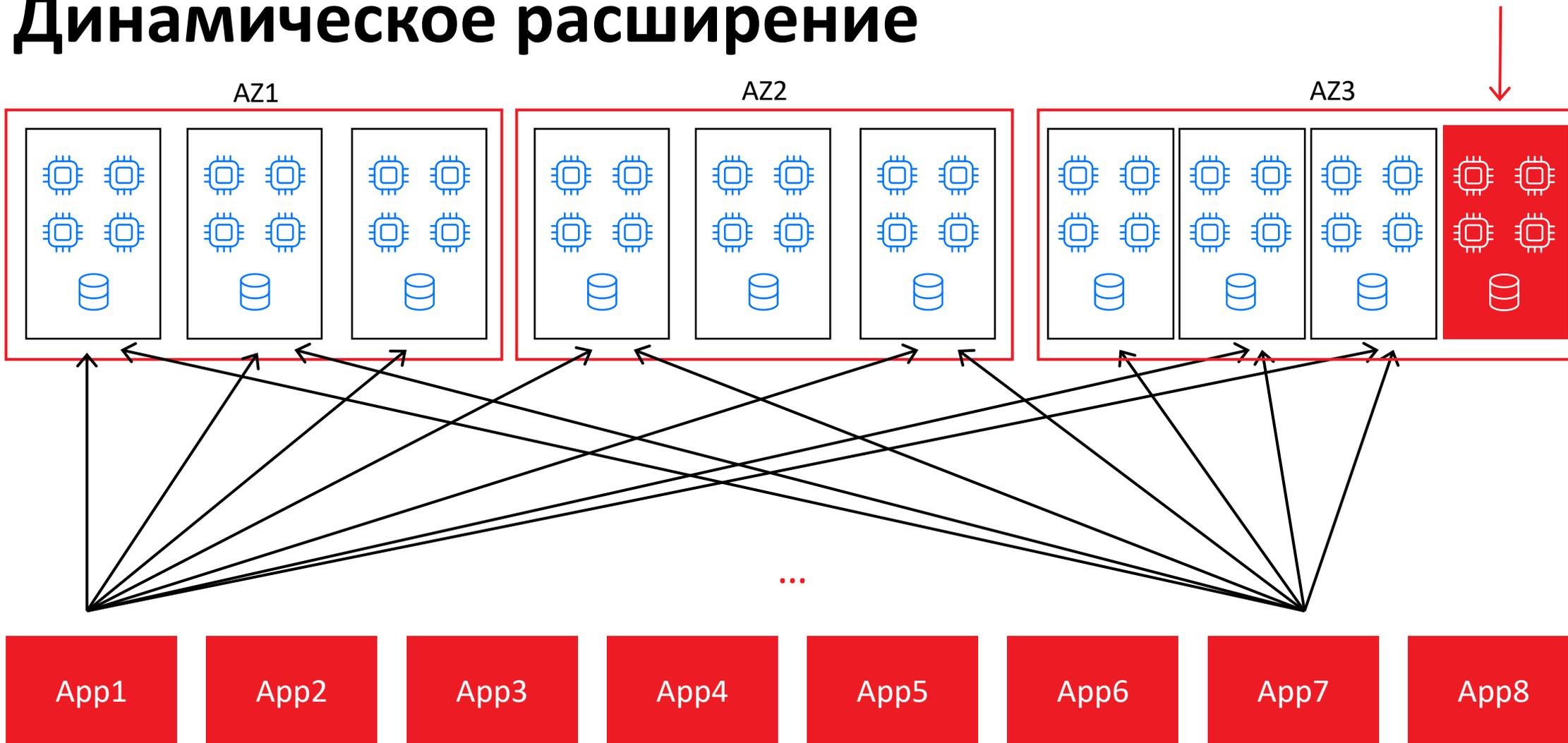
Исполнительный узел в кластере YDB

- + виртуальная машина
- + pod kubernetes
- + отдельный самостоятельный процесс

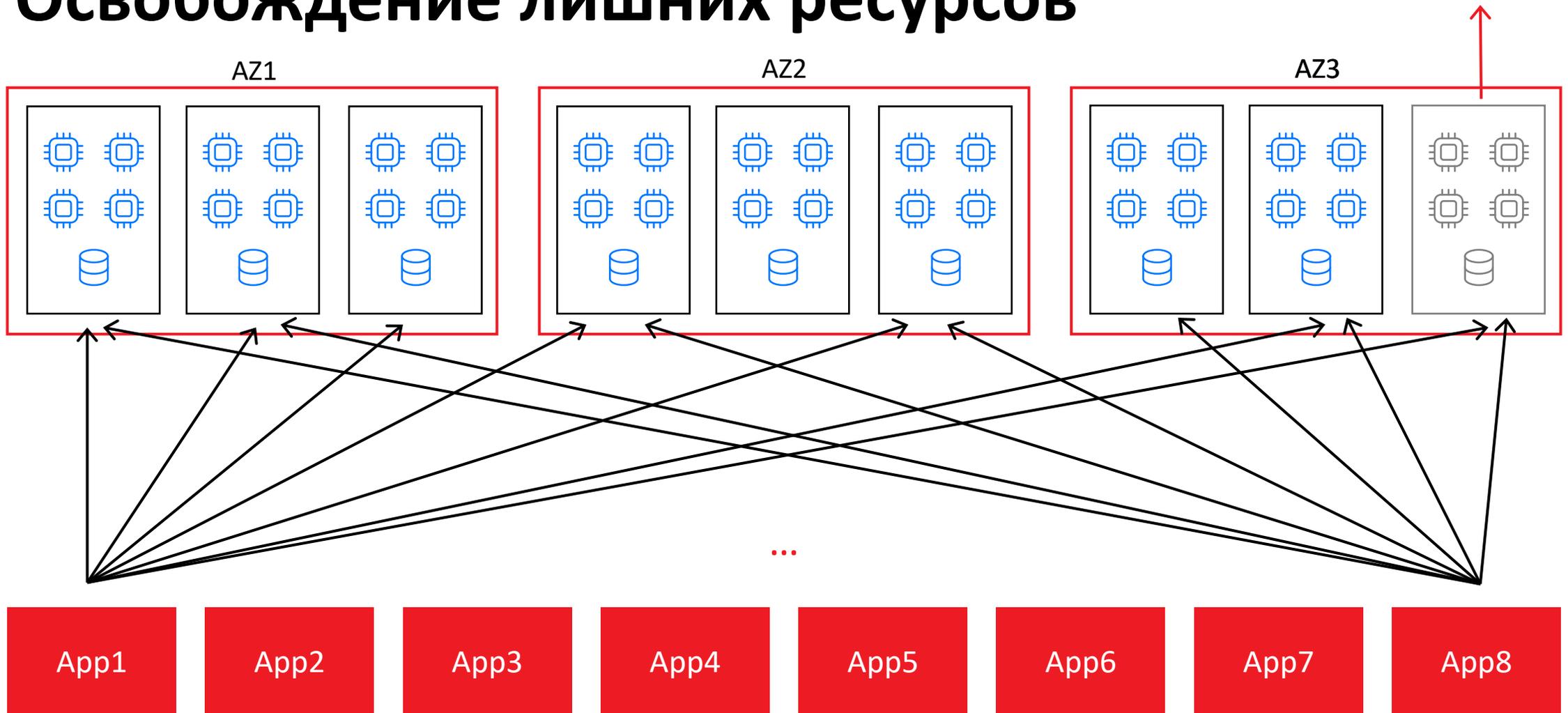
Особенности:

1. однородность
2. волатильность
3. разные для разных баз данных
4. fqdn + порт
5. grpc API

Динамическое расширение



Освобождение лишних ресурсов



Отказы в большой системе - норма



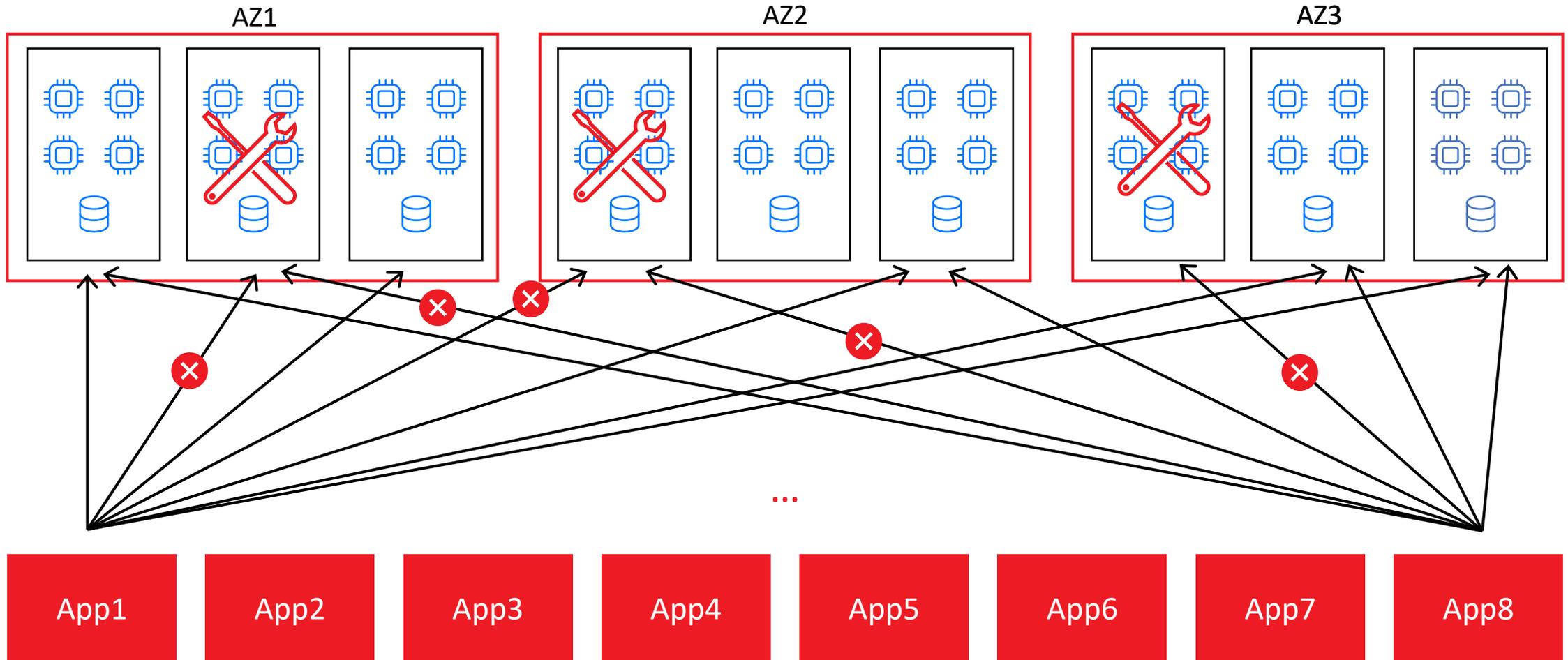
+ *Диски*

- Время наработки на отказ (MTTF) одного диска — 1.4 млн часов или 160 лет
- Время наработки на отказ ≥ 1 диска из ~ 1000 серверов по 4 диска в каждом — 180 часов или 8 дней

+ *Серверы*

- Отказы железа
- Отключение серверов/стоек для обслуживания
- Выход из строя целых стоек (ToR switch, питание)
- Отключение дата-центров (учения, питание, сеть)

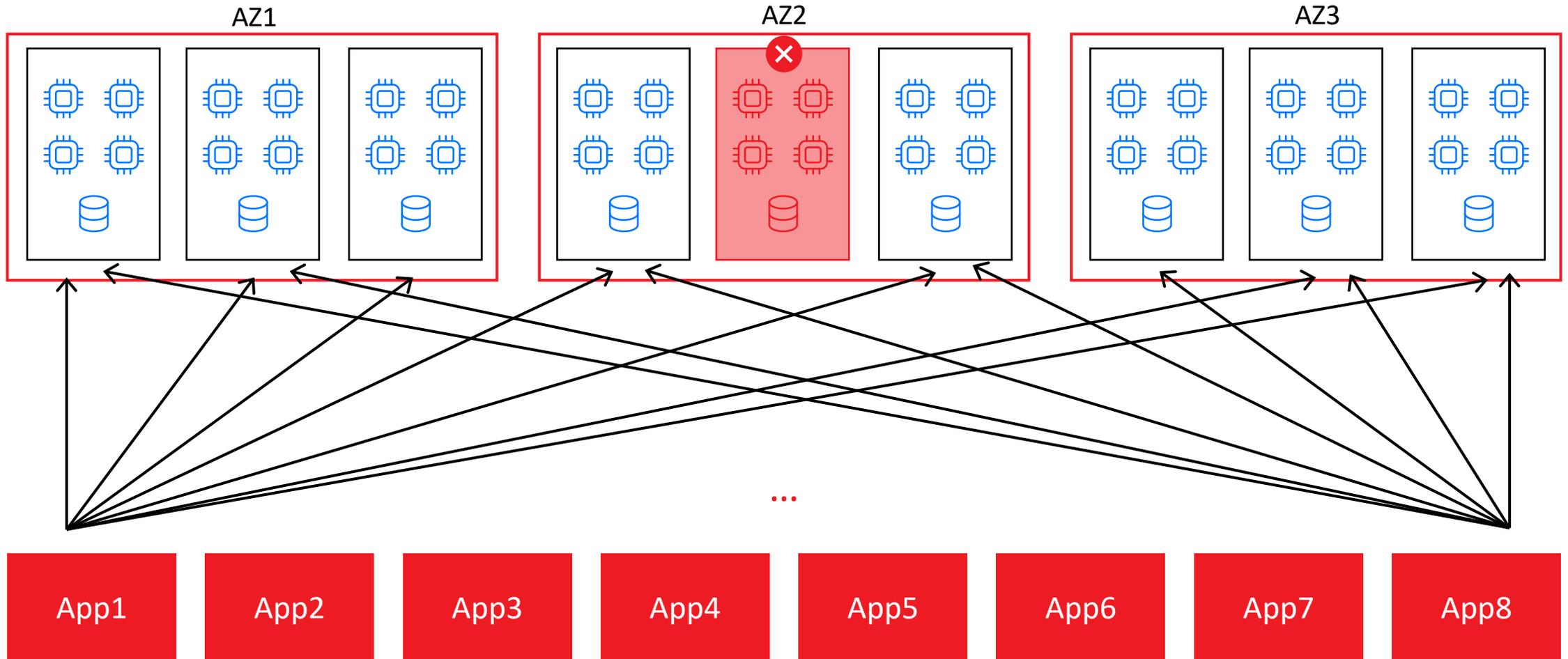
Обслуживание при режиме эксплуатации 24/7



