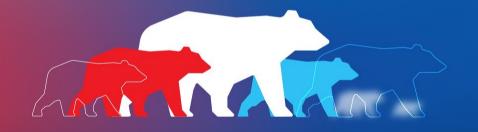
# Shardman — постгрес для кластеров. Что есть сейчас, что будет завтра

Андрей Лепихов, Postgres Professional



## **High**Load \*\* Siberia 2019

Профессиональная конференция для разработчиков высоконагруженных систем





#### Осебе

- Ph.D. в параллельных системах баз данных
- Core Developer в Postgres Professional

#### Postgres-специализация:

- WAL
- Planner
- B-tree/GiST/SP-GiST access method
- VACUUM





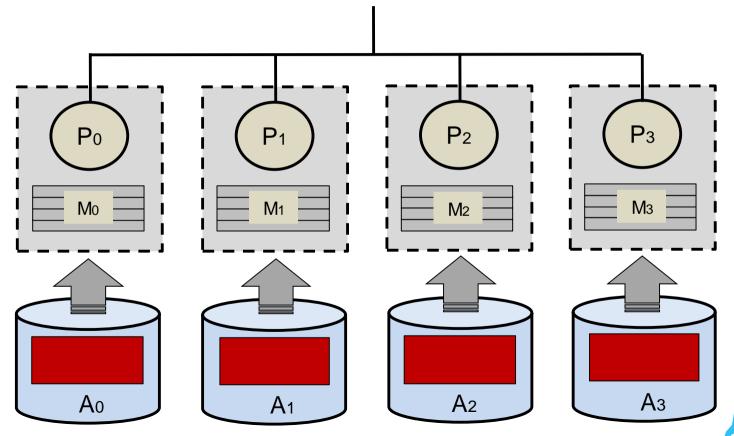
#### Кластер — для чего?

- Шардирование БД размером от 10 ТБ
- Доступность данных
- Масштабирование
  - Ускорение (*увеличить TPS*)
  - Расширяемость (*сохранить TPS на растущей БД*)



