

PGConf.Russia 2023



Пять оттенков шардинга

Обзор современных методов
горизонтального масштабирования
реляционных СУБД



Павел Конотопов

руководитель кластерной группы департамента
внедрения и технической поддержки

- 20+ лет в ИТ
- Инженер по отказоустойчивости PostgreSQL
- Последние пять лет работаю с PostgreSQL
- Последние два года работаю в Postgres Professional

Email:

p.konotopov@postgrespro.ru

Что такое шардирование

- Общепринятой дефиниции не существует
- Техника горизонтального масштабирования
- Заключается в разделении таблиц на более мелкие части
- Распределение частей таблиц по разным серверам (шардам)
- Шард содержит часть данных таблиц
- Запрос может быть выполнен параллельно
- Может быть реализовано различными способами

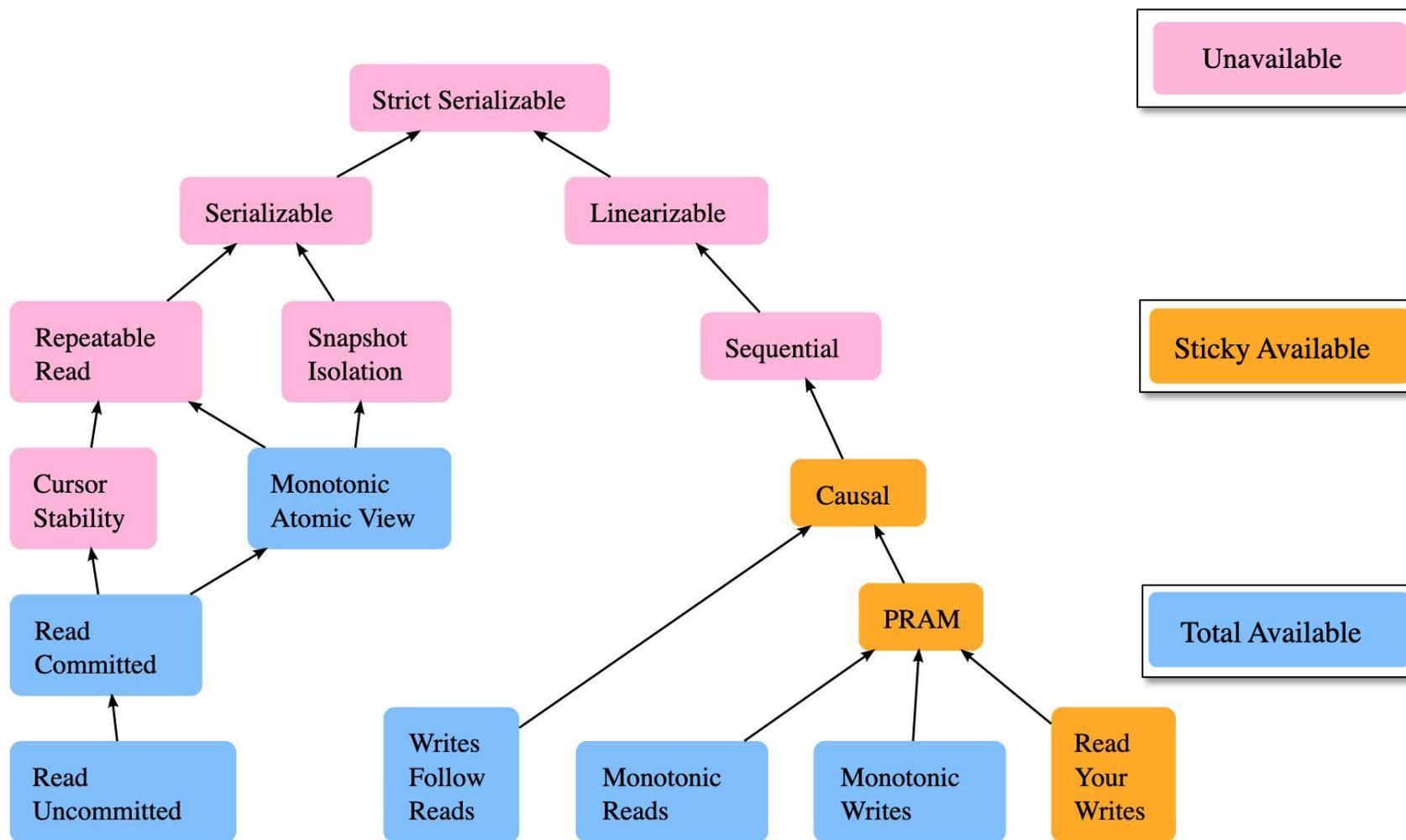
Требования

- Консистентность
- Отказоустойчивость
- Соответствует какому-то типу нагрузки
 - OLTP
 - OLAP
 - HTAP
- Удобство использования
- Простота администрирования

Терминология

- **Сериализуемость (serializability)**
 - Уровень изоляции
 - каждая операция записи или чтения будет возвращать результаты, соответствующие какой-то линейной последовательности операций, которые могут происходить одновременно
 - согласованность данных в контексте *параллельного* выполнения транзакций
- **Линеаризуемость (linearizability)**
 - Уровень консистентности
 - Тоже самое, что сериализуемость
 - Порядок должен учитывать реальное время выполнения
 - Гарантии не для транзакций, а для простых операций над одним объектом
- **Строгая сериализуемость (strict serializable)**
 - Сочетание уровня изоляции **serializability** и уровня консистентности **linearizability**

Модели консистентности



Unavailable

- Недоступно при некоторых типах сетевых сбоев. Узлы должны приостановить работу для обеспечения безопасности.