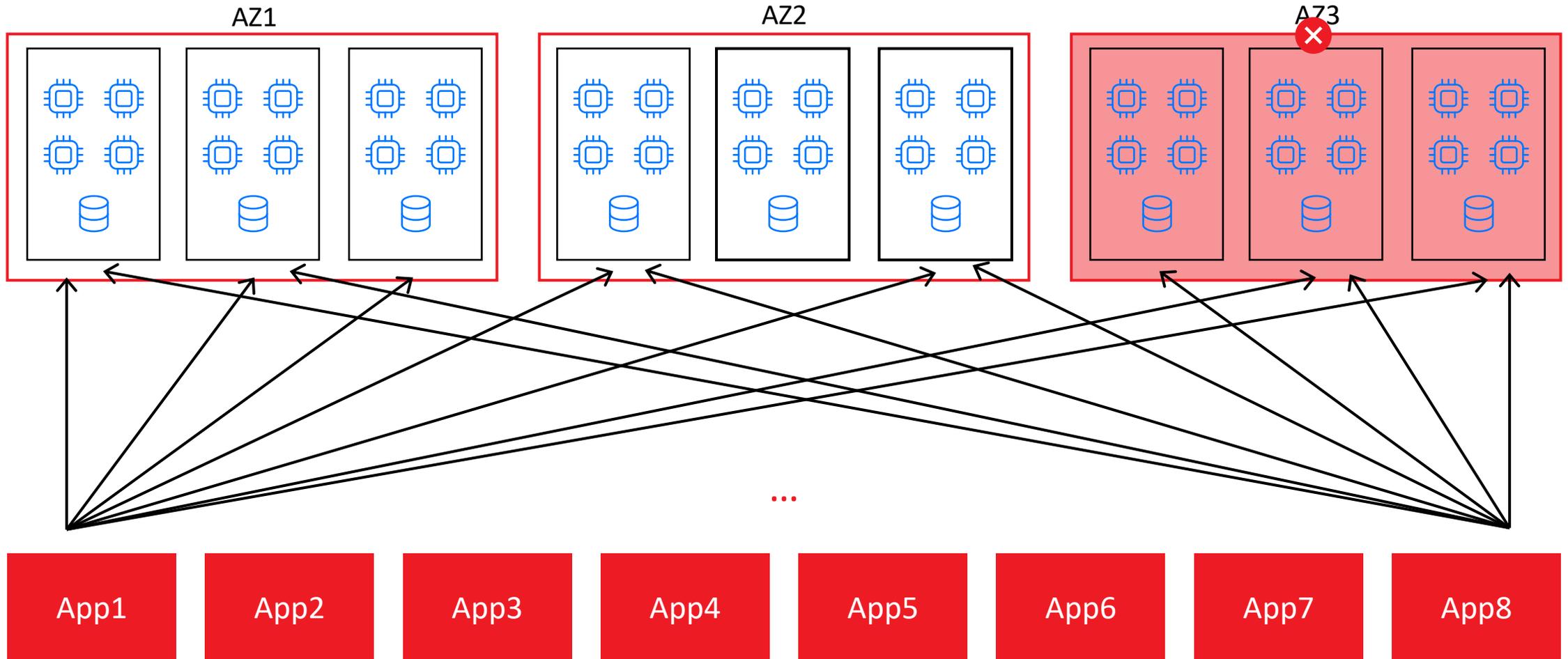
03

Требования к драйверу

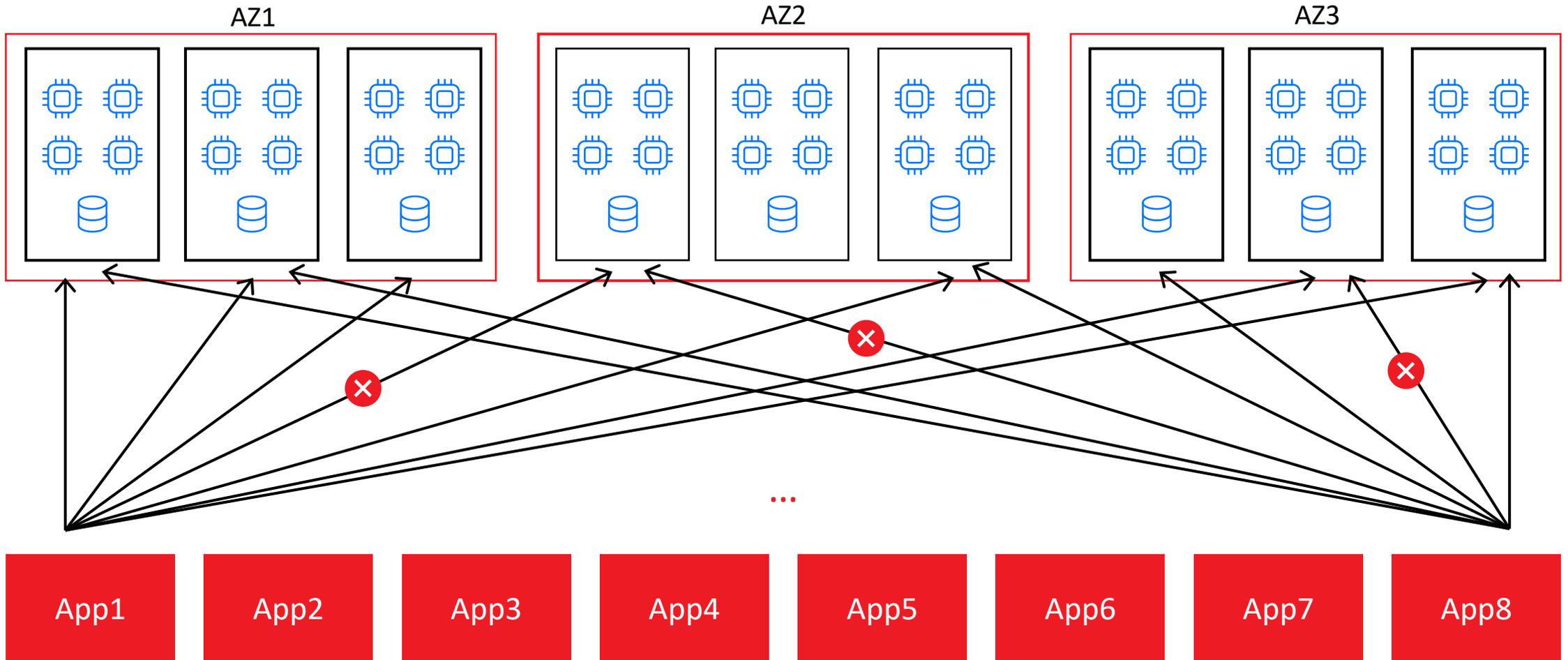
Отказ ноды



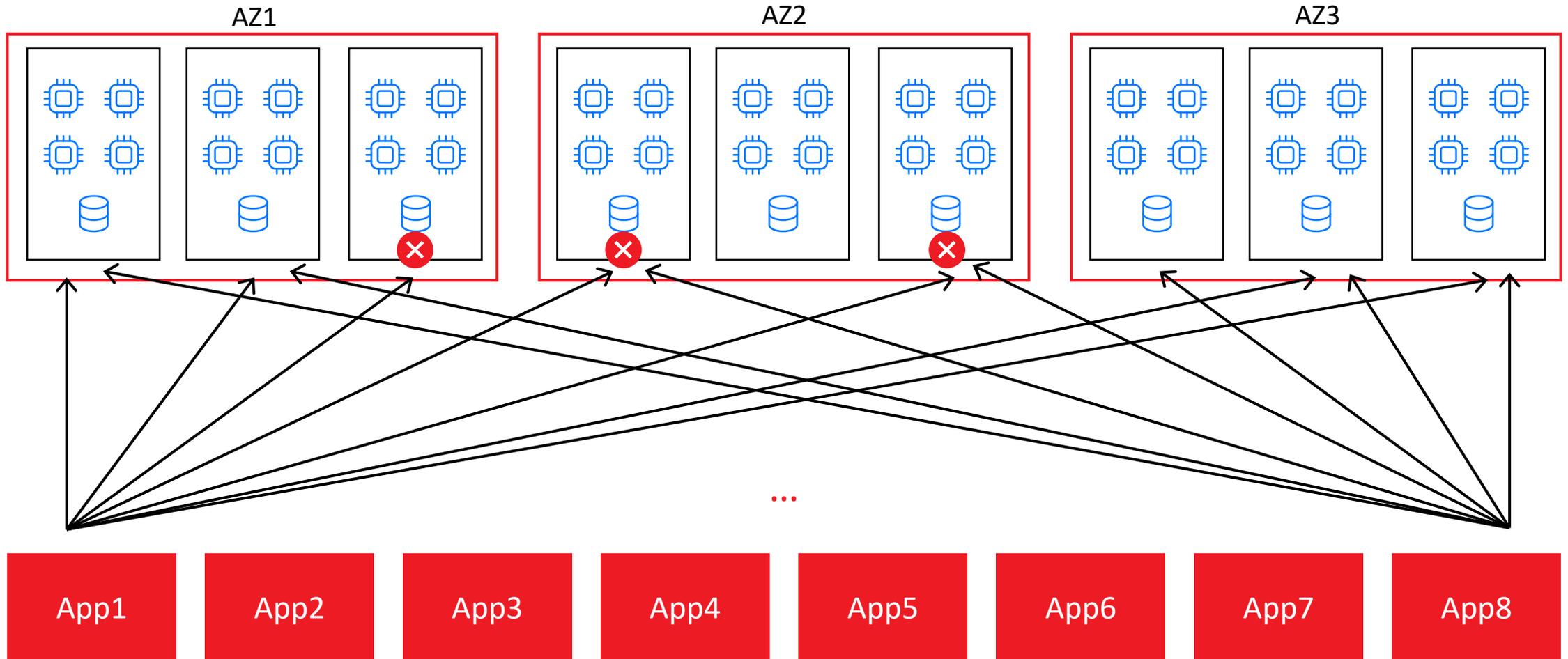
Отказ ДЦ



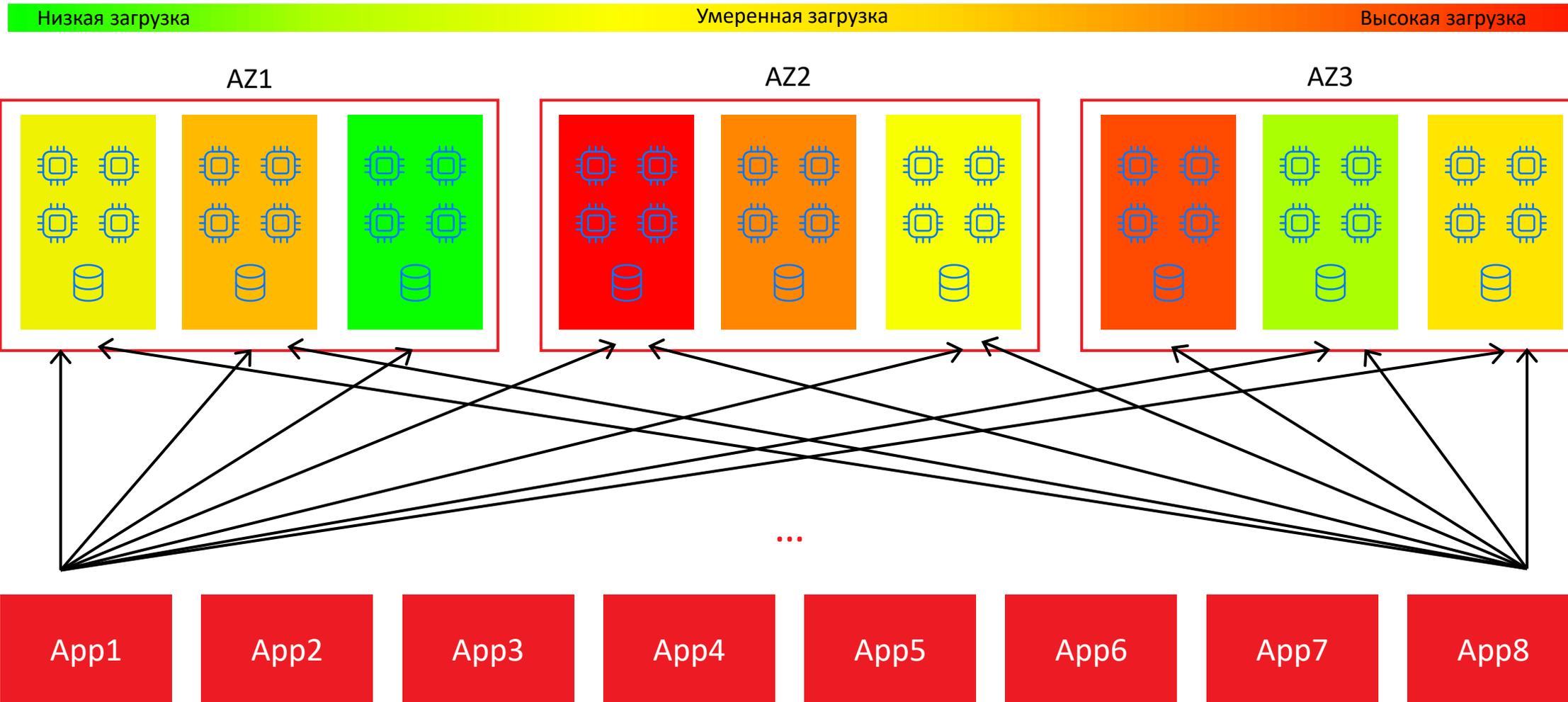
Сетевые ошибки



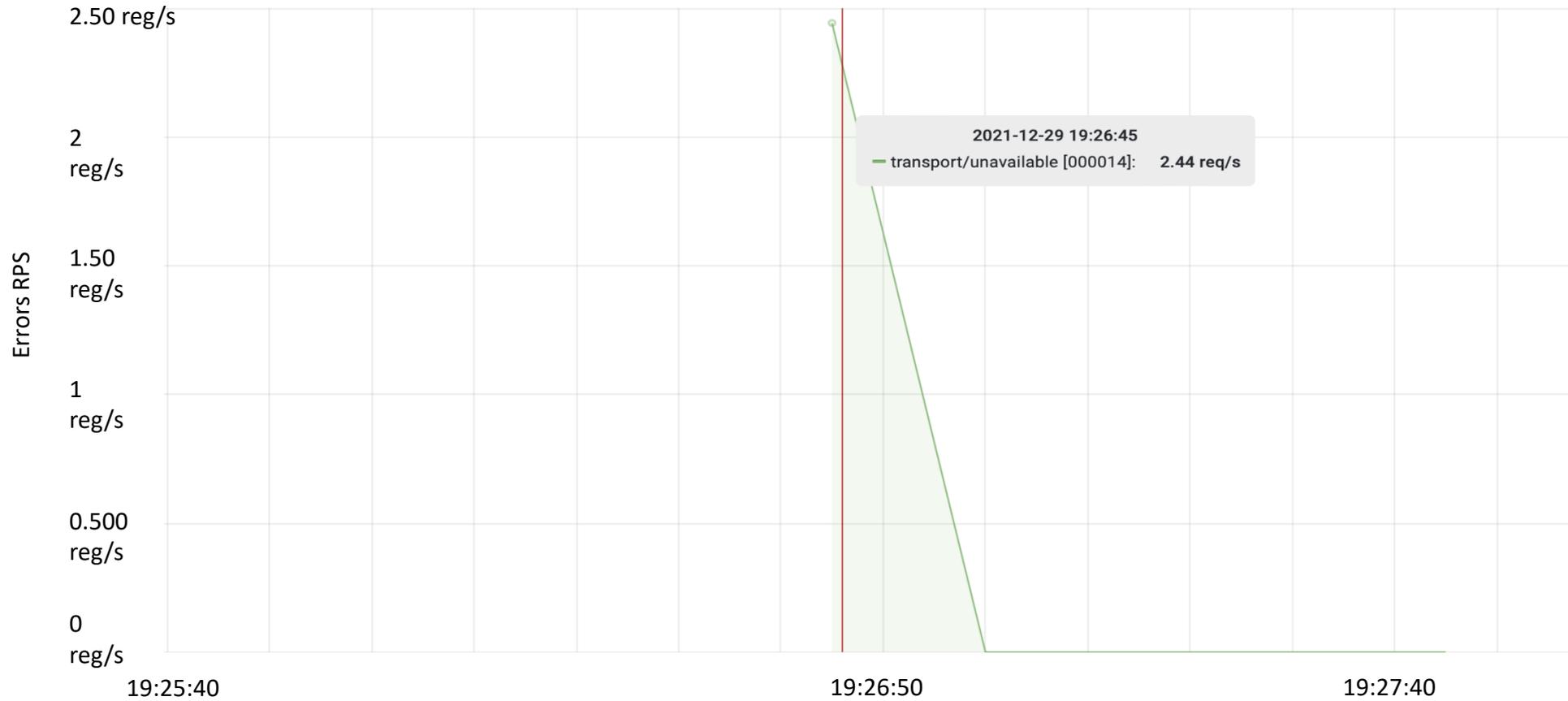
Серверные ошибки



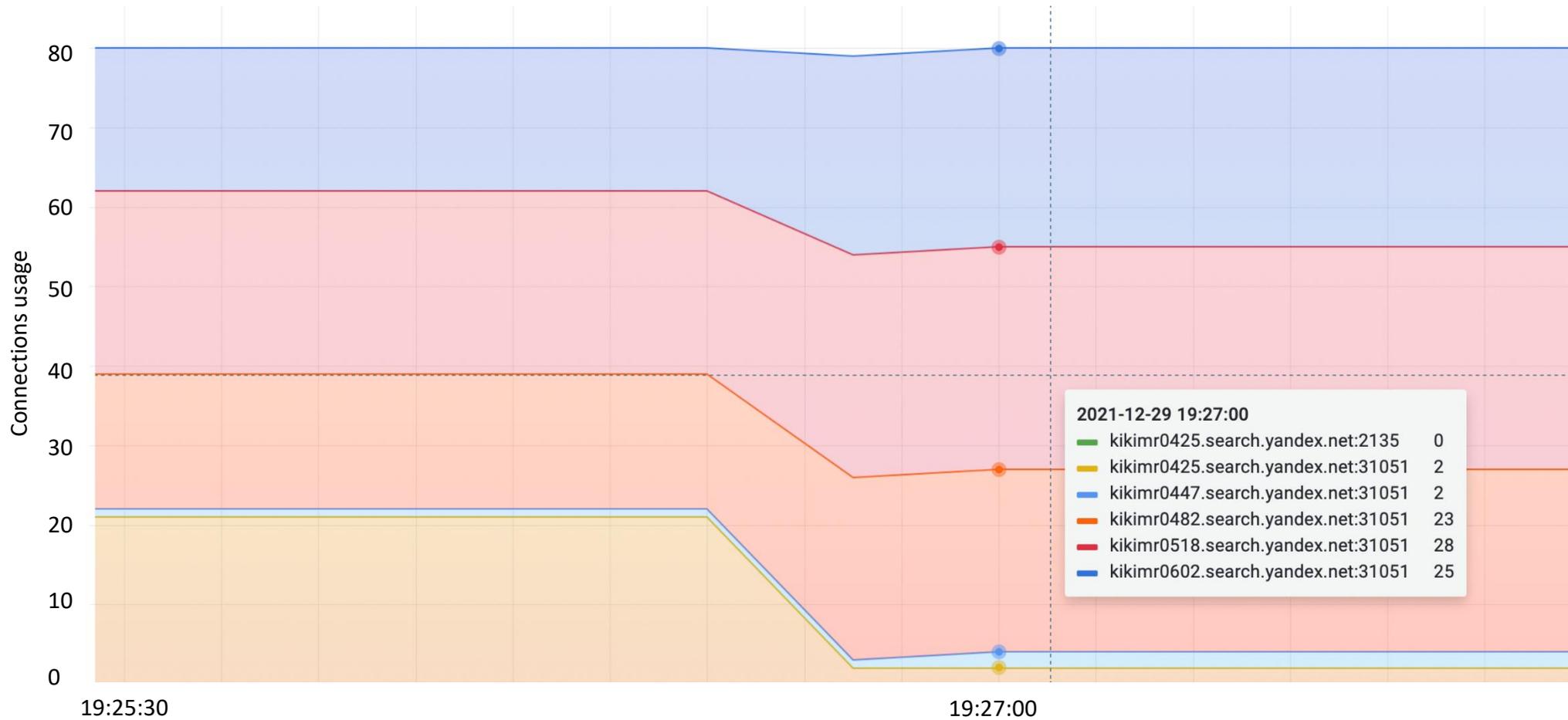
Клиентская и серверная балансировка запросов



Пессимизация соединений



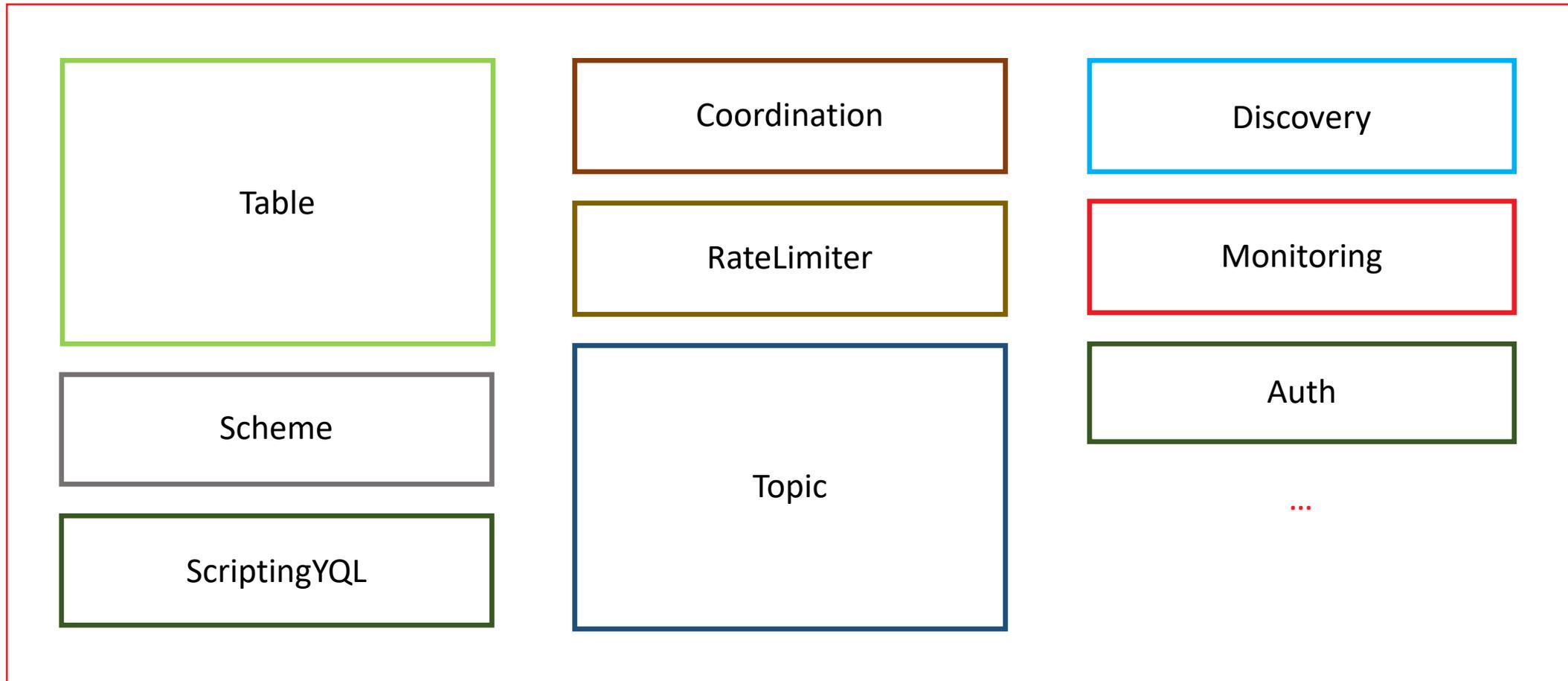