#### Утилизация кластера

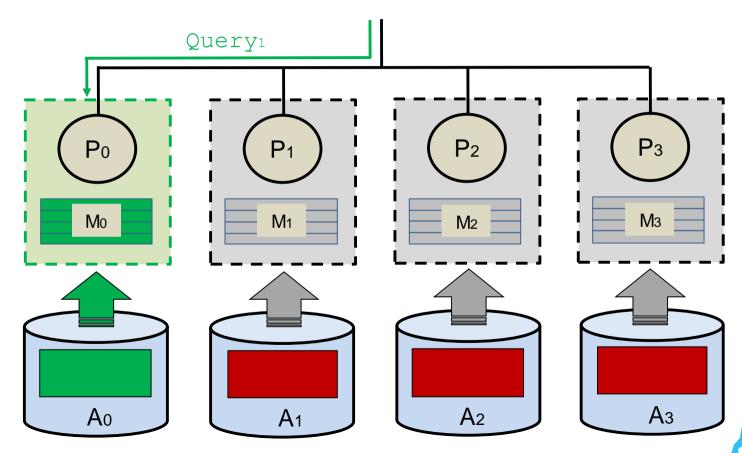






### Утилизация кластера: OLTP Postgres

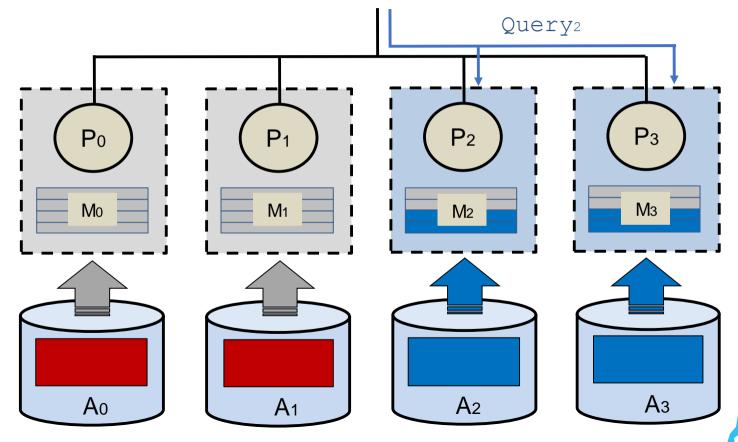






### Утилизация кластера: OLTP Postgres

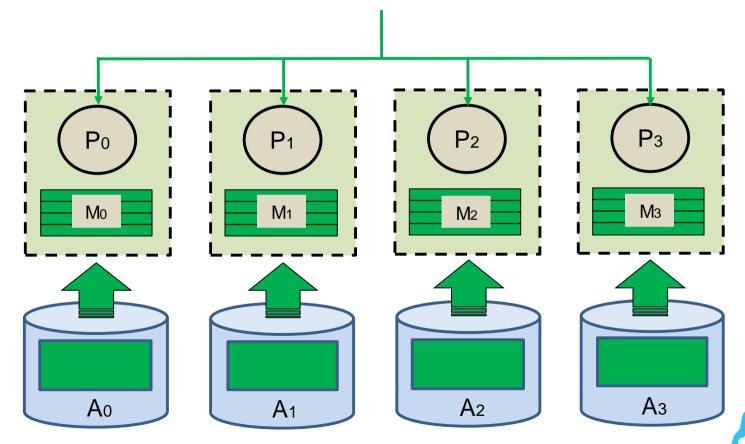






## Утилизация кластера: OLAP Postgres











- Планирование запроса
  - JOIN
  - Агрегаты
  - Сортировки
- Доступность данных
- Распределенные транзакции
- Балансировка загрузки







• Citus

https://www.citusdata.com/

Postgres-XL
 <a href="https://www.postgres-xl.org/">https://www.postgres-xl.org/</a>