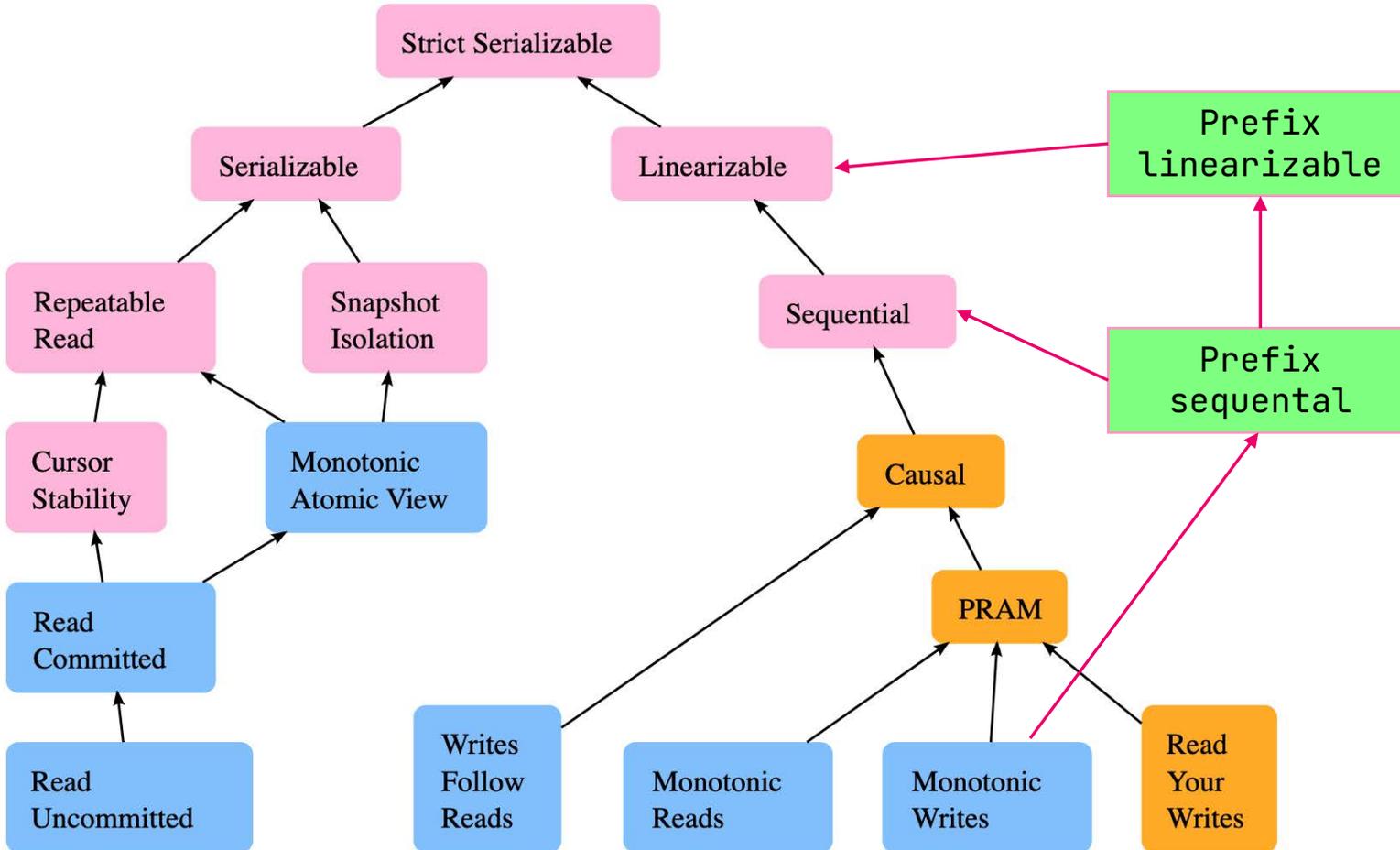
Sticky Available

- Доступен на всех узлах, не имеющих сбоев, пока клиенты общаются только с теми же серверами, а не переключаются на новые.

Total Available

- Доступен на каждом неисправном узле, даже когда сеть полностью отключена.

Модели консистентности



- **Одиночный сервер**
 - serializable
 - level linearizable
- **С синхронной репликацией (remote_apply)**
 - serializable
 - write linearizable
 - read sequential
- **С асинхронной репликацией**
 - Serializable
 - Read-write linearizable
 - Read-only prefix linearizable

Использование и администрирование

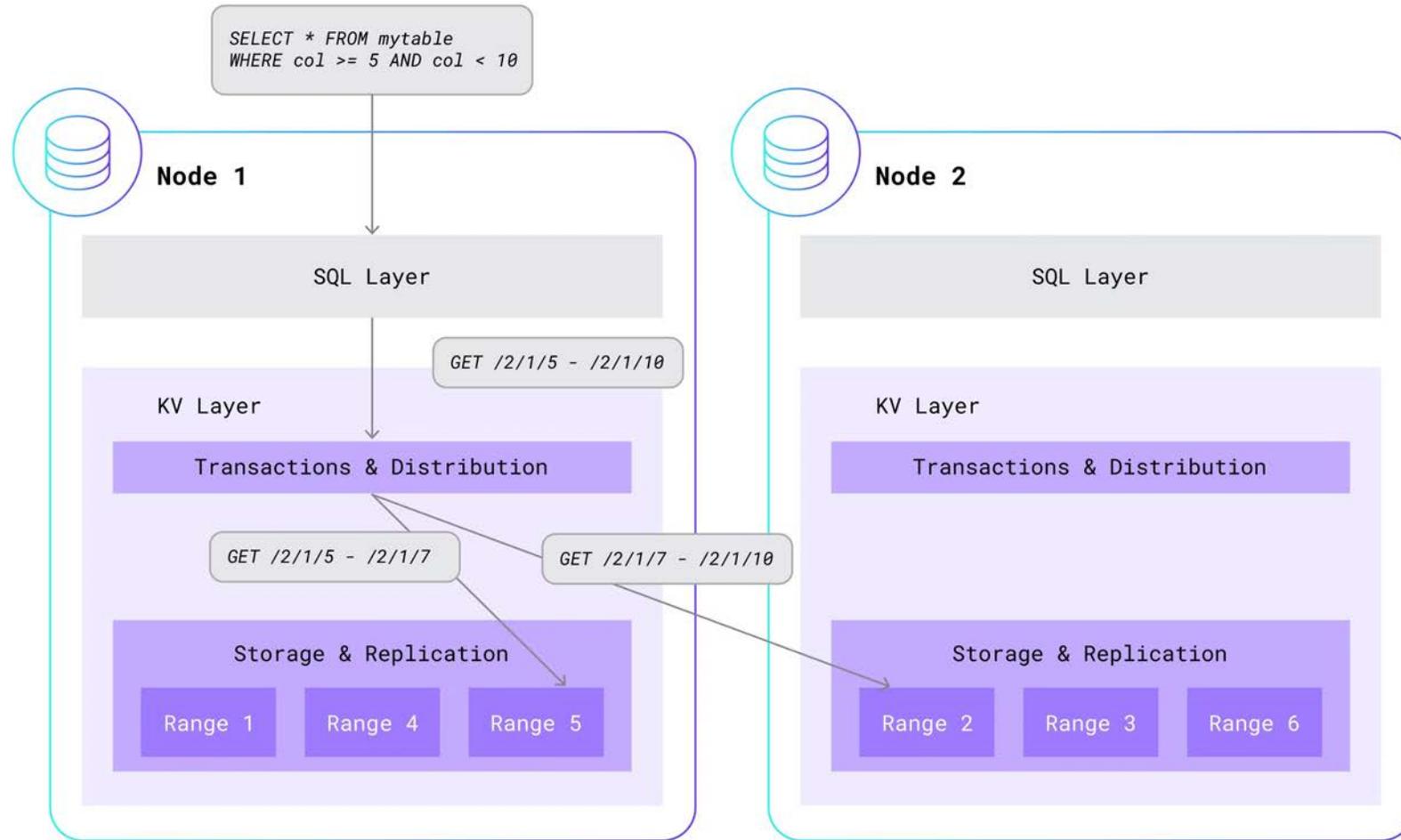
- Чем меньше ограничений тем лучше
- Чем меньше дополнительной инфраструктуры тем лучше
- Доступ в кластер через любой узел
- Минимум специального SQL диалекта
- Максимум совместимости с одиночным сервером PostgreSQL
- Не нужно думать о шардинге с точки зрения SQL

- Наличие cli-утилиты
- Мониторинг
- Удобство развертывания
 - on-premise, cloud, containers

CockroachDB

- Предназначен для геораспределенных данных
- Встроенная отказоустойчивость
 - Один бинарный файл
- Диалект SQL
 - DistSQL
 - Частично совместим с PostgreSQL
- Не основан на коде PostgreSQL
- Хранилище данных - поверх K/V RocksDB (своя реализация - Pebble)
- Предназначена OLTP, но не слишком быстро

- Модель консистентности:
 - Почти `strict serializable`
 - `stale reads` невозможен
 - `causal reverse` возможен