Пессимизация соединений



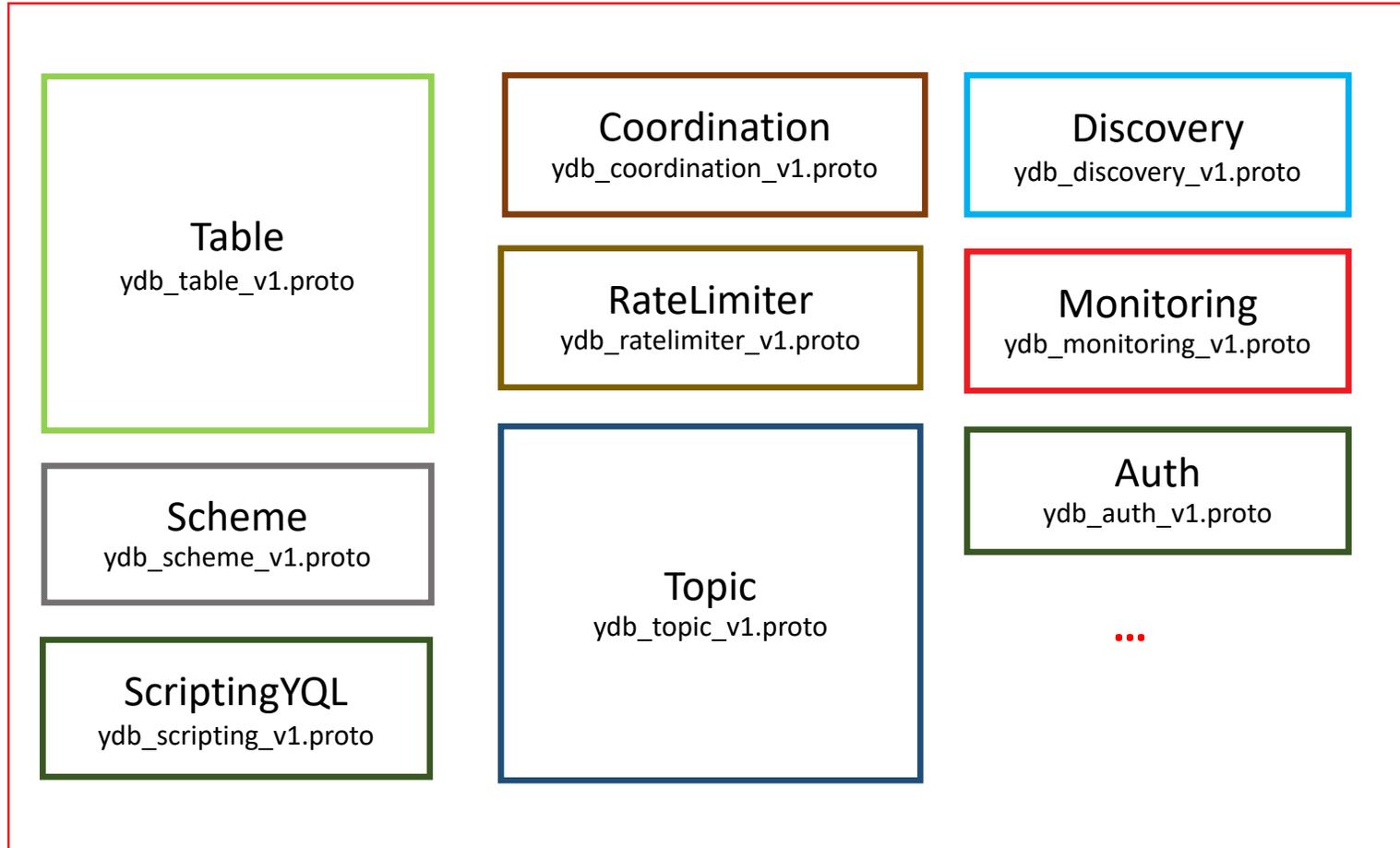
04

YDB API

gRPC-сервисы YDB на ноде YDB

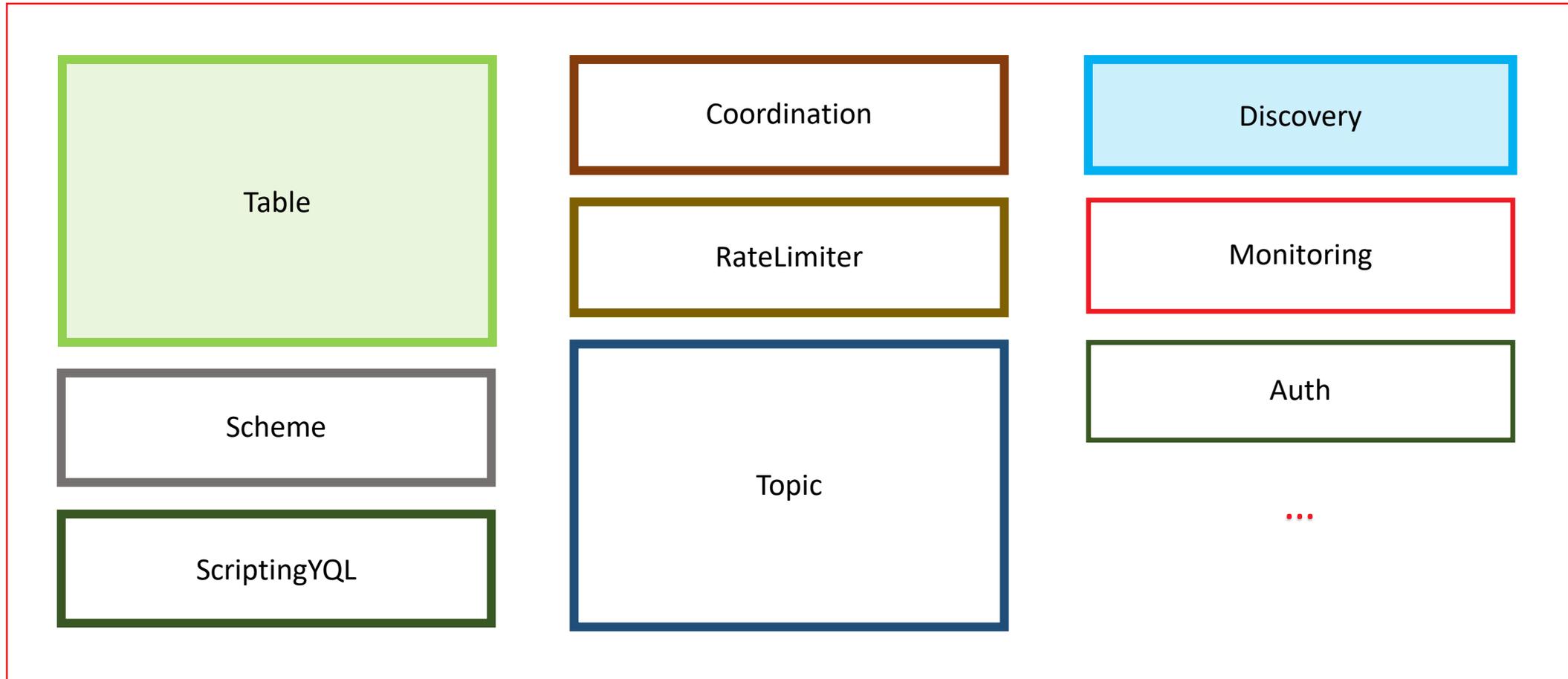


YDB API

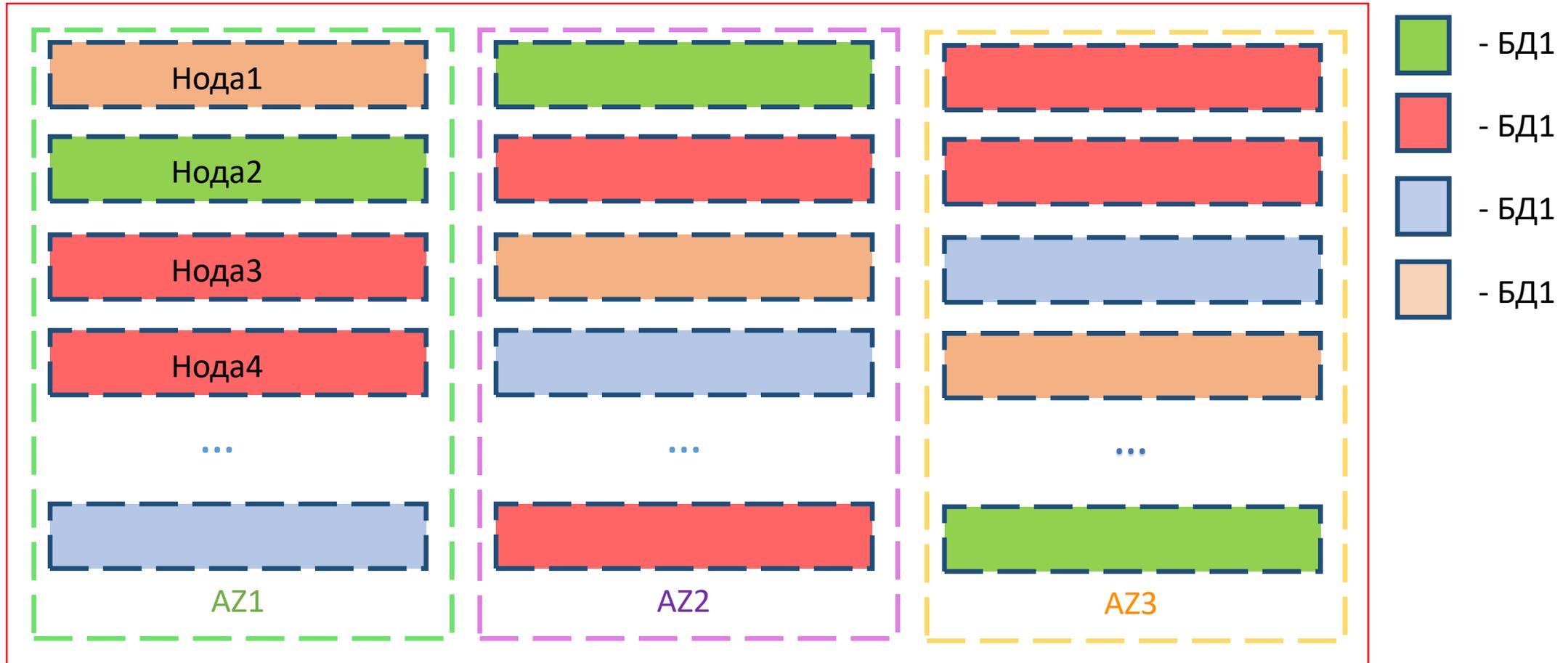


<https://github.com/ydb-platform/ydb-api-protos/>

YDB API



Сервис Discovery



Чтобы начать работать с YDB

1. Выяснить конфигурацию кластера через специальный запрос **Discovery/ListEndpoints**, указав имя базы данных
2. Подключиться напрямую к нодам YDB