• Greenplum

https://greenplum.org/

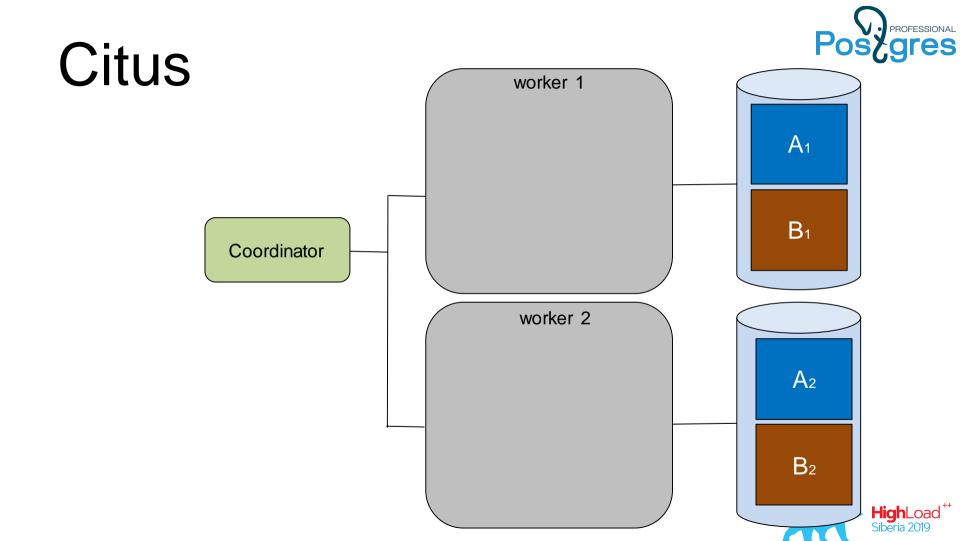


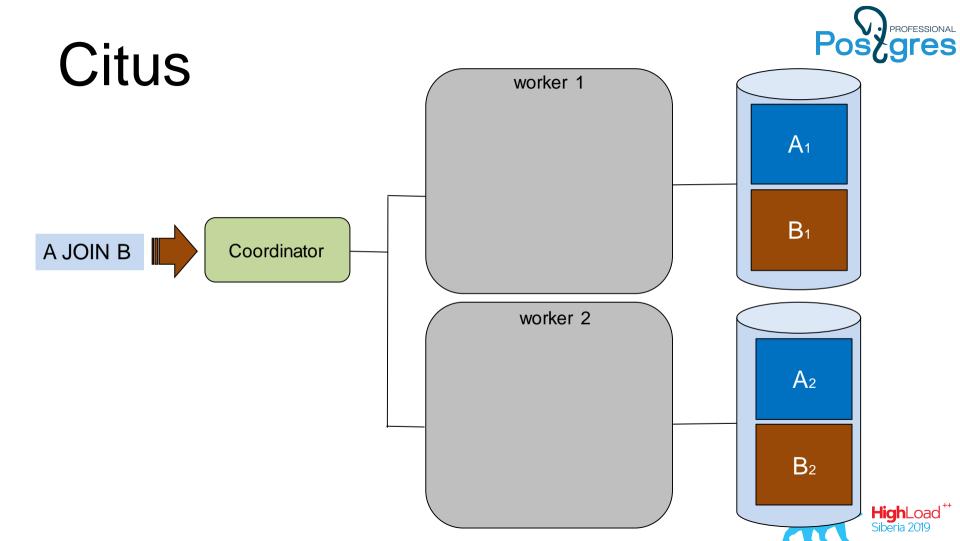
#### Citus

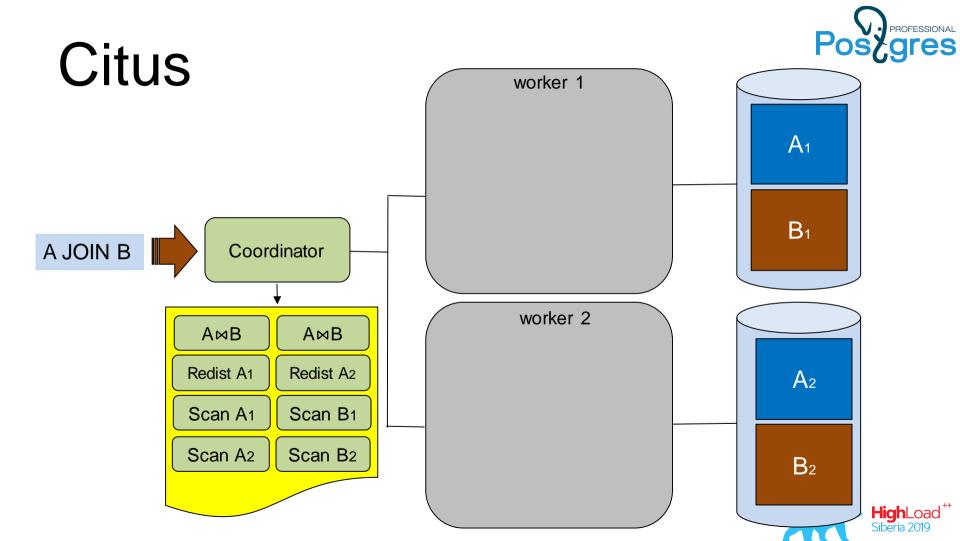


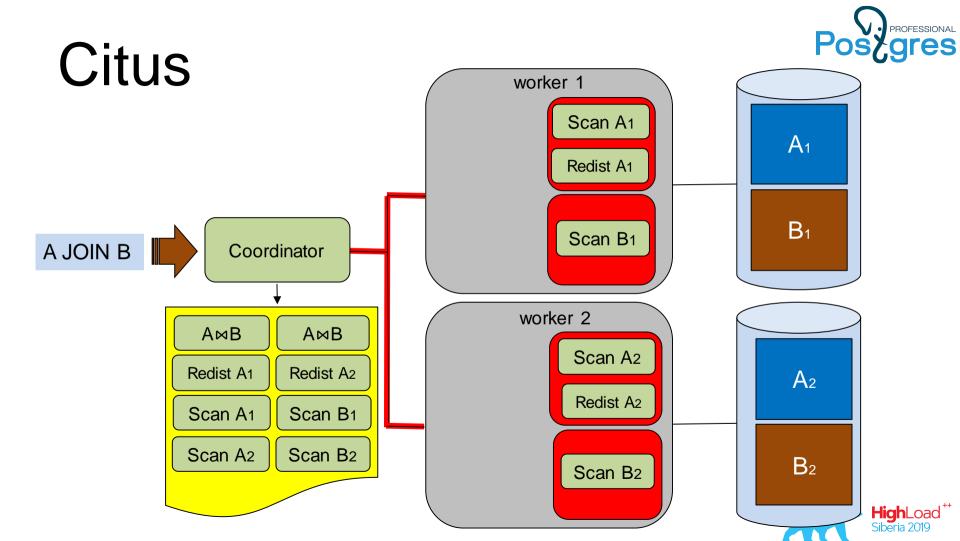
- Чистый EXTENSION
- Преобразует запрос в серии подзапросов, которые могут быть выполнены независимо каждым инстансом
- Динамически перераспределяет данные с материализацией между сериями подзапросов
- Собственный планнер
- Собственный механизм шардирования и прунинга
- Выделенный координатор