СockroachDB, ограничения

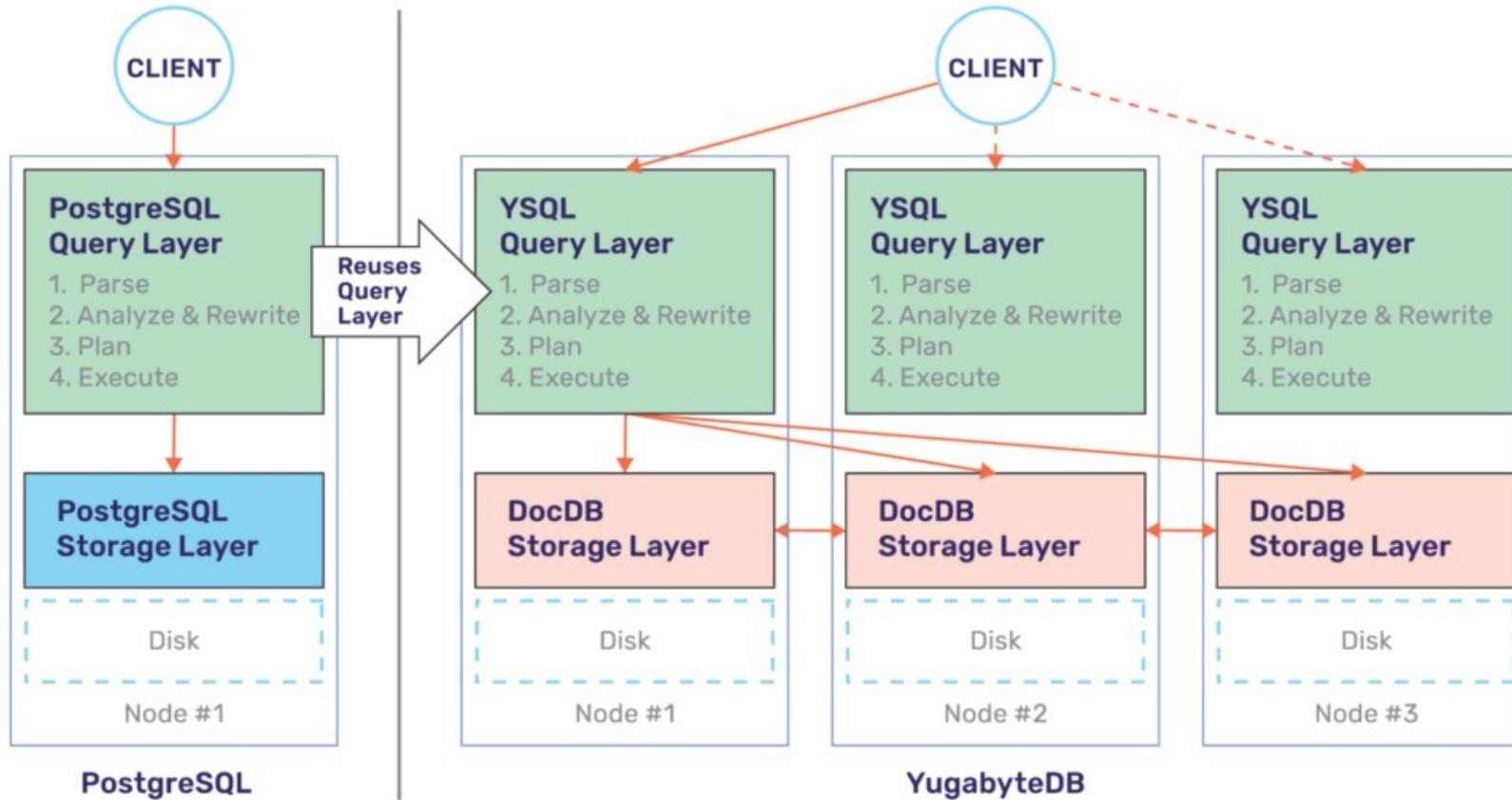
- Хранимые процедуры
 - Поддерживаются `user-defined functions`
- Триггеры
- События (`listen/notify`)
- Полнотекстовый поиск и функции
- Удаление первичного ключа
- XML функции
- Ограничение прав доступа для колонок
- Создание БД из шаблона
- Удаление одной партиции
- FDW

YugabyteDB

- Очень похожа на CockroachDB по архитектуре
- Предназначена для геораспределенных данных
- Встроенная отказоустойчивость
 - отказ от доступности в пользу согласованности
- Диалекты SQL
 - YSQL (совместим с PostgreSQL)
 - YCQL (полуреляционный язык, основан на Cassandra Query Language)
- Частично построена на коде PostgreSQL 11
 - Parser, Analyser, Planner, Executor
- Хранилище данных - поверх K/V DocDb (RocksDB)
- Предназначена для OLTP, но не слишком быстро
- отказ от доступности в пользу согласованности

- YSQL поддерживает уровни изоляции
 - Serializable, Snapshot и Read Committed
- YCQL поддерживает только изоляцию Snapshot
- Модель консистентности
 - Почти strict serializable

YugabyteDB



How YugabyteDB reuses PostgreSQL query layer

YugabyteDB, ограничения

- Кодовая база PostgreSQL 11
- Блокировки уровня таблицы
- Наследование таблиц
- Констрейнты (`exclusion`, `deferrable`)
- GiST индексы
- GIN индексы по нескольким колонкам
- События (`Listen/Notify`)
- XML функции
- Ограничения ключей `Primary/Foreign` на `foreign` таблицах
- `GENERATED ALWAYS AS STORED`
- `CREATE ACCESS METHOD`
- `CREATE SCHEMA with elements`
- DDL в транзакции

Администрирование

CLUSTER OVERVIEW

Capacity Usage

0.0%

USED CAPACITY 0 B USABLE CAPACITY 64.5 GiB

Node Status

3

LIVE NODES

Replication Status

High Availability Alert Configurations User Management

VIEW: **NODE LIST** ▾

Live Nodes

ID ▾	ADDRESS ▾	UPTIME ▾	BYTES ▾
n1	10.178.7.166:26257	4 minutes	1.7 MiB
n2	10.178.7.185:26257	4 minutes	1.7 MiB
n3	10.178.7.158:26257	4 minutes	1.7 MiB

- Dashboard
- Universes
- Metrics
- Tasks
- Alerts
- Configs
- Admin

Alert Policies
Alert Destinations
Notification Channels
Health
Maintenance Windows

Create Alert Config ▾

Hide Filter

Name:

Type: Platform Universe

Metric Name:

Destination:

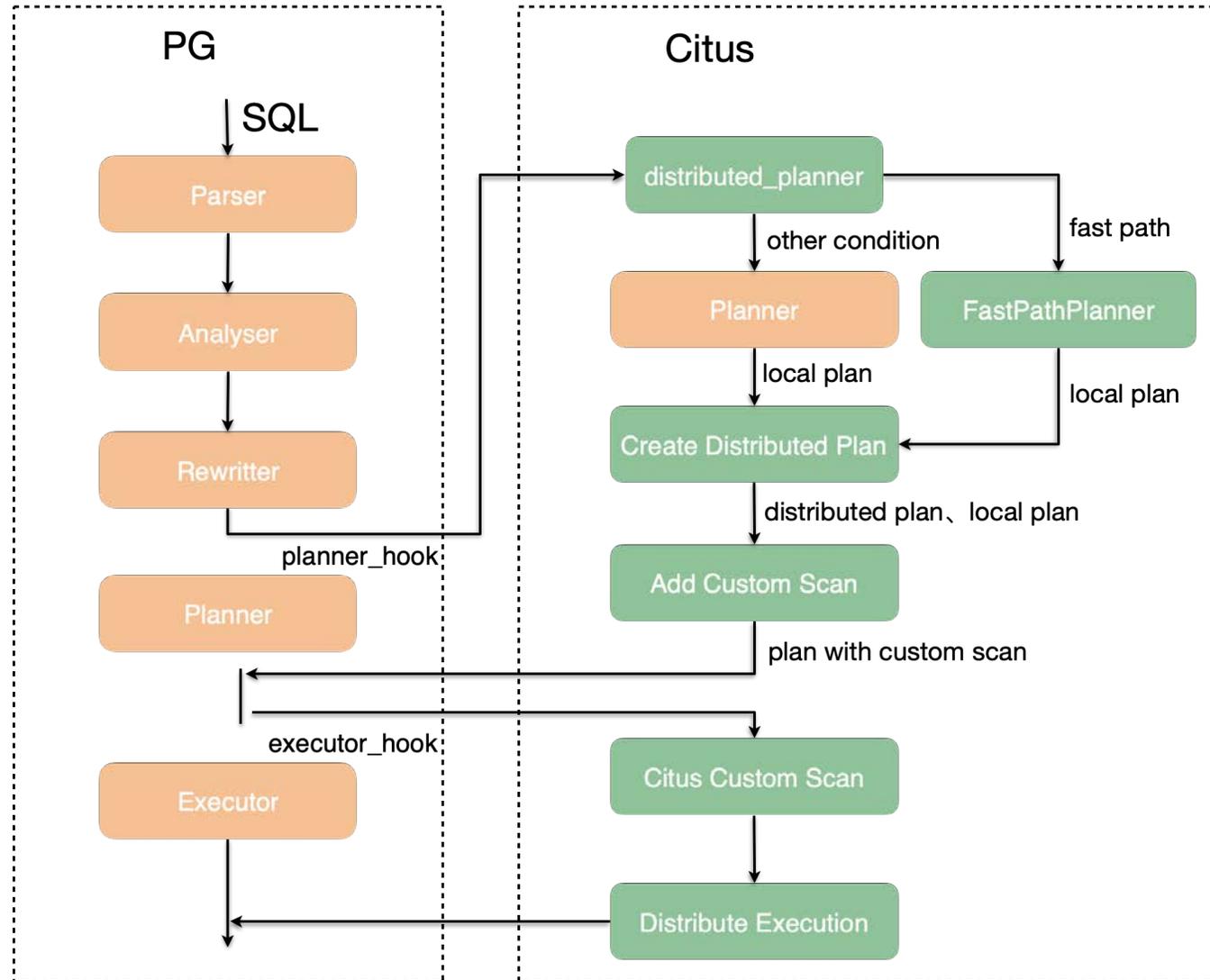
State: Active Inactive

Severity: Severe Warning

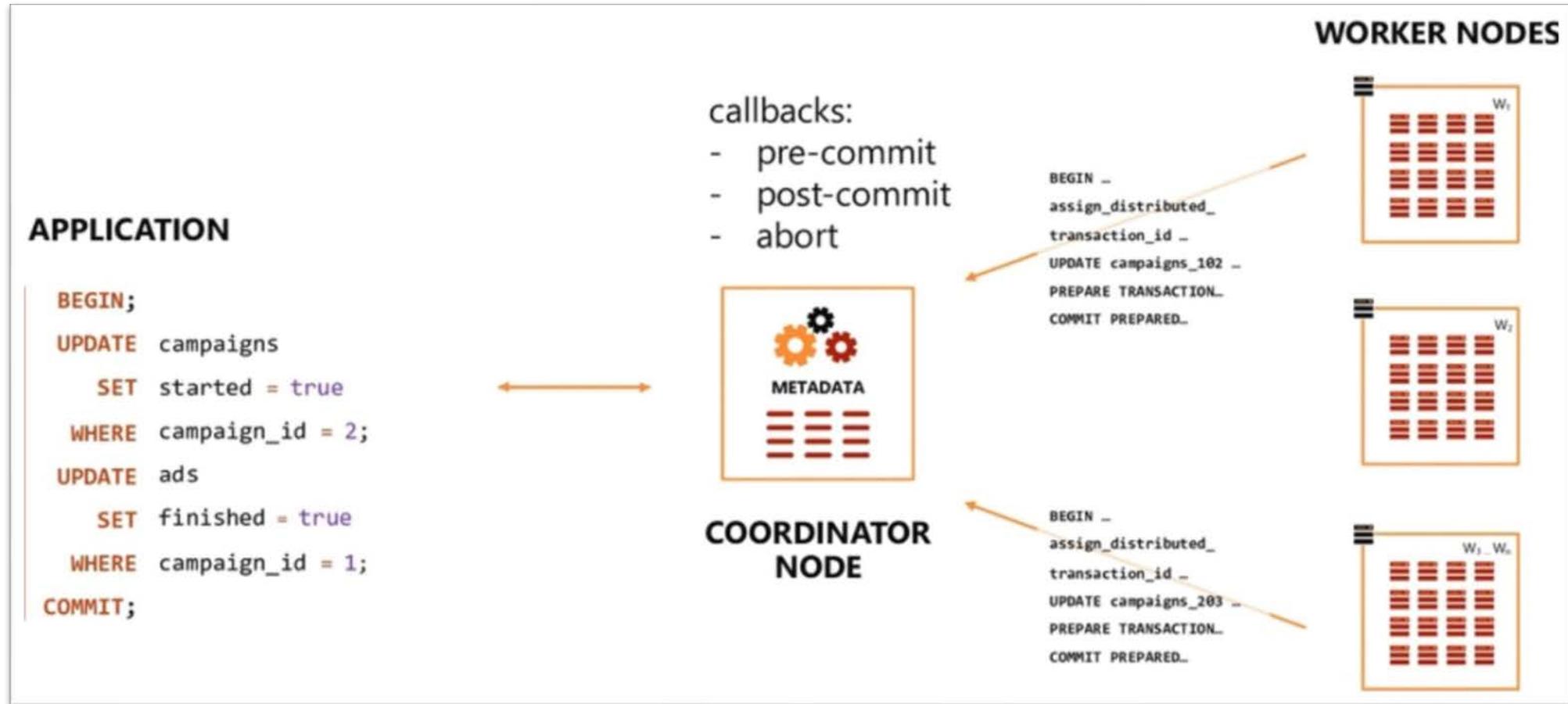
NAME ▲	TARGET UNIV...	SEVERITY ▾	DESTINATION ▾	CREATED ▾	METRIC NAME ▾	ACTIONS
Alert Channel Failed <small>PLATFORM</small>	ALL	SEVERE	MyAlertRoute-5 <small>DEFAULT</small>	08/02/2021	ALERT_NOTIFICATION_CHANNE...	...
Alert Notification Failed (Inactive) <small>PLATFORM</small>	ALL	SEVERE	MyDest	08/02/2021	ALERT_NOTIFICATION_ERROR	...
Alert Query Failed <small>PLATFORM</small>	ALL	SEVERE	MyAlertRoute-5 <small>DEFAULT</small>	08/02/2021	ALERT_QUERY_FAILED	...
Alert Rules Sync Failed <small>PLATFORM</small>	ALL	SEVERE	MyAlertRoute-5 <small>DEFAULT</small>	08/02/2021	ALERT_CONFIG_WRITING_FAILED	...
Alert Rules Sync Failed Testing <small>PLATFORM</small>	ALL	SEVERE	No destination	10/22/2021	ALERT_CONFIG_WRITING_FAILED	...
Backup Failure (Inactive) <small>UNIVERSE</small>	ALL	SEVERE	MyAlertRoute-5 <small>DEFAULT</small>	08/02/2021	BACKUP_FAILURE	...
Backup Schedule Failure <small>UNIVERSE</small>	ALL	SEVERE	MyAlertRoute-5 <small>DEFAULT</small>	08/24/2021	BACKUP_SCHEDULE_FAILURE	...
CA certificate expires in 30 days <small>UNIVERSE</small>	ALL	SEVERE	MyAlertRoute-6	11/08/2021	CLIENT_TO_NODE_CA_CERT_EX...	...

CitusDB

- Предназначен для локализованных данных
- Встроенной отказоустойчивости нет
 - Узел координатор
 - Узлы хранения данных
 - Поддержка в Patroni 3.x.x, k8s
 - Требуется дополнительная инфраструктура
 - Избыточность данных
- PostgreSQL внутри PostgreSQL
 - Начинаясь до ванильного партиционирования
 - Целились в OLAP
- Расширение PostgreSQL
 - Широко использует хуки PostgreSQL
- Предназначен для OLAP
 - Развивается в сторону OLTP
- Модель согласованности
 - Snapshot в рамках одного узла + 2PC
 - Read uncommitted – уровень изоляции
 - Linearizable – уровень консистентности



CitusDB



CitusDB, достоинства и недостатки

- Только SQL интерфейс для управления кластером
- Слабая модель согласованности
 - ACID на уровне отдельных узлов
 - Нет «честных» распределенных транзакций
- DDL только через координатор
- Нет консистентных резервных копий

- Ориентирован на OLAP, движется в сторону OLTP
 - Fastpath – быстрый разбор запроса и отправка его на нужный узел
 - Shuffle join через libpq
 - DML через любой узел (с 11 версии)
 - Нет мультиплексирования (интерконнект)
 - Требуется пулер соединений между узлами кластера