Сервис таблиц на ноде YDB

Сессия

- + однопоточная
- + stateful
- + легковесная
- + долгоживущая
- + обеспечивает серверную балансировку
- + имеет TTL

05

Жизненный цикл драйвера YDB

Исходные данные:

1) начальный Endpoint

`grpc://ydb.serverless.yandexcloud.net:2135`

`grpc://lb.etnt8n4t.ydb.mdb.yandexcloud.net:2135`

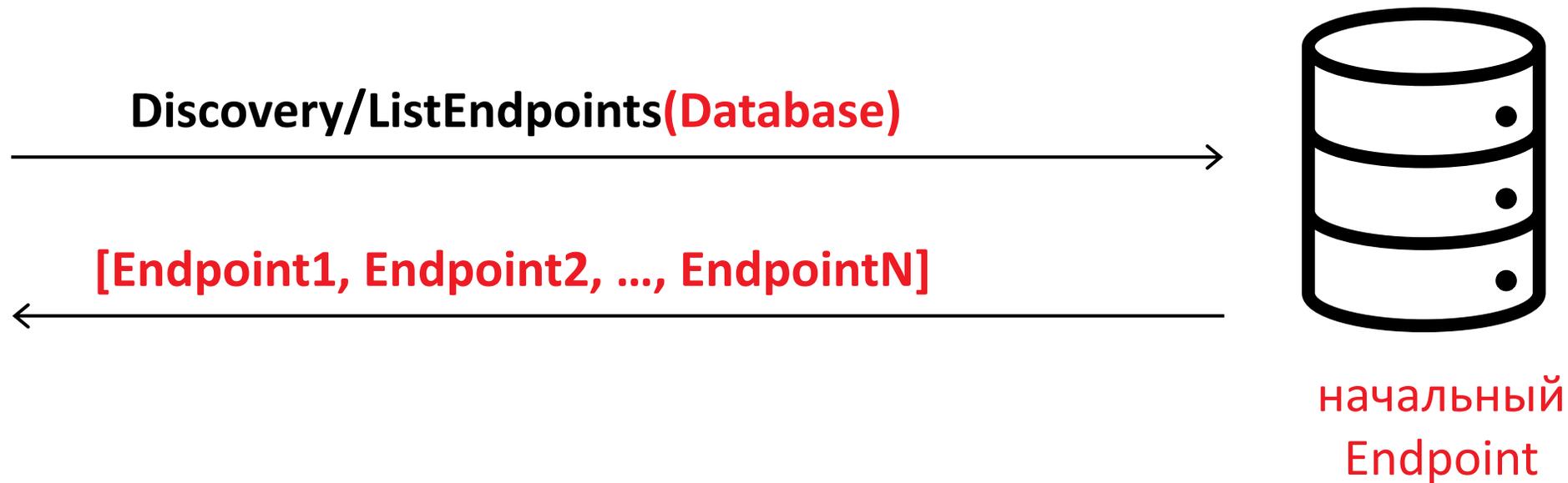
- + Прокси-сервис (балансер)
- + DNS-запись с IP-адресами нод YDB

2) имя базы данных

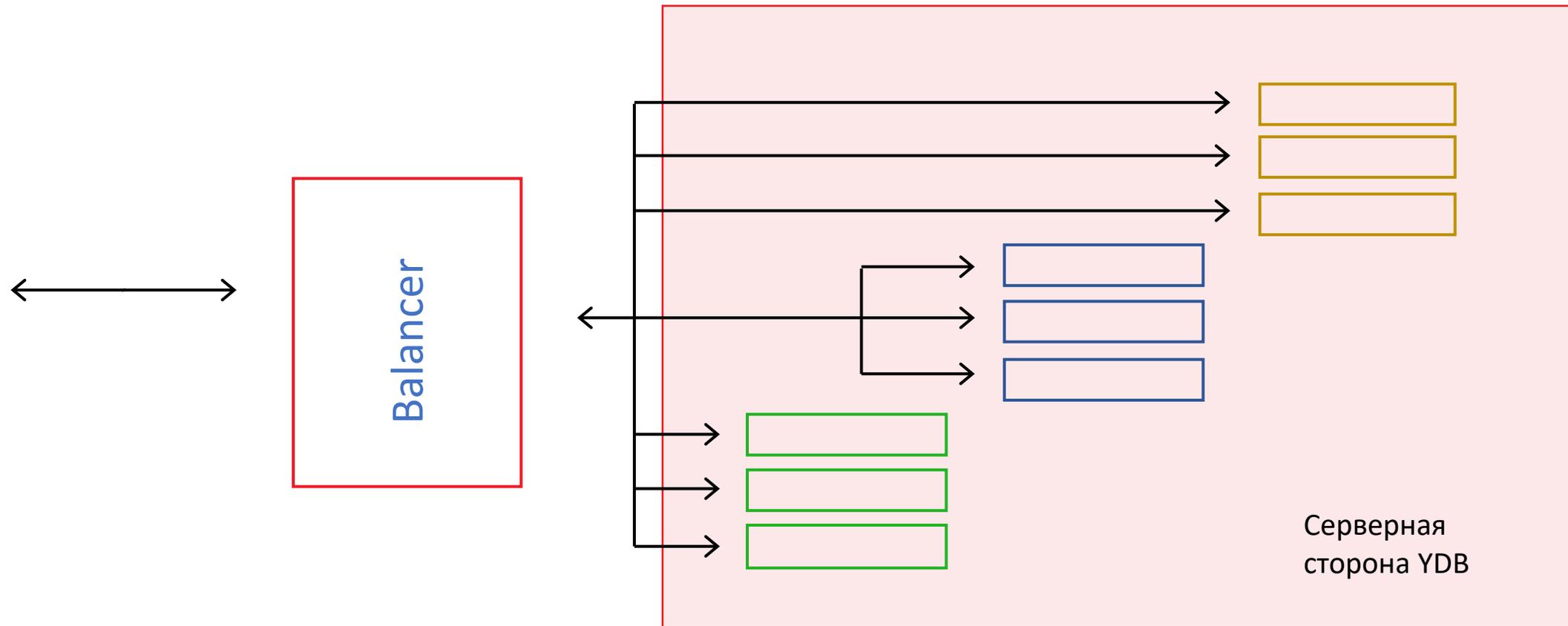
`/ru-central1/b1g8spbal1f3s/etntdt2dn44t`

1. Инициализация драйвера

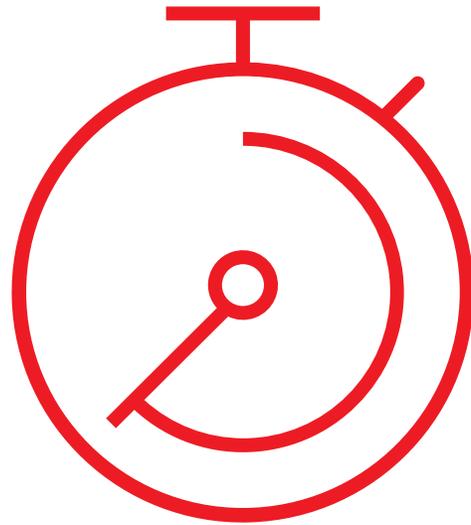
1.1 Выясняем конфигурацию кластера



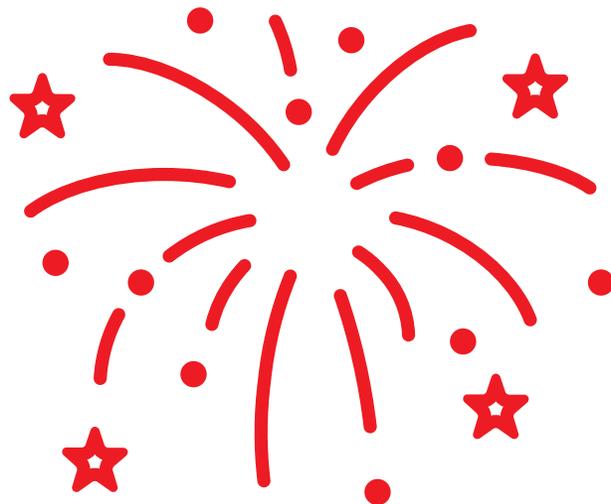
1.2 Инициализируем клиентский балансировщик



1.3 Запускаем фоновый *Discovery* и применение результатов в балансировщике



1 минута



Драйвер инициализирован

2. Выполнение табличных запросов

```
SELECT `title`  
FROM `/highload/2022/november/reports`  
WHERE speaker=$speakerName;
```

Простейшая реализация:

1. Создать сессию



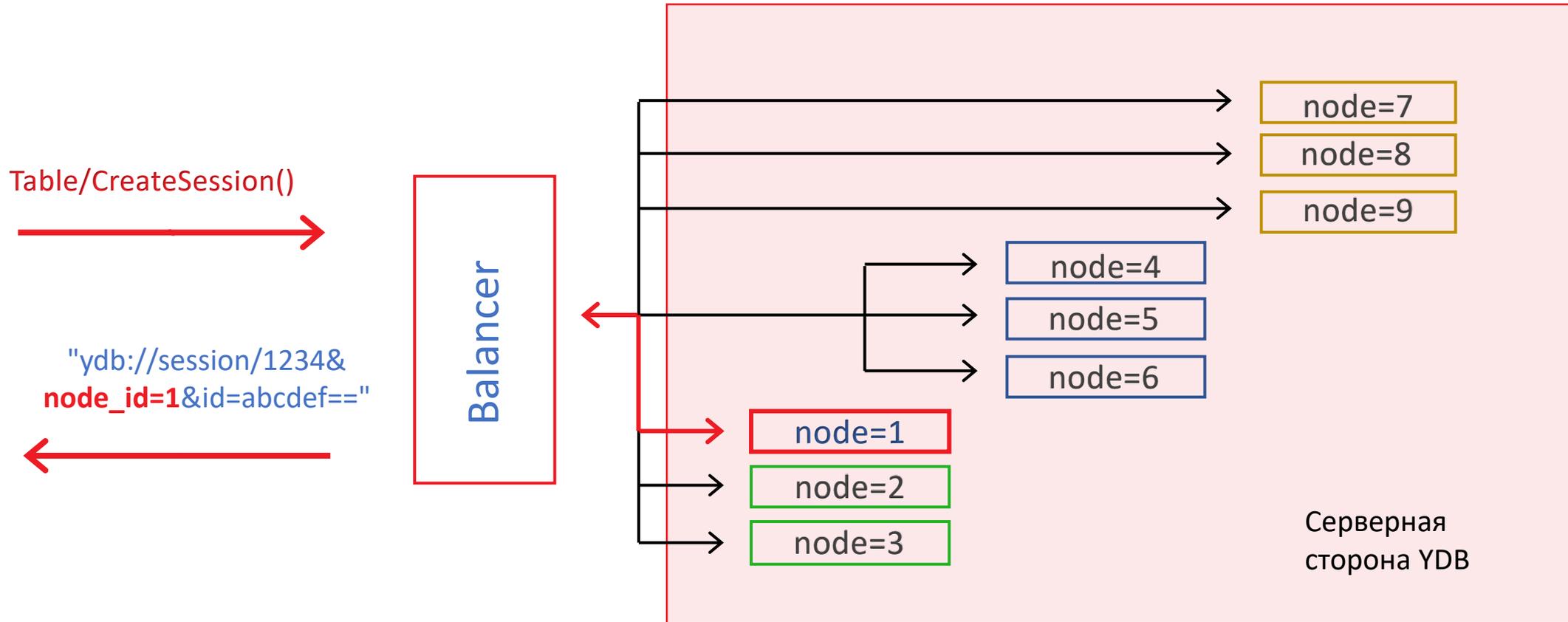
2. Выполнить запрос на сессии



3. Закрыть сессию



2.1 Создание сессии

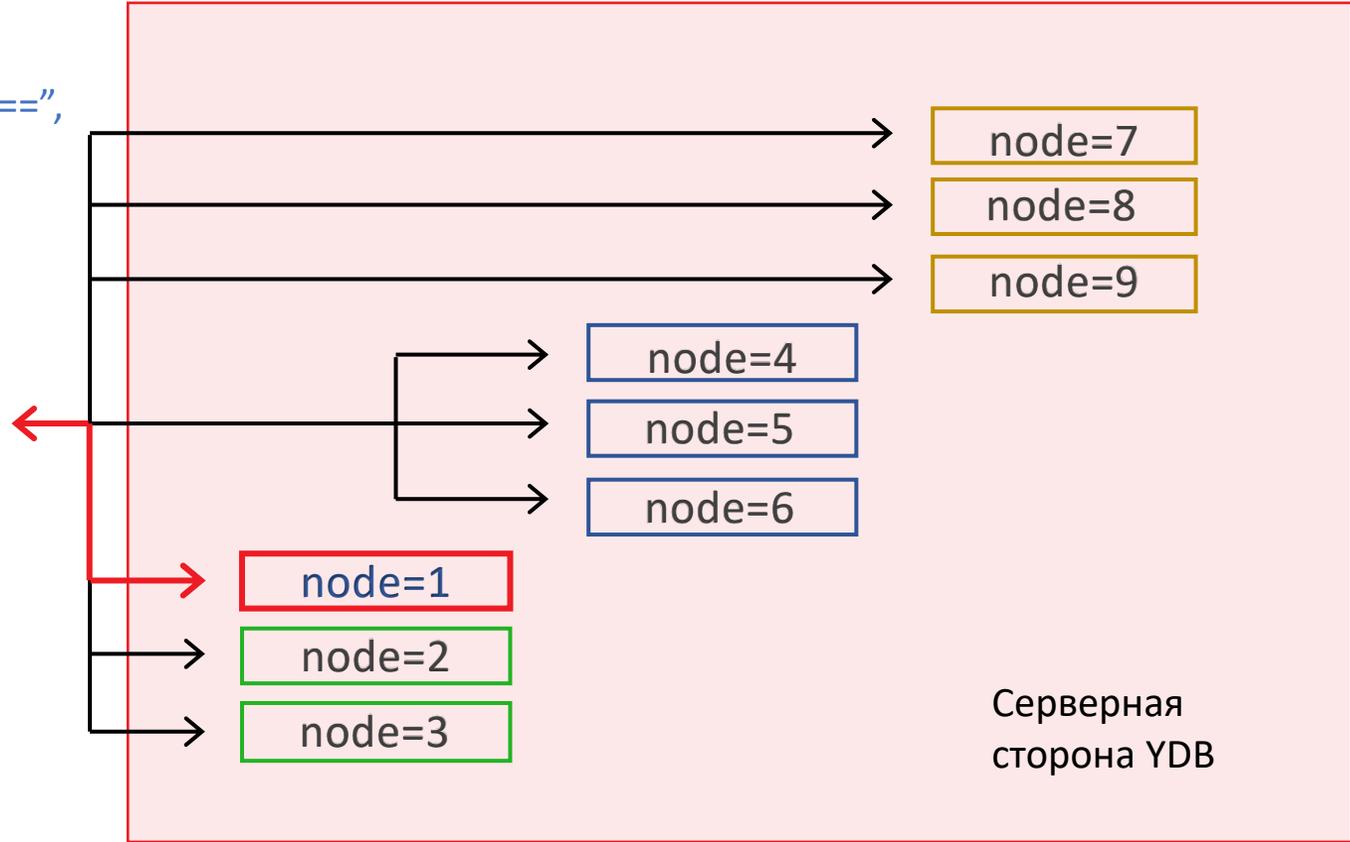
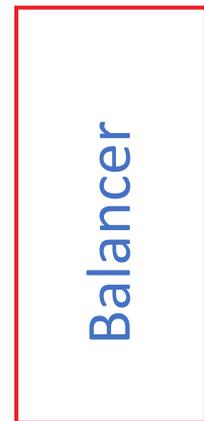