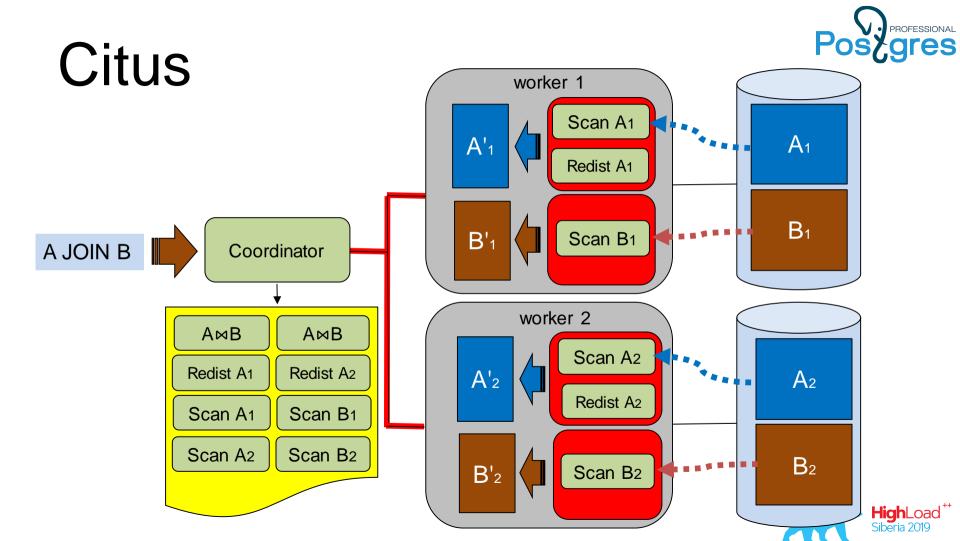


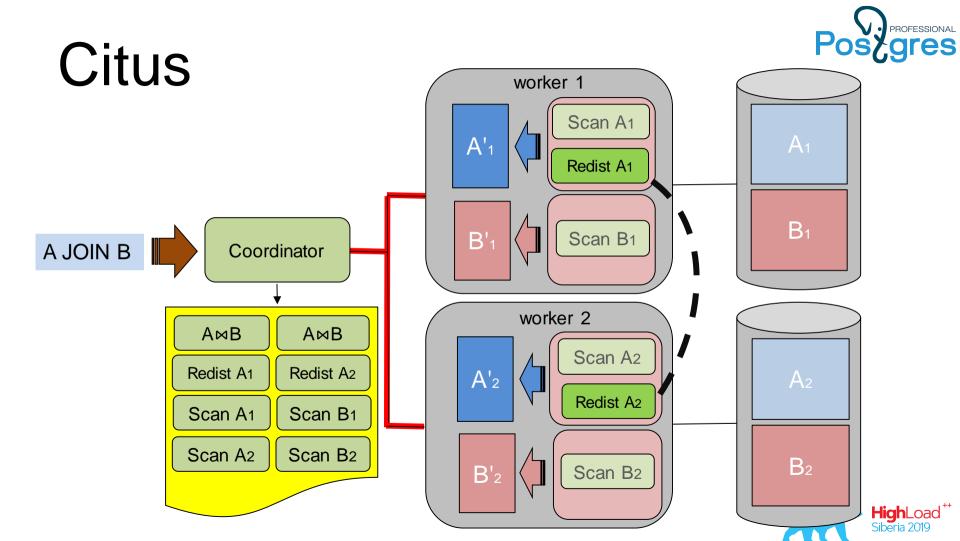


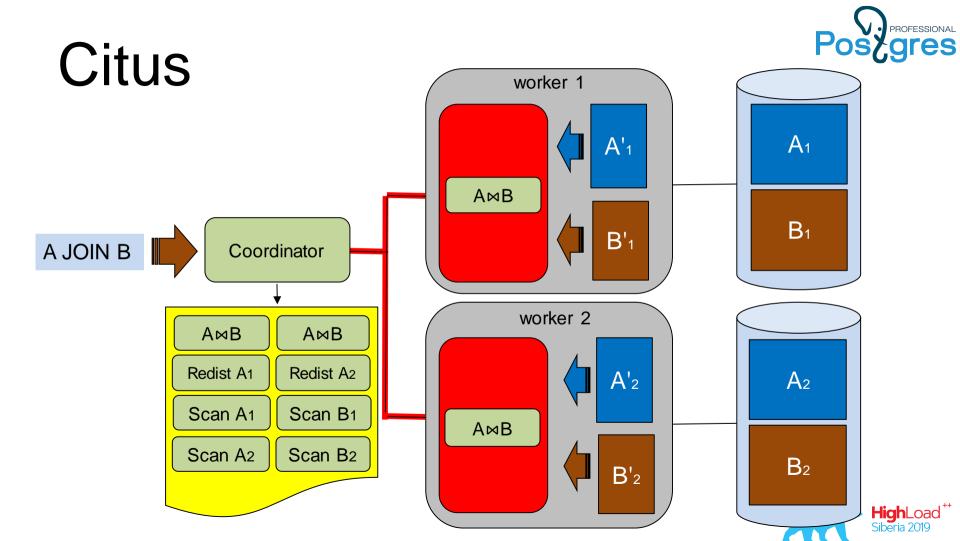












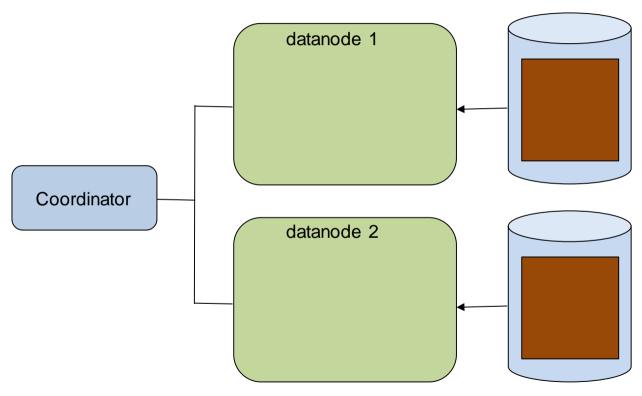
### Posegres

- Строит план один раз
- Разрезает план в местах, где требуется перераспределение данных
- Нет жесткого разделения на координатор и воркер
- Собственный планнер
- Собственный механизм шардирования и прунинга