Greenplum

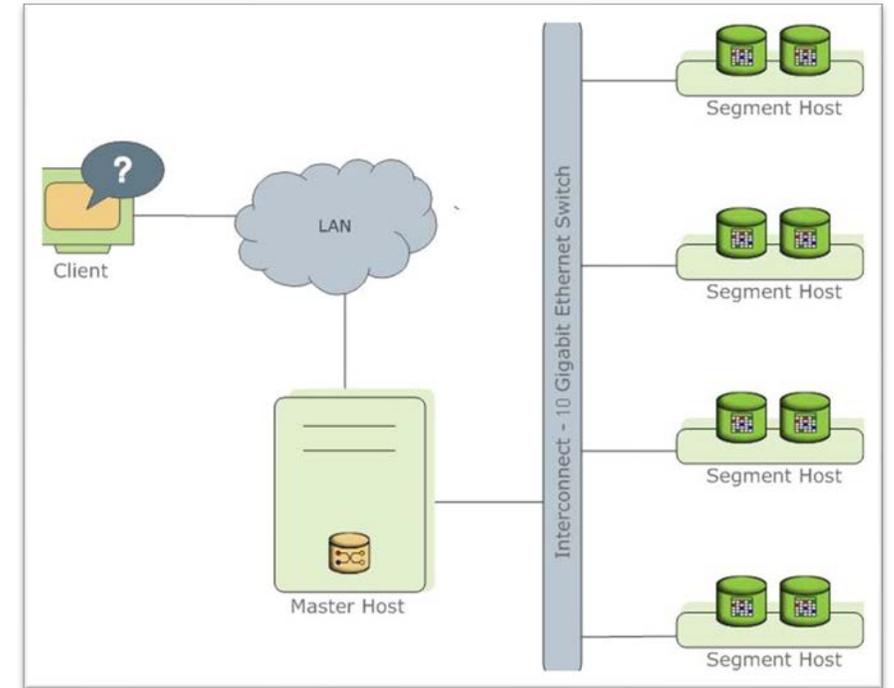
- Распределенная массивно-параллельная СУБД (MPP)
- Частичная встроенная отказоустойчивость
 - Избыточность данных
- Аналитическая система
- Форк, базируется на PostgreSQL 9.6-12
- Горизонтальное масштабирование
- Интерконнект
- Отказоустойчивость, обеспечивается созданием зеркал каждого логического сегмента и резервного мастер-сервера.
- Запросы – через координатор
- Возможность выполнять локальные и распределённые JOIN.

- Два уровня изоляции транзакций
 - Read committed и repeatable read

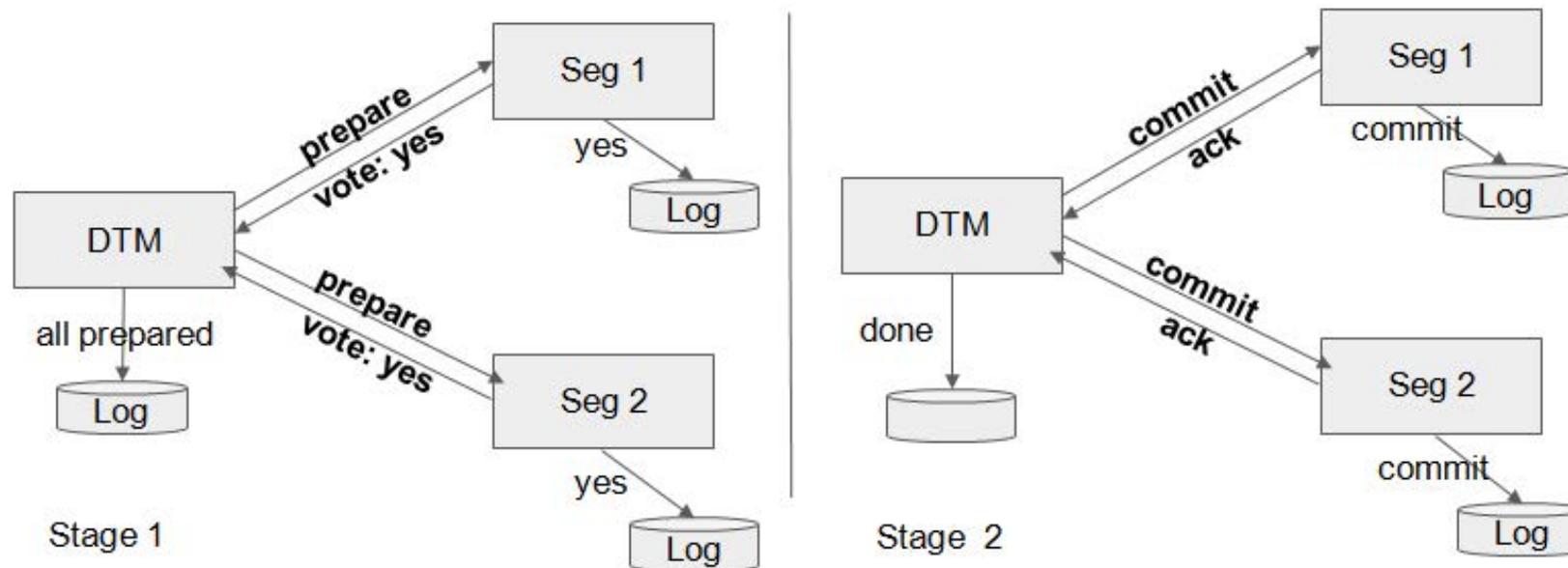
- Модель согласованности
 - Snapshot isolation
 - Linearizable

Greenplum

- **Мастер-сервер**
 - точка входа в Greenplum
 - координирует свою работу с сегментами
 - распределяет нагрузку между сегментами
 - не содержит пользовательских данных
- **Сервер-сегмент**
 - Хранит и обрабатывает данные
 - Отдает результаты мастеру
- **Интерконнект (внутрикластерный протокол)**
 - сетевое соединение для связи между отдельными экземплярами PostgreSQL



Cross node transaction: two phase commit



Greenplum, достоинства и недостатки

- Координатор строит распределенные планы выполнения запроса
- Распределенные транзакции через менеджер транзакций (DTM)
- Потеря координатора

- Хорошо поддержаны разные типы JOIN
 - Загрузка узлов кластера вычислительной работой
- cī для управления кластером, миграции данных, мониторинга
- Параллельный бекап-восстановление кластера
 - Access share lock – запрет изменений
 - Инкрементальный бекап
- Встроенная отказоустойчивость сегментов

SPQR/PGCAT

- Middleware
- Решают проблему для простых OLTP запросов
- Решает проблему роутинга и балансировки
- Пулер соединений
- Мониторинг реплик

Ванильный PostgreSQL

- Ванильное партицирование + `postgres_fdw`
- Ручное управление
 - Создать партицированную таблицу на всех узлах
 - Настроить `foreign servers`
 - Настроить чередование секций на каждом сервере
 - Связать секции через `postgres_fdw`
- Будет работать медленно
 - плохо работает `push-down` запросов
 - используется неоптимальный протокол `libpq` для интерконнекта
 - нет мультиплексирования, большое количество подключений
- Нет 2PC через `fdw` (идет работа)
- Модель консистентности
 - `Snapshot` в рамках одного узла
 - `Read uncommitted`
 - Возможно `linearizable`

Shardman

- Прозрачное горизонтальное масштабирование
- OLTP нагрузка
- Доступ в кластер через любой узел
- Избыточность данных (отказоустойчивость)
- Целостность данных (распределенные транзакции)

Shardman

- Расширение
- Секционирование
- `postgres_fdw`
- Мультиплексирование TCP соединений
- Физическая и логическая репликация
- CSN – `commit sequence number` (Clock-SI)