2.2 Выполнение запроса

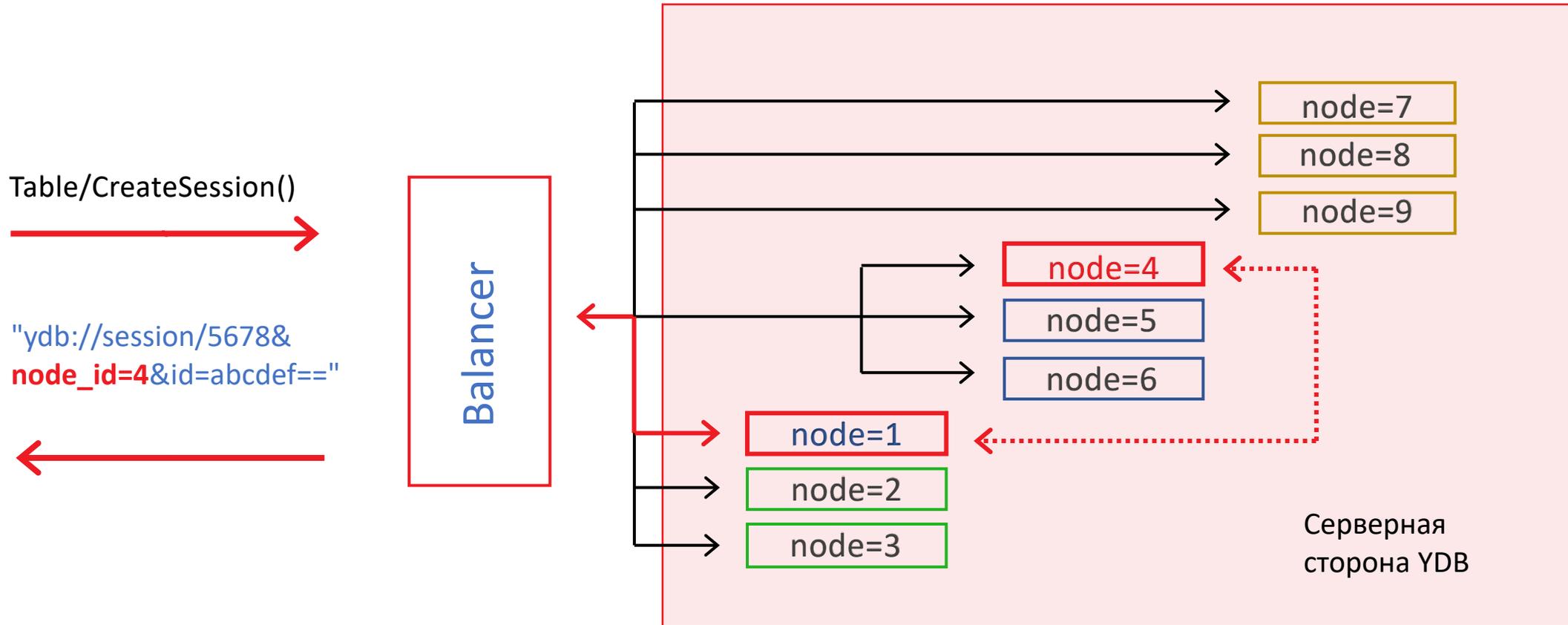
```
Table/ExecuteDataQuery(  
  "ydb://session/1234&node_id=1&id=abcdef==",  
  query,  
  $speakerName="Мясников Алексей",  
)
```



title="Просто о сложном:
как работает драйвер
распределенной
базы данных YDB"



2.1 Сессия может быть создана на другой ноде

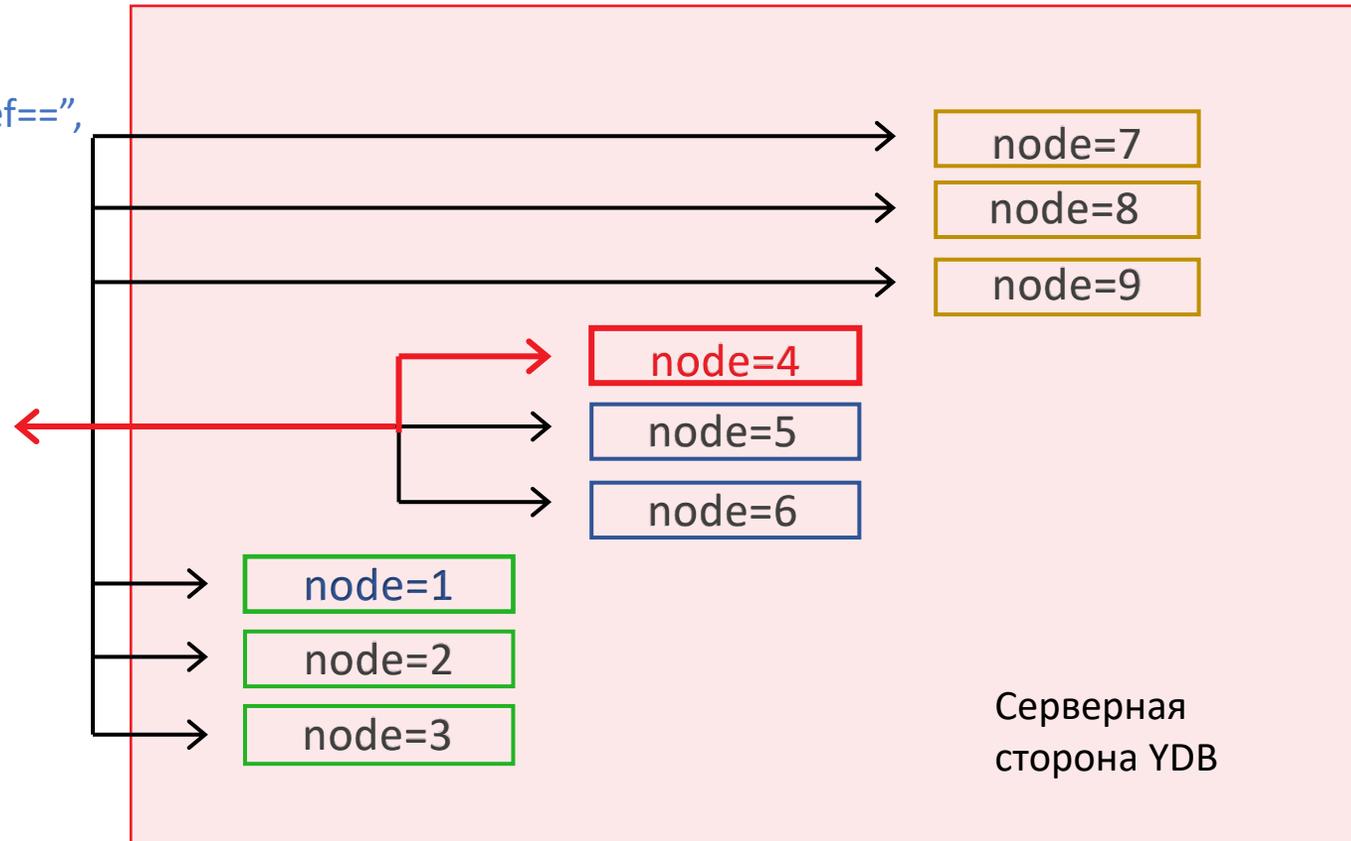
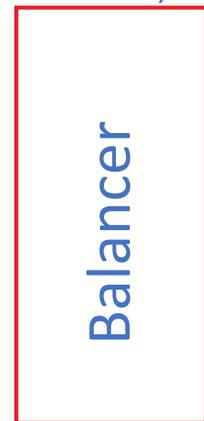


2.2 Запросы на сессии следует отправлять на «правильную» ноду

```
Table/ExecuteDataQuery(  
  "ydb://session/1234&node_id=4&id=abcdef==",  
  query,  
  $speakerName="Мясников Алексей",  
)
```



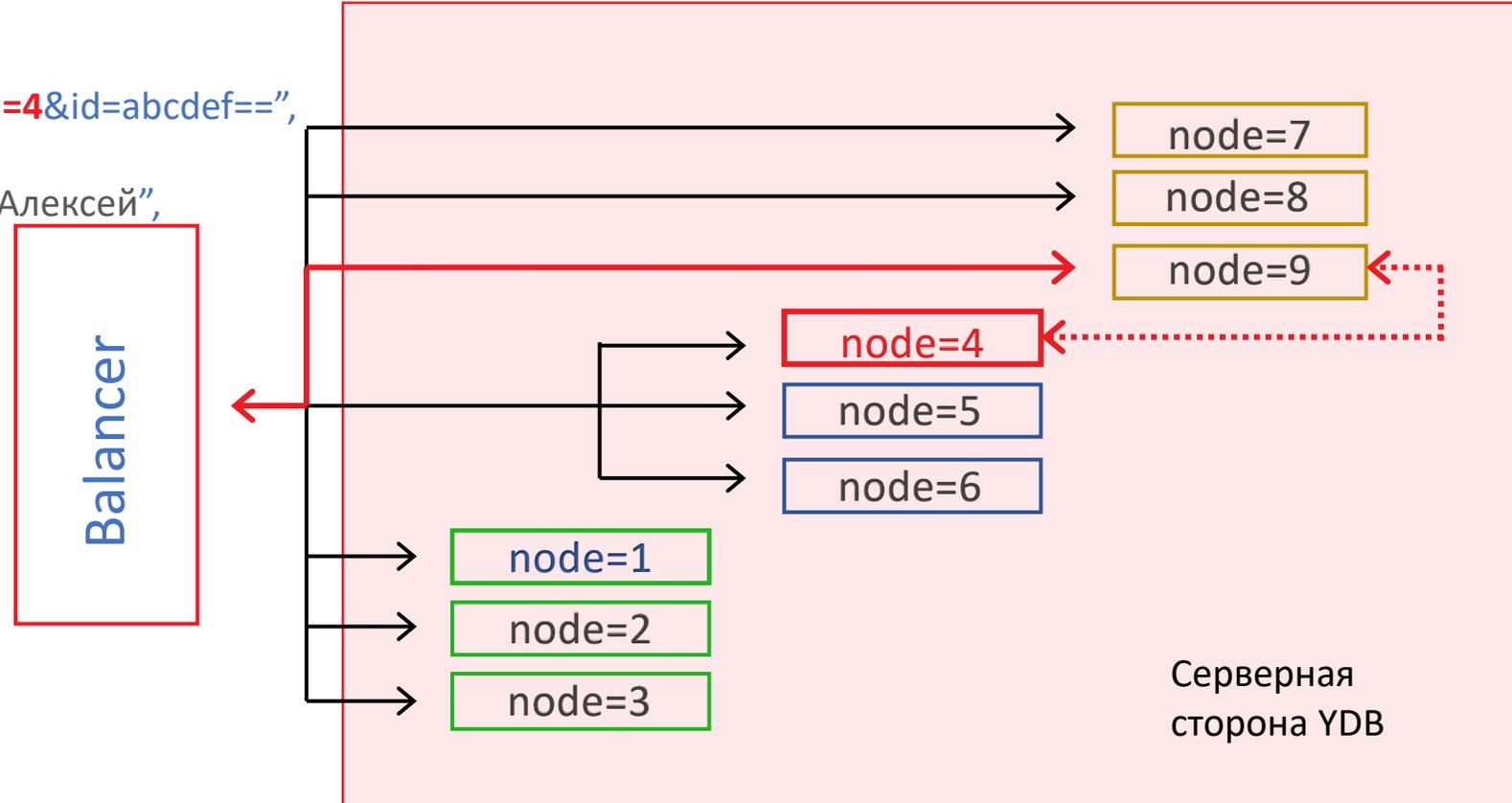
title="Просто о сложном:
как работает драйвер
распределенной
базы данных YDB"



2.2 Если запрос ушел на «неправильную» ноду

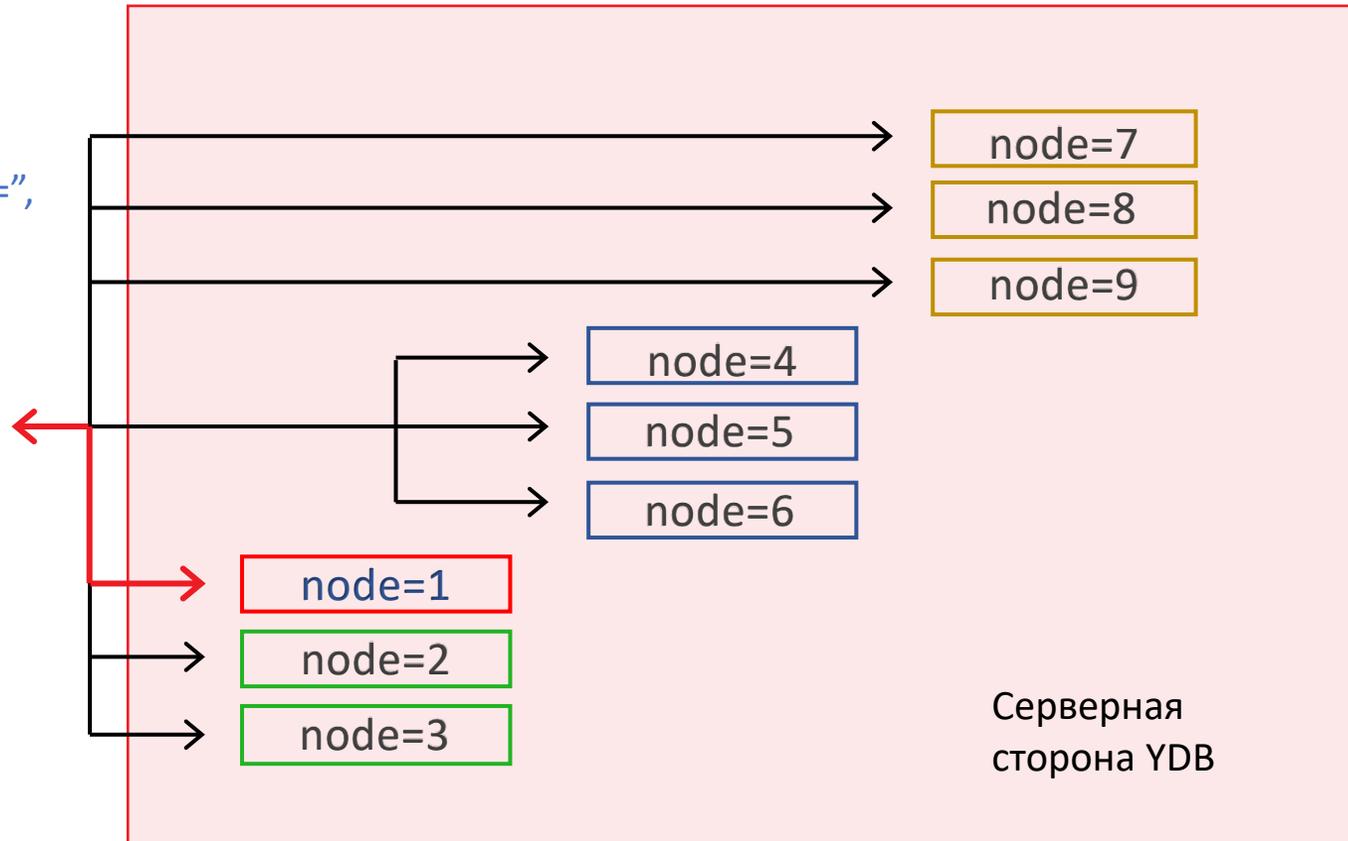
```
Table/ExecuteDataQuery(  
  "ydb://session/1234&node_id=4&id=abcdef==",  
  query,  
  $speakerName="Мясников Алексей",  
)
```

title="Просто о сложном:
как работает драйвер
распределенной
базы данных YDB"