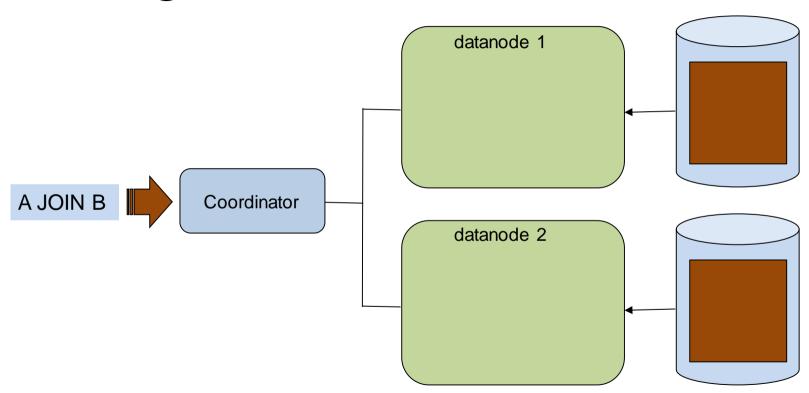






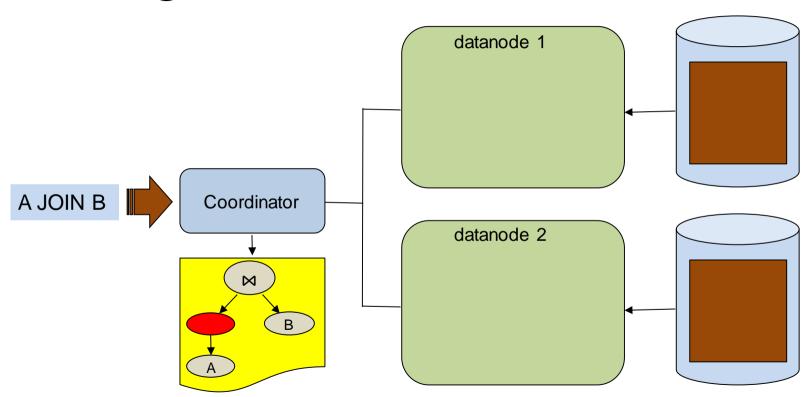






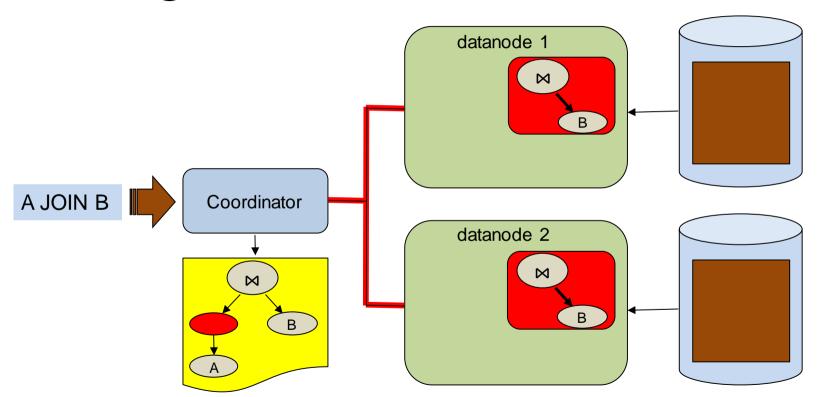






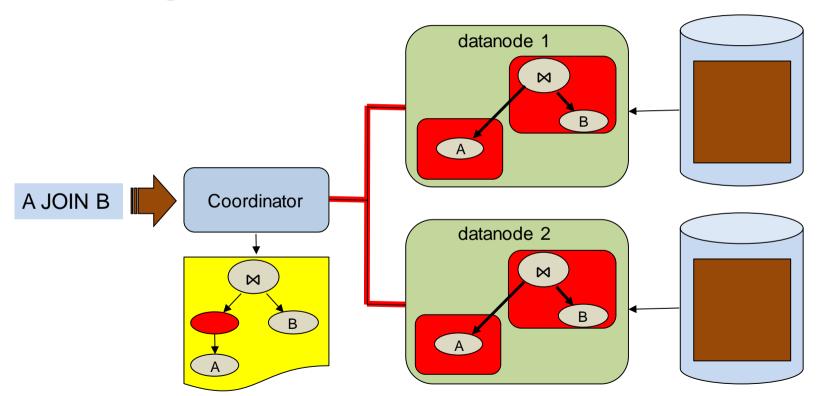