Shardman

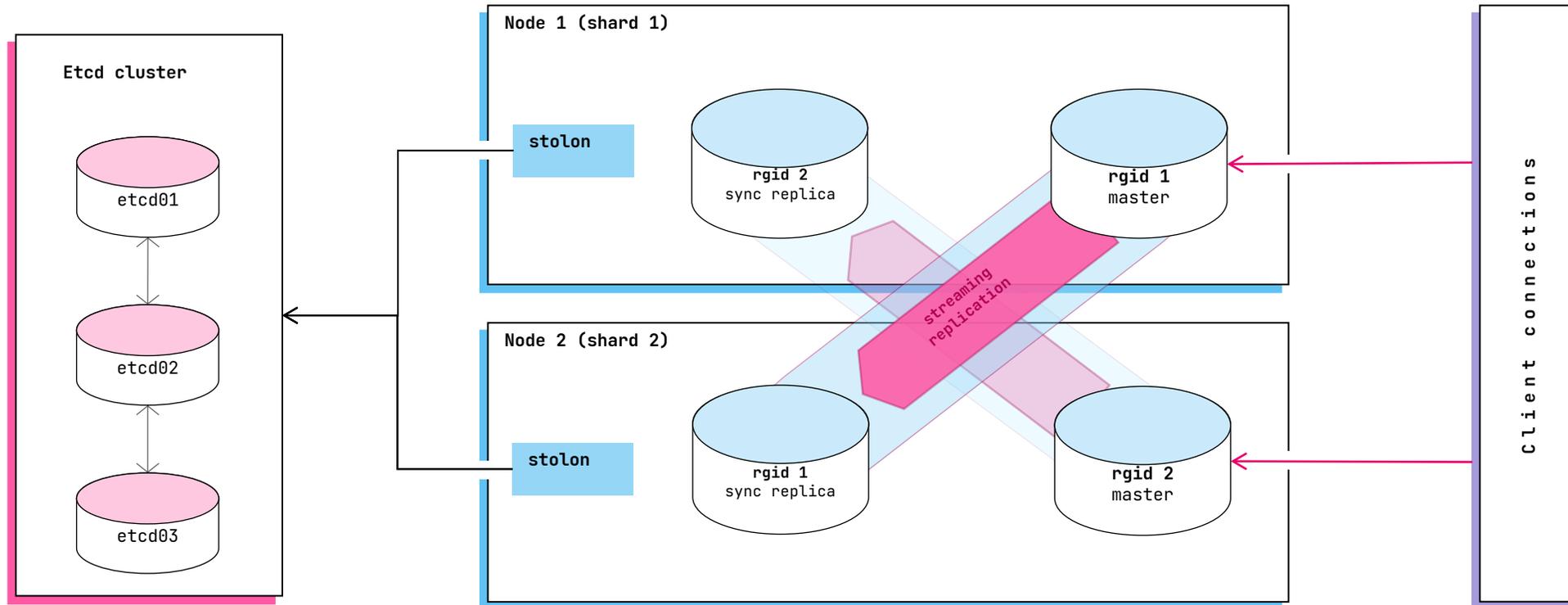
- Расширение
- Интегрированы возможности Postgres Pro Enterprise
 - CFS, PTRACK, pgpro_stats, AQO – адаптивный оптимизатор запросов
- Push-down запросов
 - Ассиметричный JOIN (партиция+глобальная таблица)
 - Агрегатные функции
 - Partition wise join (партиция + партиция)
- Доступ через любой из узлов кластера
 - Мультиплексирование соединений
- Консистентный бекап-восстановление
 - Инкрементальное, дельта и PTRACK
- Любая топология кластера
- Единая точка управления – cli утилита
- Модель консистентности
 - Изоляция Repeatable read (snapshot isolation)
 - Консистентность
 - write-only linearizable
 - read-write и read-only – prefix sequential

Shardman

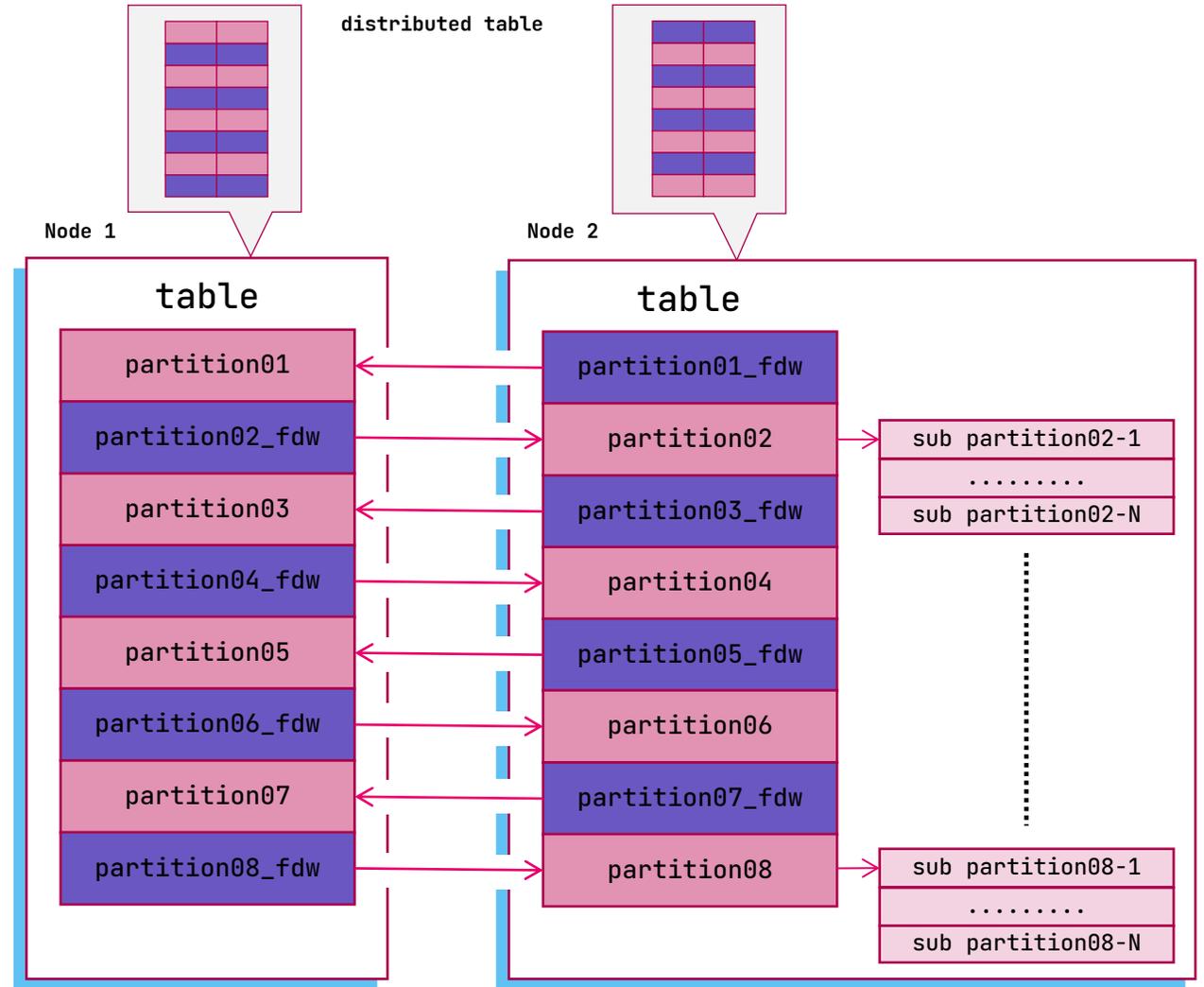
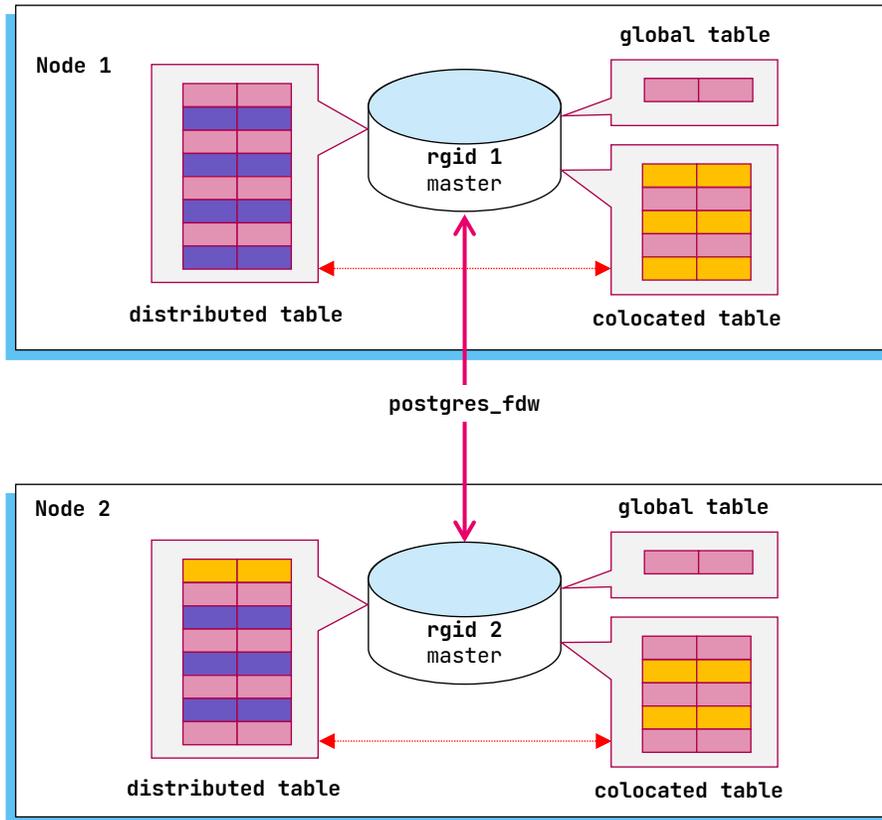
- PostgreSQL 14
- Stolon
- Расширения Shardman и postgres_fdw
- Etcd