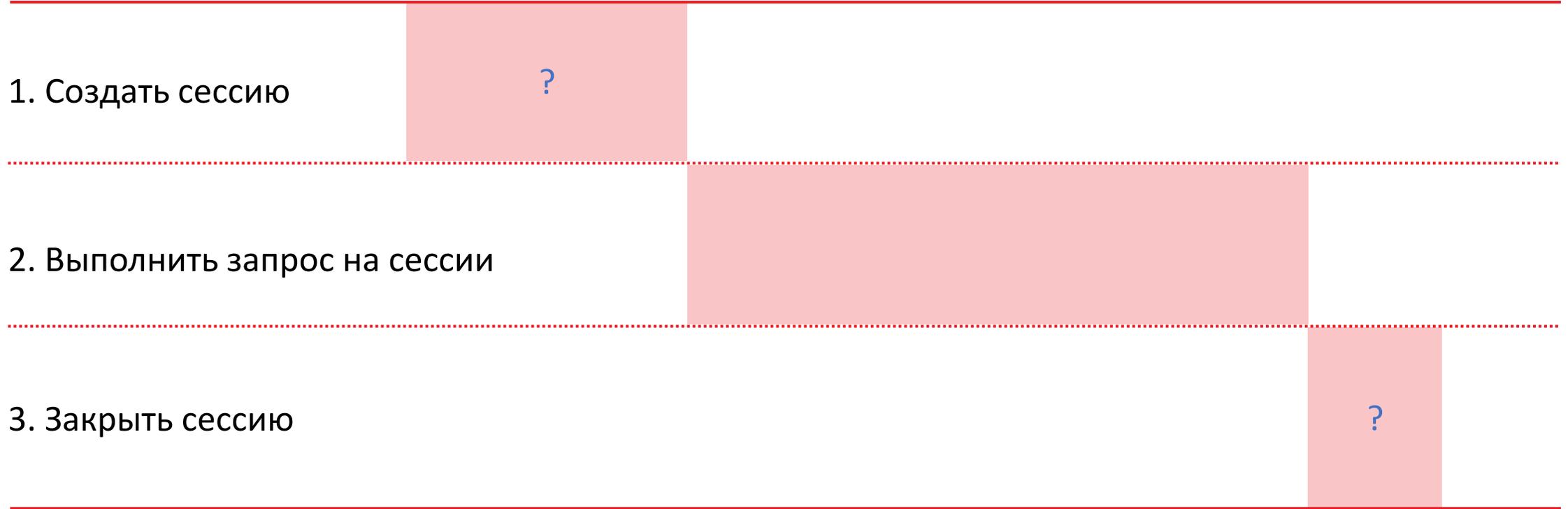


2.3 Закрывать сессию

```
Table/CloseSession(  
"ydb://session/1234&node_id=1&id=abcdef==",  
)
```



Что можно оптимизировать?



3. Пул сессий

3.1 Работа с пулом:

1. Создать сессию



2. Выполнить запрос на сессии



3. ~~Закреть сессию~~ Положить в пул



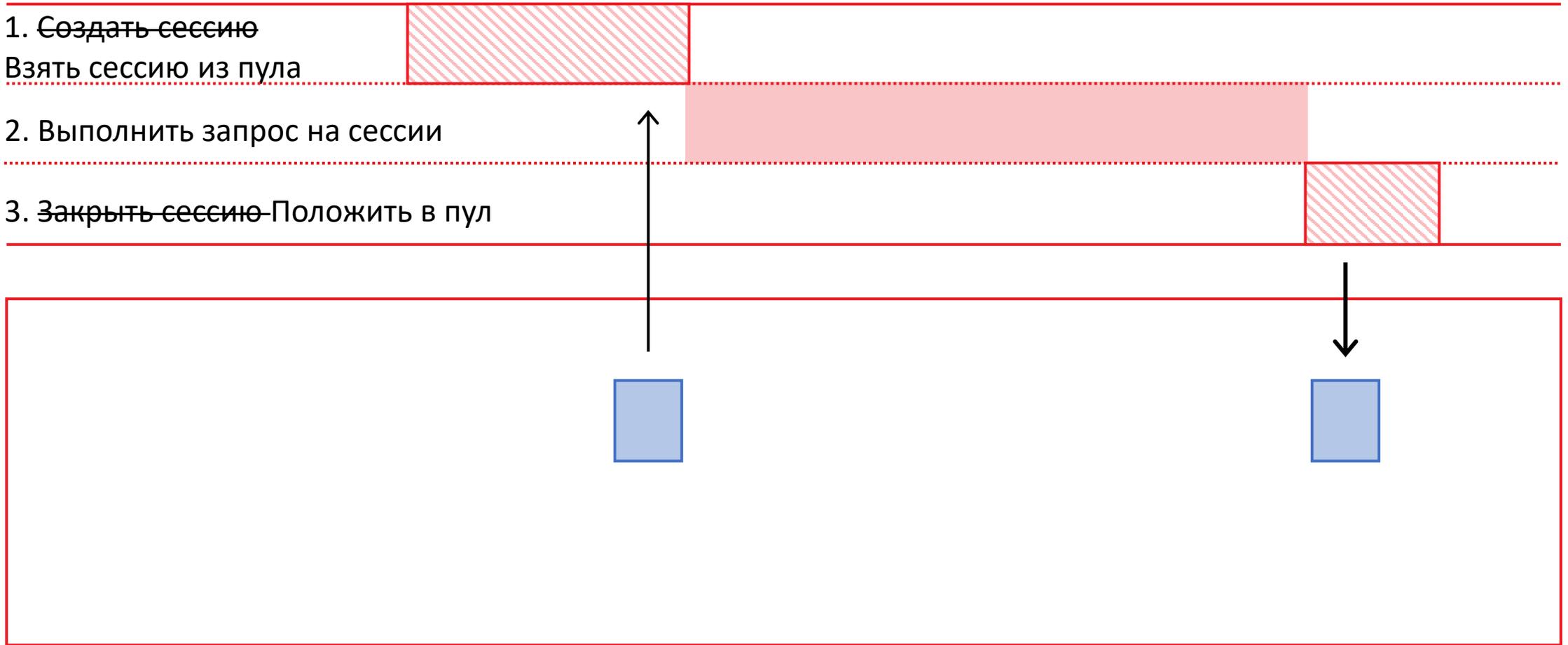
3.1 Работа с пулом:

1. Создать сессию

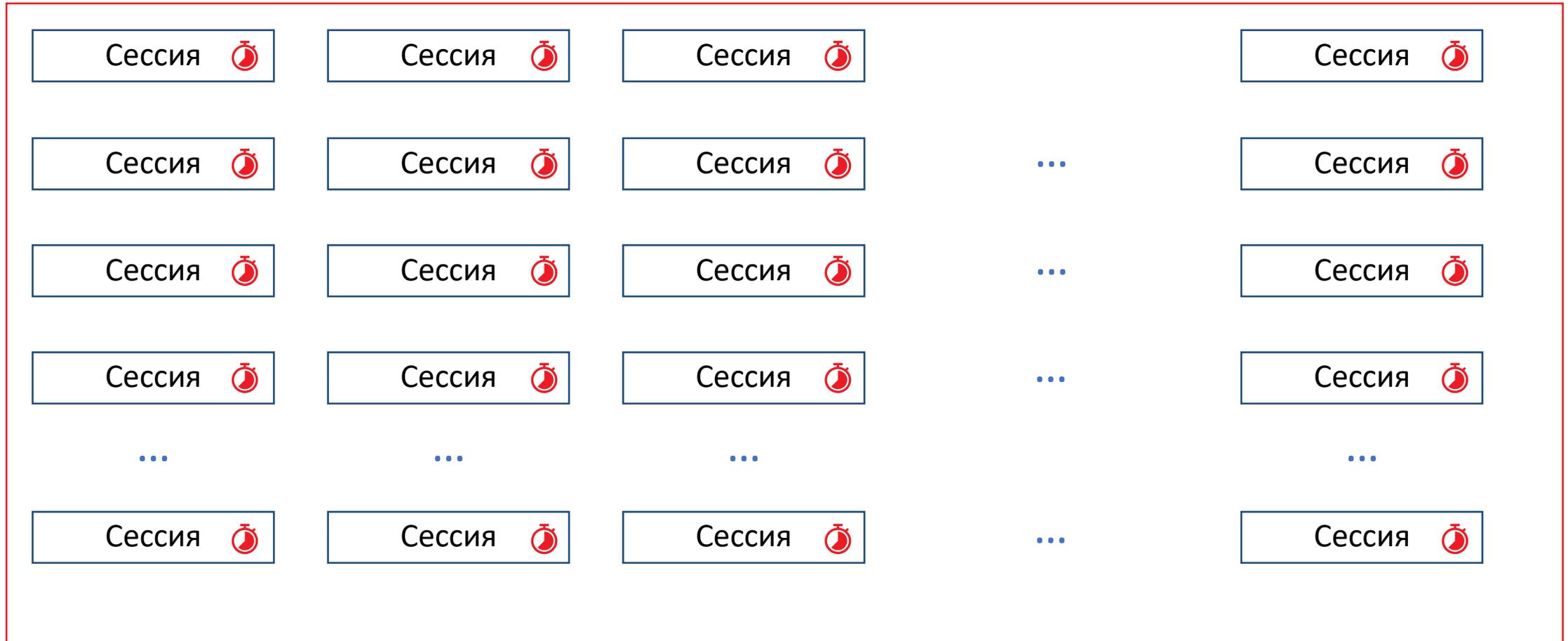
Взять сессию из пула

2. Выполнить запрос на сессии

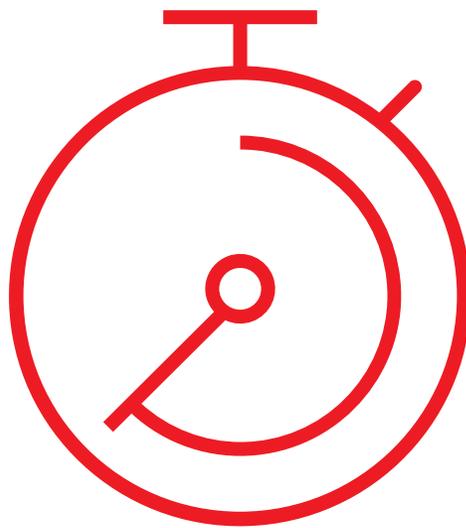
3. Закрыть сессию Положить в пул



Время жизни сессии



3.2 Фоновый *KeepAlive* для простаивающих сессий



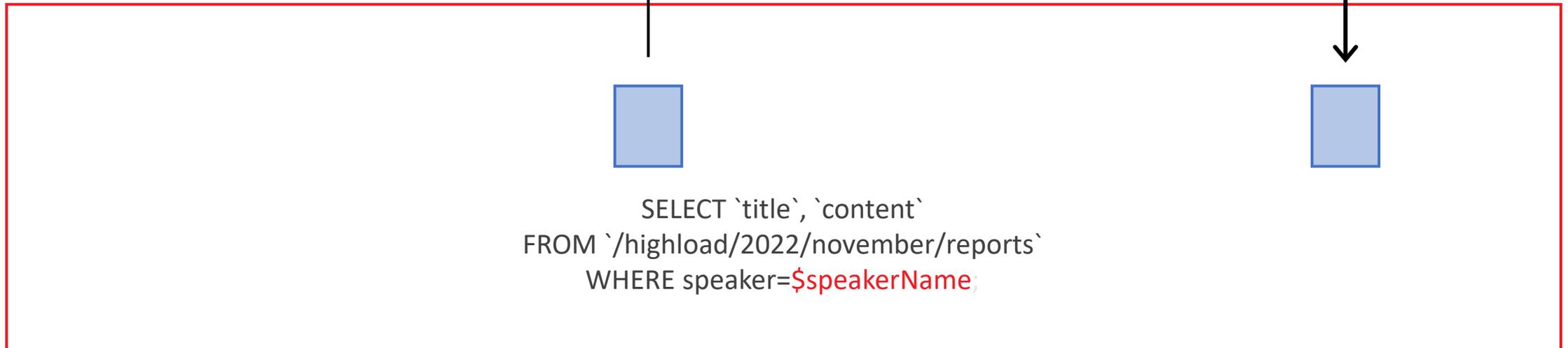
Что еще оптимизировать?

1. Создать сессию

Взять сессию из пула

2. Выполнить запрос на сессии

3. Закрыть сессию Положить в пул



Что еще оптимизировать?