- shardmand
сервис, отвечающий за кластерное ПО
- shardmanctl
утилита управления кластером

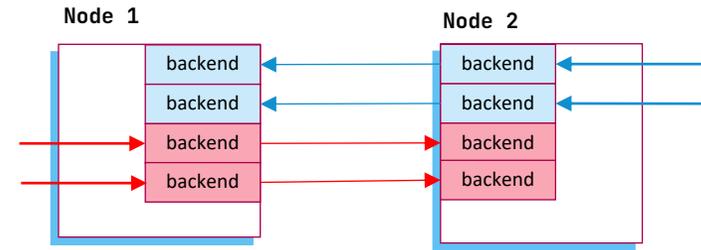
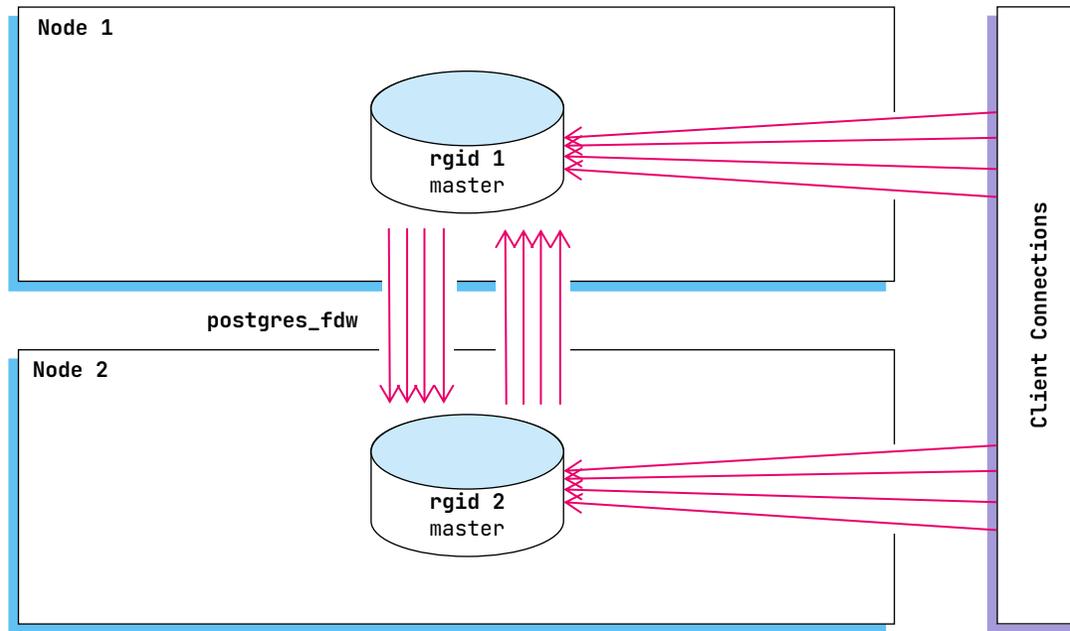
Shardman



Shardman

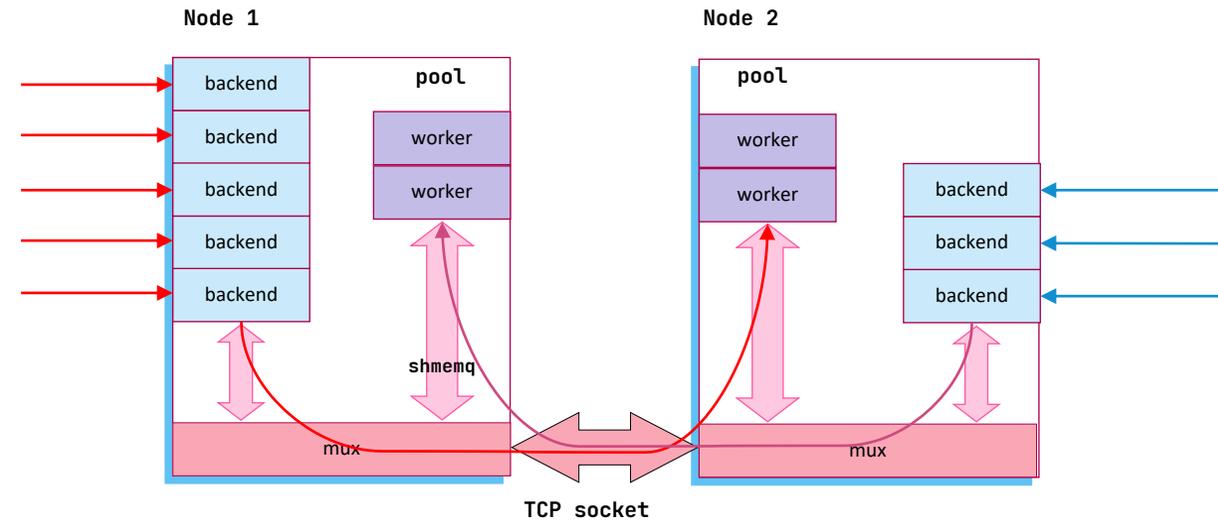
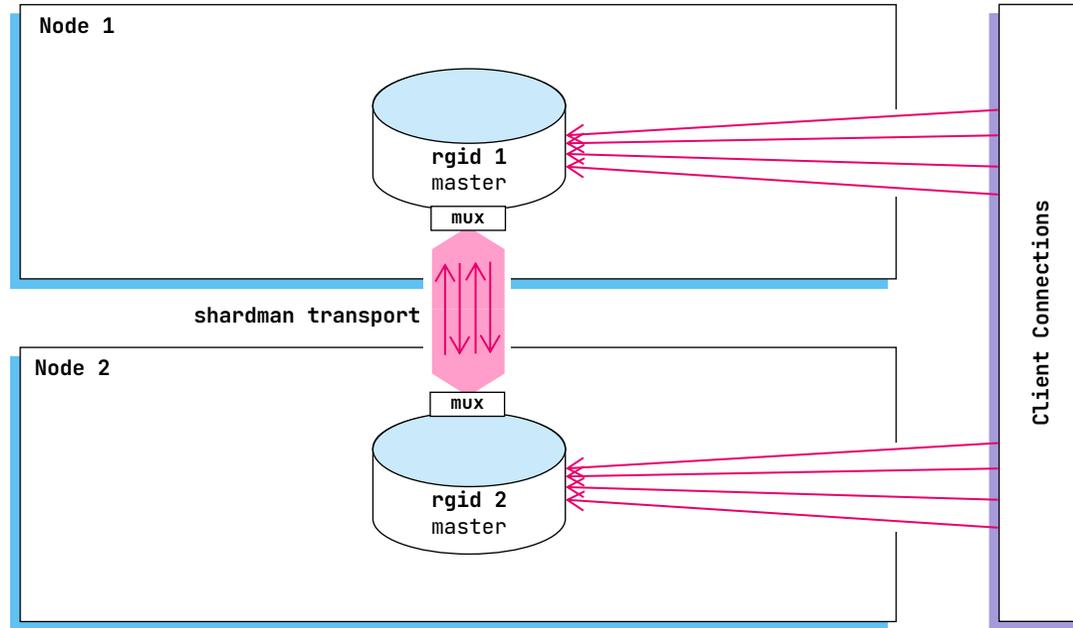