1. Создать сессию

Взять сессию из пула

Компиляция запроса

2. Выполнить запрос на сессии

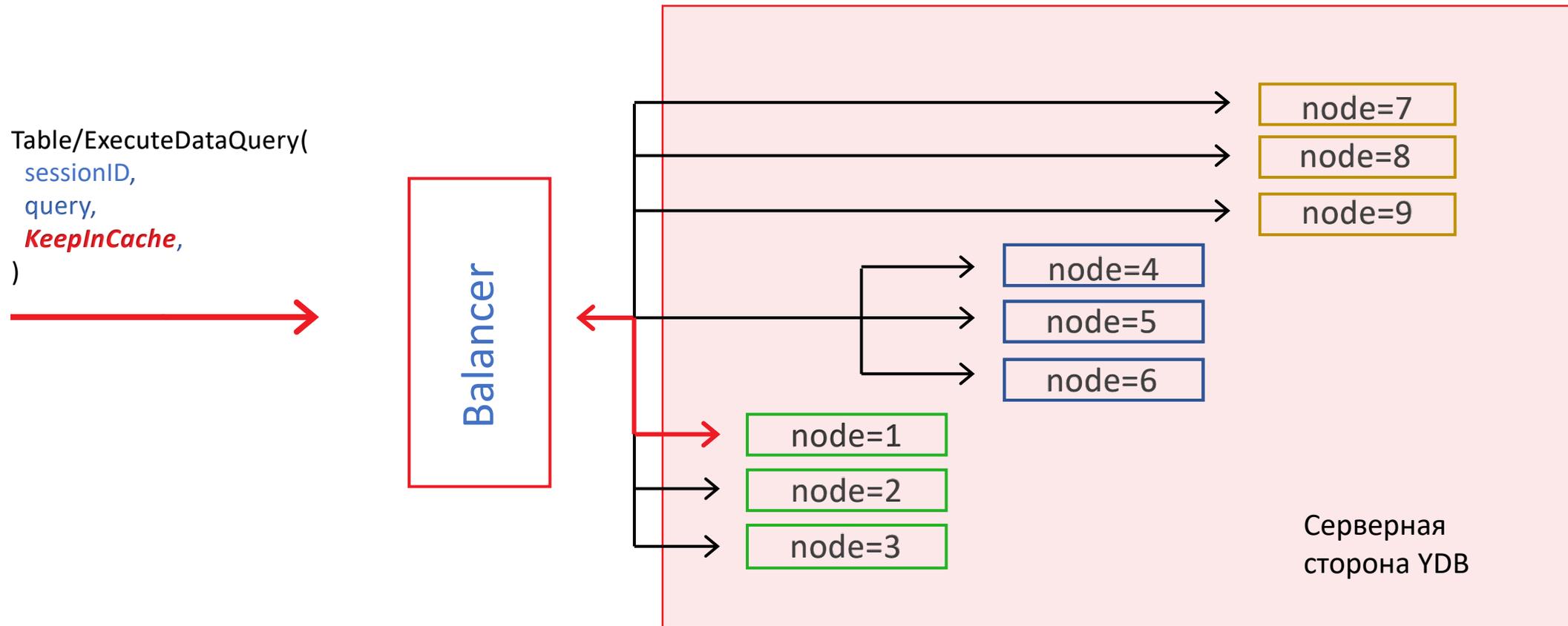
3. Закрыть сессию-Положить в пул

Выполнение запроса

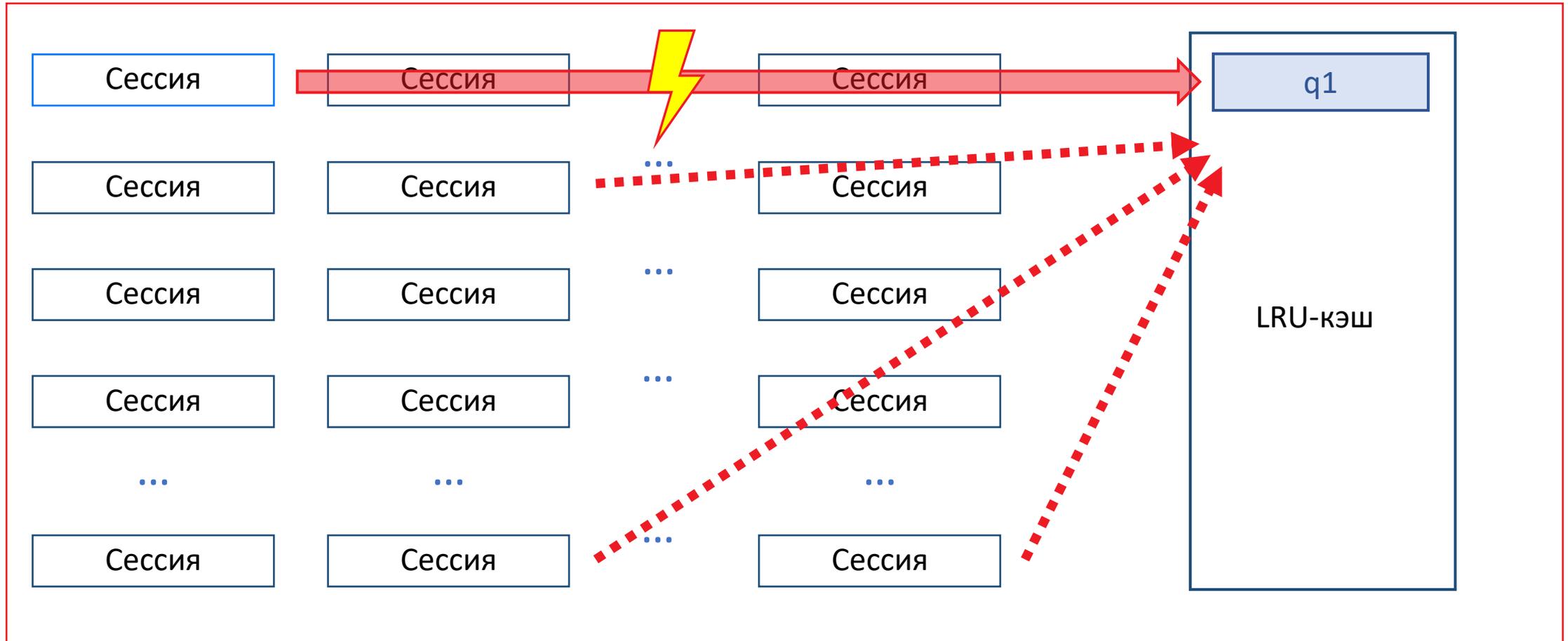
```
SELECT `title`, `content`  
FROM `/highload/2022/november/reports`  
WHERE speaker=$speakerName;
```

4. Кэширование результатов компиляции запроса

Флаг кэширования результатов компиляции запроса



Запрос закэширован на ноде 1



Первый запрос на ноде 1

1. Создать сессию
(взять из пула)

2. Выполнить запрос на сессии

3. Вернуть сессию в пул

Повторный запрос на ноде 1

1. Создать сессию
(взять из пула)

2. Выполнить запрос на сессии

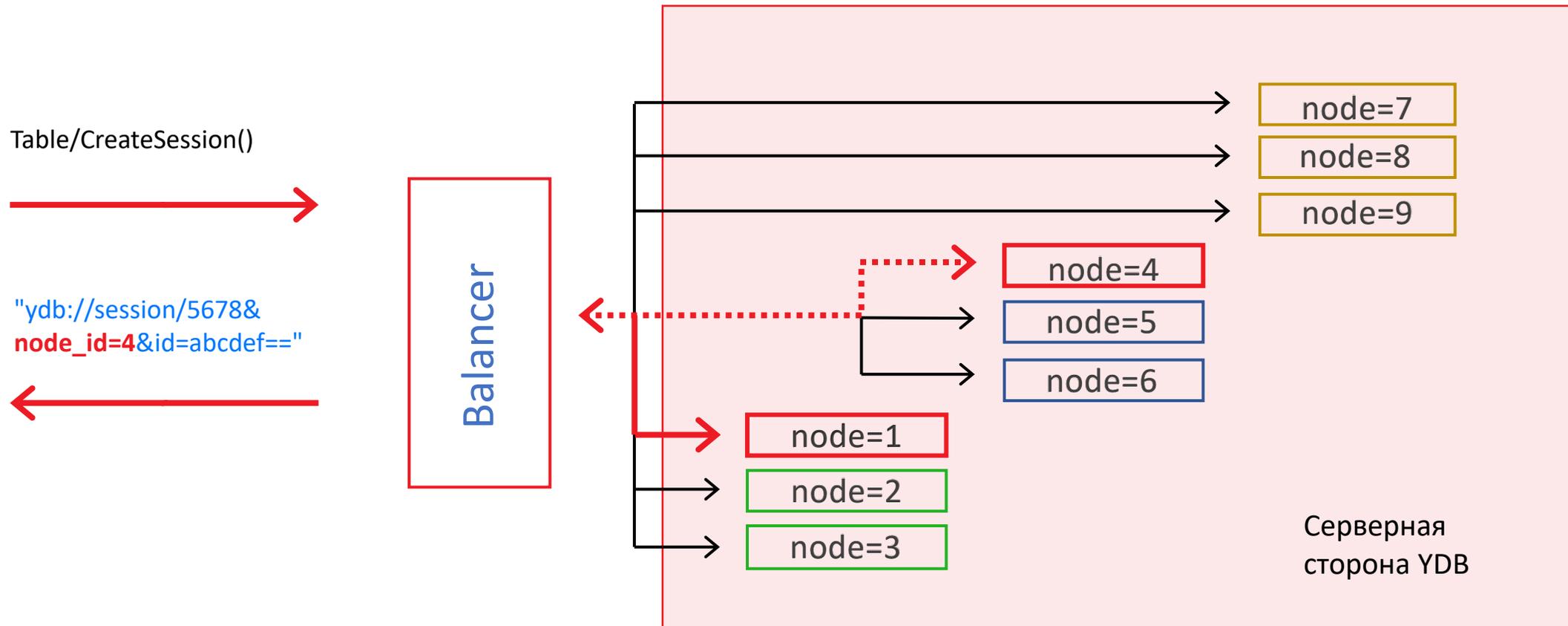
3. Вернуть сессию в пул

Флаг кэширования запросов

- + уменьшает общее время выполнения запроса
- + упрощает клиентский код
- + защищает от рестартов нод
- + защищает от вымывания серверного кэша
- + помогает равномерно использовать все ноды базы для запросов
- + включен по дефолту для всех запросов с параметрами

5. Серверная балансировка

5.1 Сессия может быть создана на другой ноде

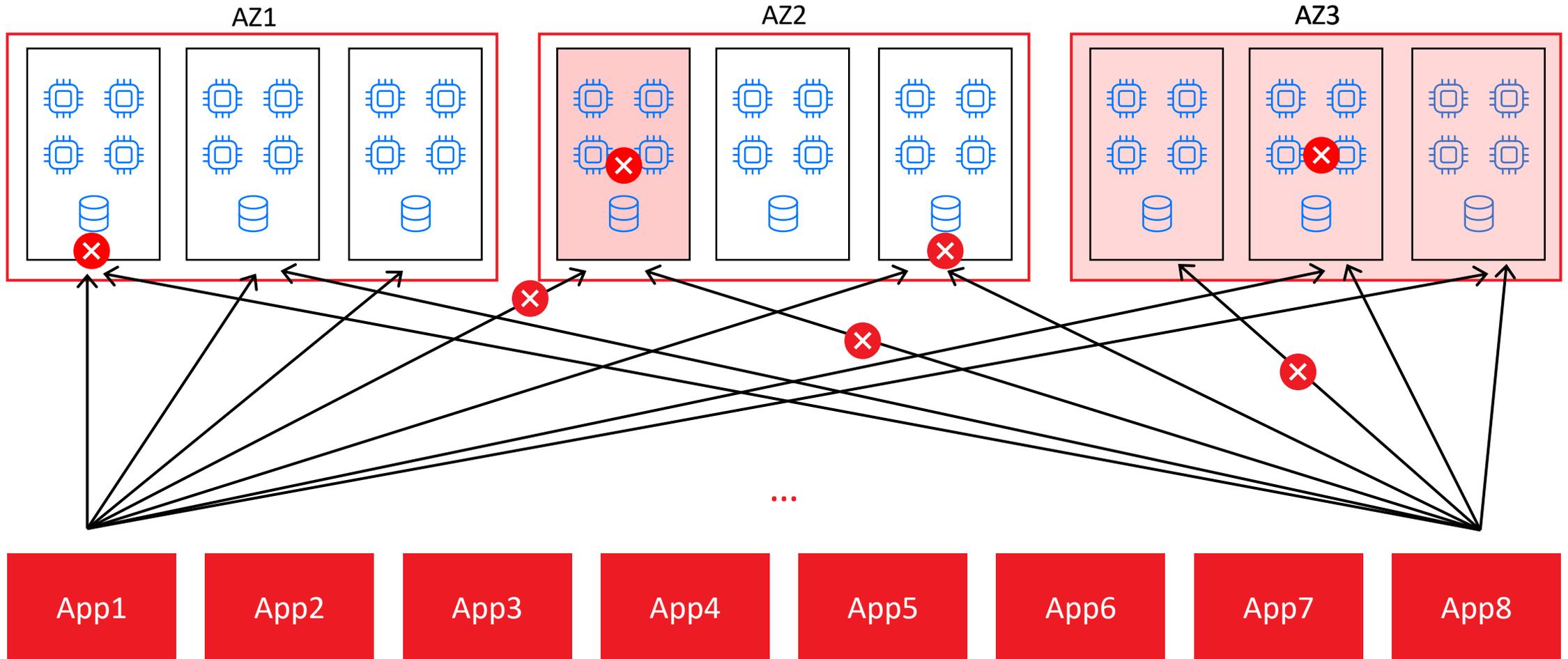


5.2 Сессии могут быть закрыты сервером принудительно



6. Обработка ошибок

Ожидаемые ошибки



Наивный ретраер

```
func retry(ctx, db, query) result {
    for {
        s, err := db.Table().CreateSession(ctx)
        if err != nil {
            continue
        }
        result, err := s.Execute(ctx, query)
        if err != nil {
            s.Close()
            continue
        }
        s.Close()
        return result
    }
}
```

Ретраер с лимитом попыток

```
func retry(ctx, db, query) result {
    for i := 0; i < 10; i++ {
        s, err := db.Table().CreateSession(ctx)
        if err != nil {
            continue
        }
        result, err := s.Execute(ctx, query)
        if err != nil {
            s.Close()
            continue
        }
        s.Close()
        return result, nil
    }
    return nil, fmt.Errorf("no progress")
}
```

Ретраер с экспоненциальной задержкой

```
func retry(ctx, db, query) result {
    for i := 0; i < 10; i++ {
        if i > 0 {
            time.Sleep(math.Pow(2, i) * time.Millisecond)
        }
        s, err := db.Table().CreateSession(ctx)
        if err != nil {
            continue
        }
        result, err := s.Execute(ctx, query)
        if err != nil {
            s.Close()
            continue
        }
        s.Close()
        return result, nil
    }
    return nil, fmt.Errorf("no progress")
}
```

Транспортные и серверные ошибки

Canceled

Unauthenticated

Internal

☐ - транспортные ошибки

Unknown

ResourceExhausted

Unavailable

☐ - серверные ошибки

InvalidArgument

FailedPrecondition

DataLoss

DeadlineExceeded

Aborted

Unimplemented

AlreadyExists

OutOfRange

PermissionDenied

BAD_REQUEST

GENERIC_ERROR

UNAUTHORIZED

TIMEOUT

INTERNAL_ERROR

BAD_SESSION

ABORTED

PRECONDITION_FAILED

CANCELLED

[Ссылка](#)

UNAVAILABLE

ALREADY_EXISTS

UNDETERMINED

OVERLOADED

NOT_FOUND

UNSUPPORTED

SCHEME_ERROR

SESSION_EXPIRED

SESSION_BUSY

Ретраибельные ошибки

Canceled	Unauthenticated	Internal	 - транспортные ошибки  - серверные ошибки  - можно безопасно ретраить  - можно ретраить, если операция идемпотентная
Unknown	ResourceExhausted	Unavailable	
InvalidArgument	FailedPrecondition	DataLoss	
DeadlineExceeded	Aborted	Unimplemented	
AlreadyExists	OutOfRange		
PermissionDenied			

BAD_REQUEST	GENERIC_ERROR		
UNAUTHORIZED	TIMEOUT		
INTERNAL_ERROR	BAD_SESSION		
ABORTED	PRECONDITION_FAILED	CANCELLED	Ссылка
UNAVAILABLE	ALREADY_EXISTS	UNDETERMINED	
OVERLOADED	NOT_FOUND	UNSUPPORTED	
SCHEME_ERROR	SESSION_EXPIRED	SESSION_BUSY	

Следует удалить сессию

Canceled	Unauthenticated	Internal	 - транспортные ошибки  - серверные ошибки  - сессия более непригодна
Unknown	ResourceExhausted	Unavailable	
InvalidArgument	FailedPrecondition	DataLoss	
DeadlineExceeded	Aborted	Unimplemented	
AlreadyExists	OutOfRange		
PermissionDenied			
BAD_REQUEST	GENERIC_ERROR		
UNAUTHORIZED	TIMEOUT		
INTERNAL_ERROR	BAD_SESSION		
ABORTED	PRECONDITION_FAILED	CANCELLED	Ссылка
UNAVAILABLE	ALREADY_EXISTS	UNDETERMINED	
OVERLOADED	NOT_FOUND	UNSUPPORTED	
SCHEME_ERROR	SESSION_EXPIRED	SESSION_BUSY	

С «медленной» экспоненциальной задержкой

Canceled	Unauthenticated	Internal	 - транспортные ошибки  - серверные ошибки  - экспоненциальная задержка между попытками от 1с
Unknown	ResourceExhausted	Unavailable	
InvalidArgument	FailedPrecondition	DataLoss	
DeadlineExceeded	Aborted	Unimplemented	
AlreadyExists	OutOfRange		
PermissionDenied			

BAD_REQUEST	GENERIC_ERROR	
UNAUTHORIZED	TIMEOUT	
INTERNAL_ERROR	BAD_SESSION	
ABORTED	PRECONDITION_FAILED	CANCELLED
UNAVAILABLE	ALREADY_EXISTS	UNDETERMINED
OVERLOADED	NOT_FOUND	UNSUPPORTED
SCHEME_ERROR	SESSION_EXPIRED	SESSION_BUSY

[Ссылка](#)

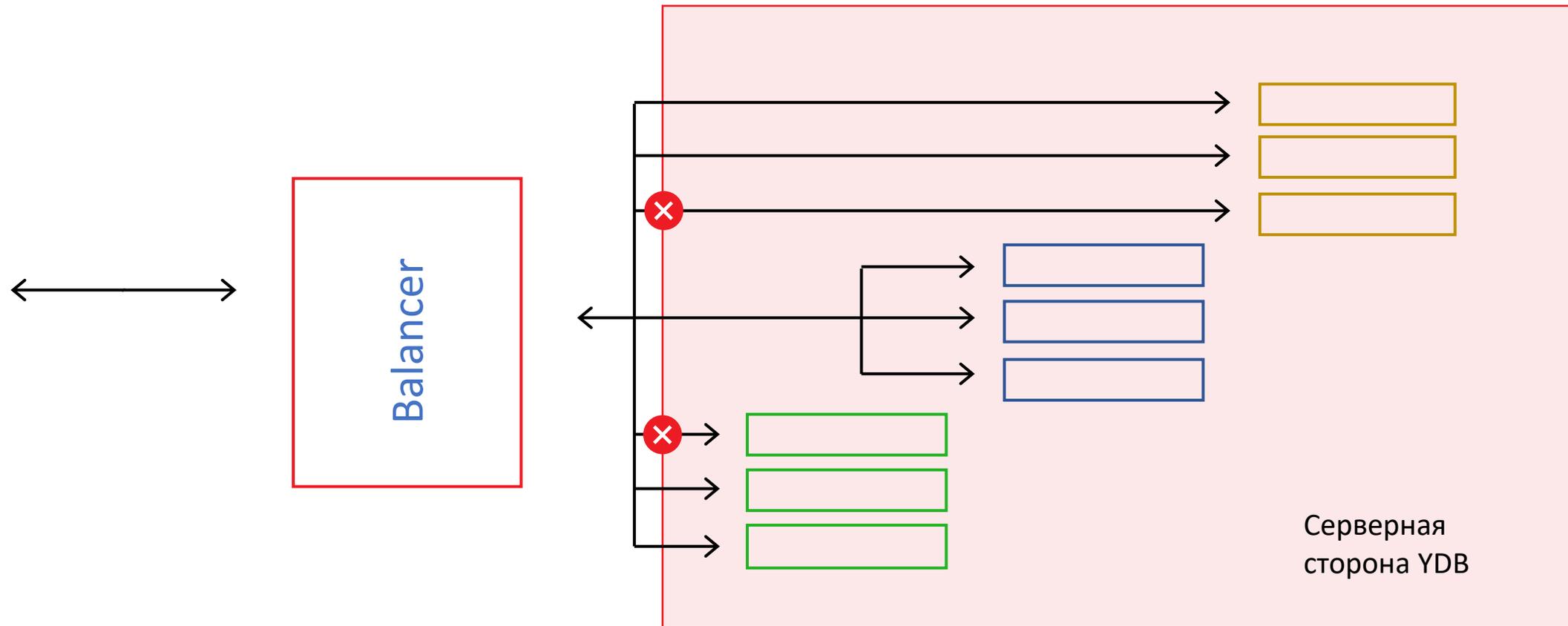
С «быстрой» экспоненциальной задержкой

Canceled	Unauthenticated	Internal	 - транспортные ошибки  - серверные ошибки  - экспоненциальная задержка между попытками от 5мс
Unknown	ResourceExhausted	Unavailable	
InvalidArgument	FailedPrecondition	DataLoss	
DeadlineExceeded	Aborted	Unimplemented	
AlreadyExists	OutOfRange		
PermissionDenied			

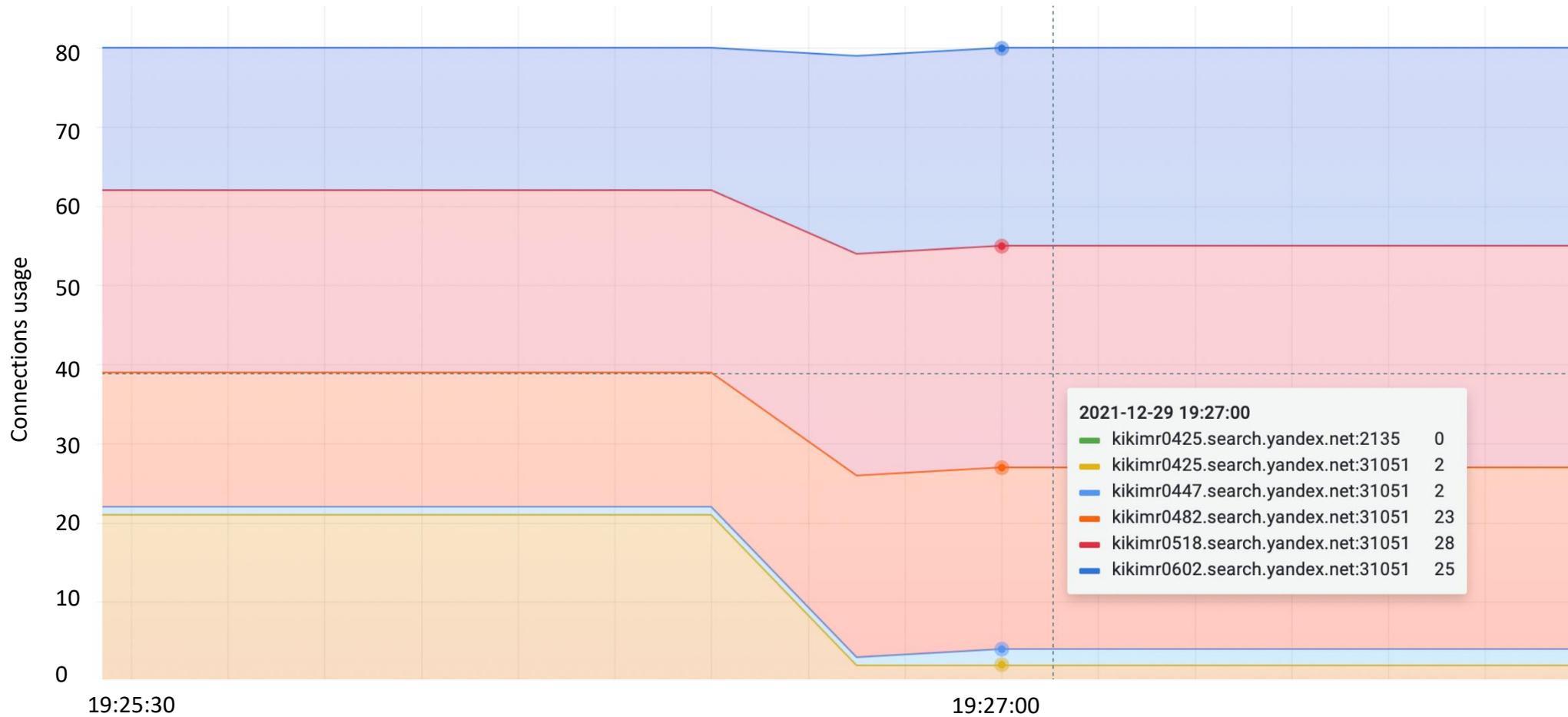
BAD_REQUEST	GENERIC_ERROR	
UNAUTHORIZED	TIMEOUT	
INTERNAL_ERROR	BAD_SESSION	
ABORTED	PRECONDITION_FAILED	CANCELLED
UNAVAILABLE	ALREADY_EXISTS	UNDETERMINED
OVERLOADED	NOT_FOUND	UNSUPPORTED
SCHEME_ERROR	SESSION_EXPIRED	SESSION_BUSY

[Ссылка](#)

Пессимизация соединений



Пессимизация соединений



Пессимизация соединений – путь в один конец?

- + Discovery/ListEndpoints === source of truth
- + Force re-discovery, если пессимизировано более 50% соединений

«Умные» ретраеры в драйвере YDB

```
err := db.Table().Do(ctx, func(ctx, session) error {
    result, err := s.Execute(ctx, query)
    if err != nil {
        return err
    }
    var title, content string
    for result.NextResultSet(ctx) {
        for result.NextRow() {
            if err := result.Scan(&title, &content); err != nil {
                return err
            }
            log.Println(title, content)
        }
    }
    return result.Err()
}, table.WithIdempotent(true))
```

Retry operation возвращает ошибку для обработки на стороне драйвера YDB

```
err := db.Table().Do(ctx, func(ctx, session) error {
    result, err := s.Execute(ctx, query)
    if err != nil {
        return err
    }
    var title, content string
    for result.NextResultSet(ctx) {
        for result.NextRow() {
            if err := result.Scan(&title, &content); err != nil {
                return err
            }
            log.Println(title, content)
        }
    }
    return result.Err()
}, table.WithIdempotent(true))
```

Сущности, порожденные волатильной сессией, также волатильные

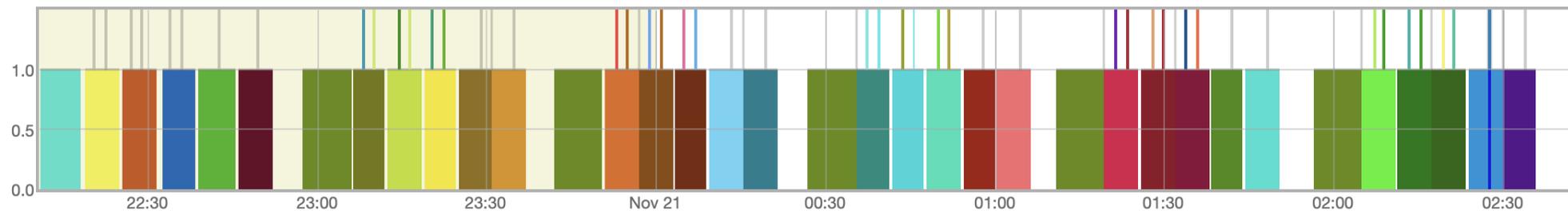
```
var resultOutOfRetryOperation result
err := db.Table().Do(ctx, func(ctx, session) error {
    result, err := session.Execute(ctx, query)
    if err != nil {
        return err
    }
    resultOutOfRetryOperation = result
    var title, content string
    for result.NextResultSet(ctx) {
        for result.NextRow() {
            if err := result.Scan(&title, &content); err != nil {
                return err
            }
            log.Println(title, content)
        }
    }
    return result.Err()
}, table.WithIdempotent(true))
```

Существенные параметры ретраера

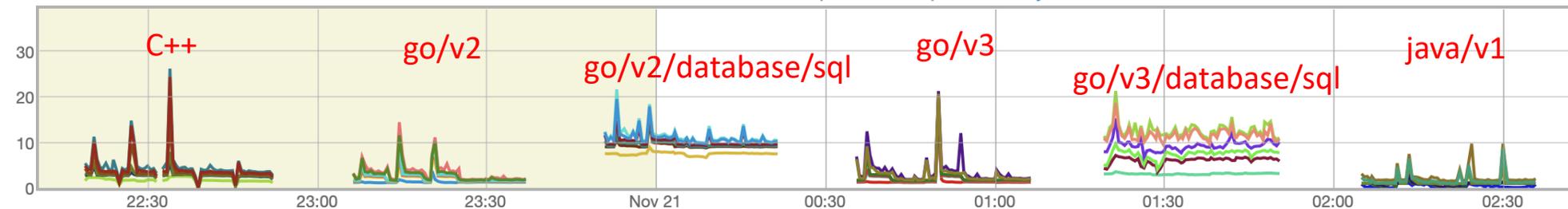
```
err := db.Table().Do(ctx, func(ctx, session) error {
    result, err := s.Execute(ctx, query)
    if err != nil {
        return err
    }
    var title, content string
    for result.NextResultSet(ctx) {
        for result.NextRow() {
            if err := result.Scan(&title, &content); err != nil {
                return err
            }
            log.Println(title, content)
        }
    }
    return result.Err()
}, table.WithIdempotent(true))
```

Тестирование надежности драйверов YDB

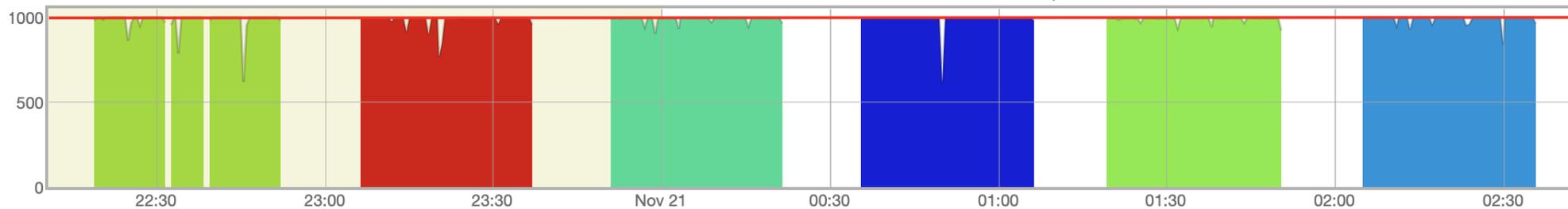
Tasks



Requests Ok replies latency measured on client, ms



Successful requests



Драйвер распределенной базы данных YDB

- + умеет инициализироваться
- + запускает фоновый процесс актуализации состояния YDB
- + реализует клиентскую балансировку запросов
- + осуществляет привязку сессий и нод YDB
- + имеет пул сессий и выполняет фоновый KeepAlive
- + корректно обрабатывает ошибки
- + имеет ретраеры
- + поддерживает пессимизацию соединений
- + готов к серверной балансировке

Feature parity

Feature	C++	Python	Go	Java	NodeJS	C#	Rust
Поддержка SSL/TLS (системные сертификаты)	+	+	+	+	+	+	+
Поддержка SSL/TLS (кастомные сертификаты)	+	+	+	+	+		-
Возможность настроить/включить GRPC KeepAlive (фоновое поддержание живости соединения)	+	+	+	?			-
Регулярный прогон тестов SLO на последней версии кода	+	+/-	+	+	+/-	-	-
Шаблоны Issue в GitHub	-	?	+	-	+		-
Клиентская балансировка							
Инициализация балансировщика через Discovery/ListEndpoints	+	+	+	+	+	+	+
Открытие клиентской балансировки (все запросы в начальный Endpoint)	+/-	-	+	-			
Фоновый Discovery/ListEndpoints (раз в минуту)	+	+	+	+	+	+	+

[Ссылка на таблицу](#)

05 • Жизненный цикл драйвера YDB (все запросы в начальный Endpoint)

Яндекс



HighLoad++
2022

Наши официальные репозитории



<https://github.com/ydb-platform/ydb-go-sdk>



<https://github.com/ydb-platform/ydb-nodejs-sdk>



<https://github.com/ydb-platform/ydb-rs-sdk>



<https://github.com/ydb-platform/ydb-python-sdk>



<https://github.com/ydb-platform/ydb-java-sdk>



<https://github.com/ydb-platform/ydb-dotnet-sdk>



<https://github.com/ydb-platform/ydb-php-sdk>

Алексей Мясников

старший разработчик в команде YDB

✉ asmyasnikov@yandex-team.ru

📍 <https://t.me/asmyasnikov>



<https://clck.ru/32kaXr>



HighLoad⁺⁺
2022

Яндекс

Голосуйте за доклад