Shardman



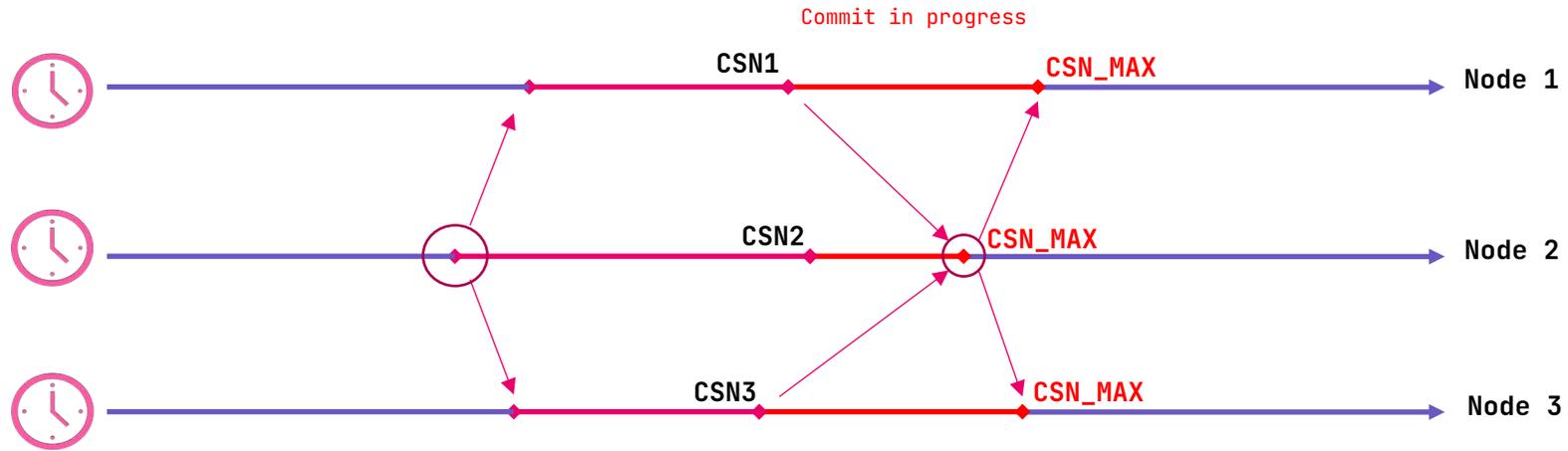
Общее число подключений в кластере – $M \times N$
 M – количество подключений,
 N – количество узлов

Shardman



Общее число подключений кластере - $M + (N \times N)$
 Общее число процессов в кластере - $M + (N \times W)$
 M - количество подключений,
 N - количество узлов,
 W - количество воркеров

Shardman



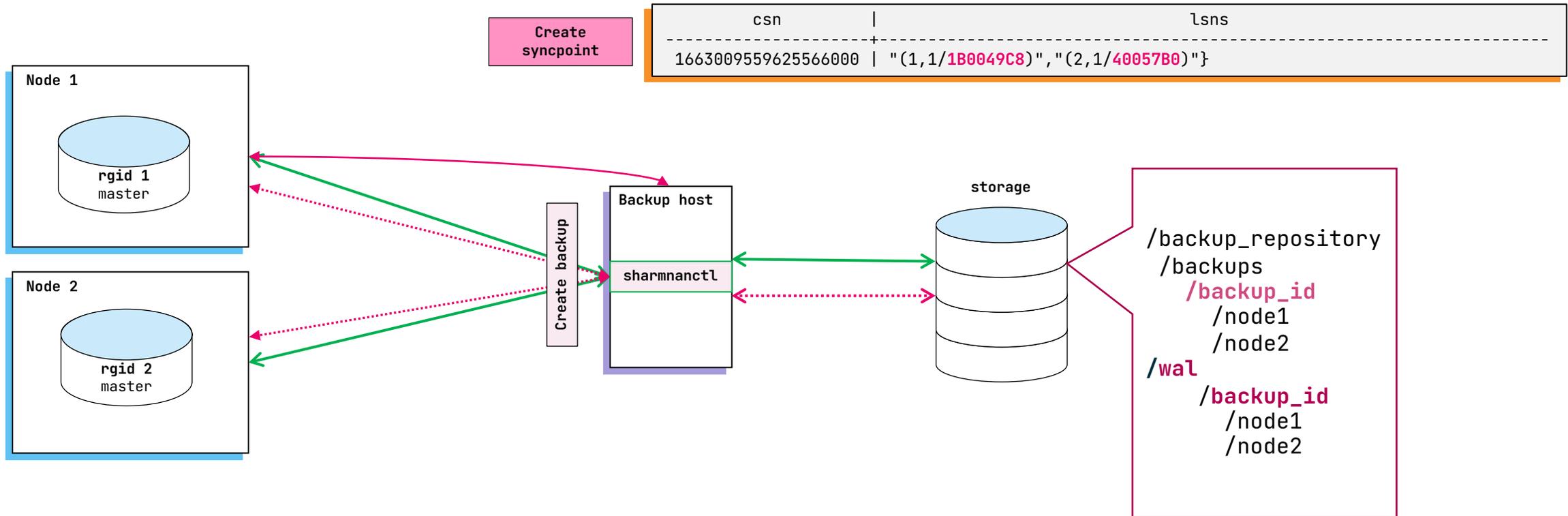
$$CSN_MAX = \text{MAX}(CSN1, CSN2, CSN3)$$

CSN - время
 CSN(LSN_NODE1, LSN_NODE2, LSN_NODE3)

```
# CALL shardman.create_syncpoint(NULL, NULL);
```

csn	lsns
1663009559625566000	{"(3,0/FD005D60)", "(1,1/1B0049C8)", "(2,1/40057B0)"}

Shardman



Shardman

- `shardmanctl` – cli утилита управления
`cluster/config/init/shards/nodes/rebalance/
backup/load/upgrade/update/status`
- Инициализация
- Добавление и удаление узлов
- Ребалансировка
- Бекап и восстановление
- Обновление до новой версии
- и т.д.

Shardman, ограничения

- Распределение данных – только хэшу
- Все формы ALTER TABLE невозможны, кроме:
 - Изменение владельца таблицы,
 - SET/DROP NOT NULL,
 - ADD/DROP COLUMN, кроме type serial (автоинкремент)
 - Добавление и удаления индексов
- Нельзя создавать распределенные временные таблицы
- Нельзя изменить количество секций распределенных таблиц
- Невозможно переименовывать глобальные роли
- Только один тип хранения данных – таблицы PostgreSQL
- Не позволяет использовать реплики для чтения данных
- Транспорт работает только для SELECT

Таблица сравнения

	CockroachDB	YugabyteDB	Greenplum	CitusDB	Shardman
Уровень изоляции	Почти strict serializable	Почти strict serializable	snapshot isolation	read uncommitted	snapshot isolation
Уровень консистентности	linearizable	linearizable prefix sequential	linearizable	linearizable	linearizable prefix sequential
Доступ в кластер	Любой узел	Любой узел	Координатор	Координатор (DDL), Любой узел (DML)	Любой узел
Отказоустойчивость	Да/Auto HA	Да/Auto HA	На уровне датанод/Auto HA/Mirroring	Нет/Patroni/Auto HA	Да
Тип нагрузки	OLTP	OLTP	OLAP	OLAP->HTAP	OLTP->HTAP
Удобство администрирования	UI/cli	UI/cli	cli	SQL интерфейс	cli
Наличие дополнительной инфраструктуры	Не требуется	Не требуется	Не требуется	Нет/DCS	DCS

Ссылки

- **Consistency models**
 - <https://jepsen.io/consistency>
 - <https://www.cockroachlabs.com/blog/consistency-model>
- **CockroachDB**
 - <https://www.cockroachlabs.com/docs>
- **YugabyteDB**
 - <https://docs.yugabyte.com>
- **CitusDB**
 - <https://docs.citusdata.com>
- **SPQR**
 - <https://github.com/pg-sharding/spqr>
- **Pgcat**
 - <https://github.com/postgresml/pgcat>
- **Shardman**
 - <http://repo.postgrespro.ru/doc/pgprosm/14.7.1/en/html>
 - <http://repo.postgrespro.ru/pgprosm-14>
 - <https://github.com/pkonotopov/shardman-docker>

Q & A

Спасибо!