### Posegres Posegres



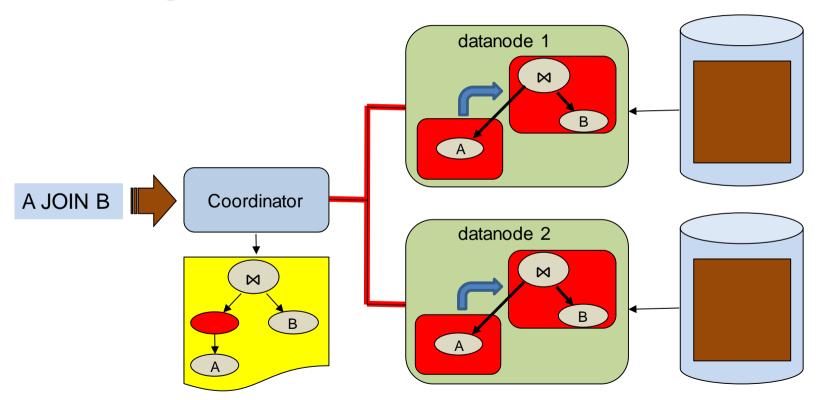


#### Post green



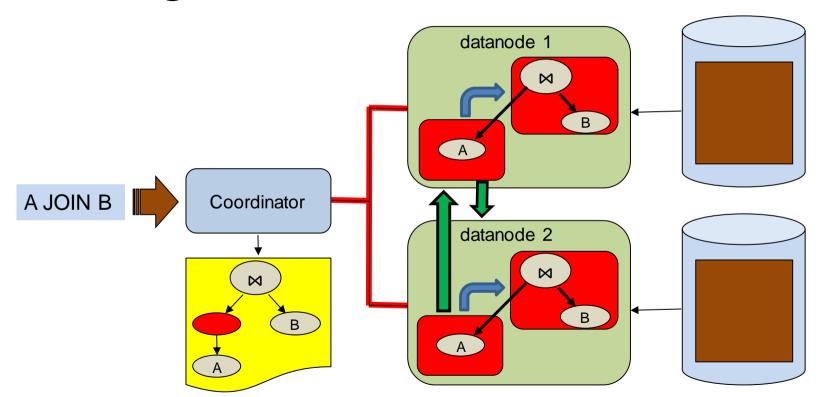








### Posegres Professional



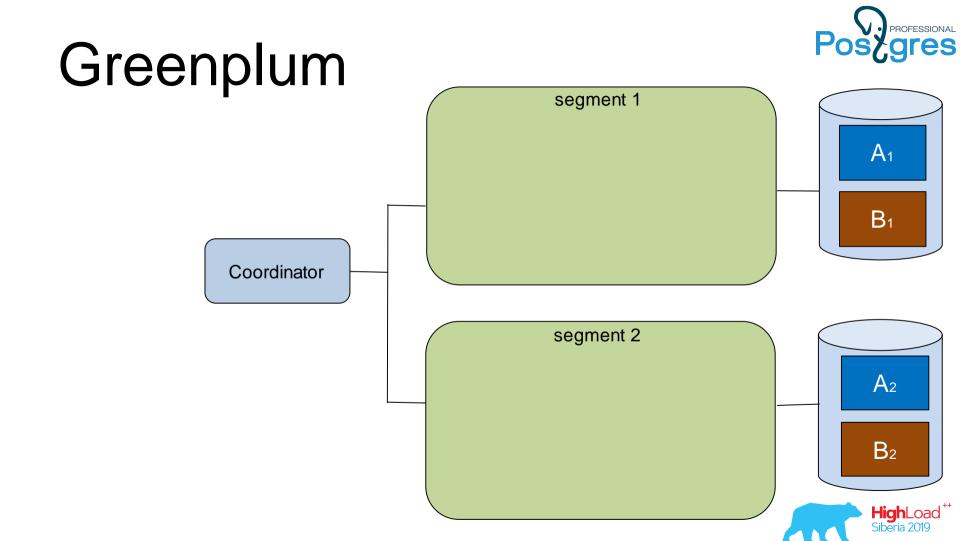


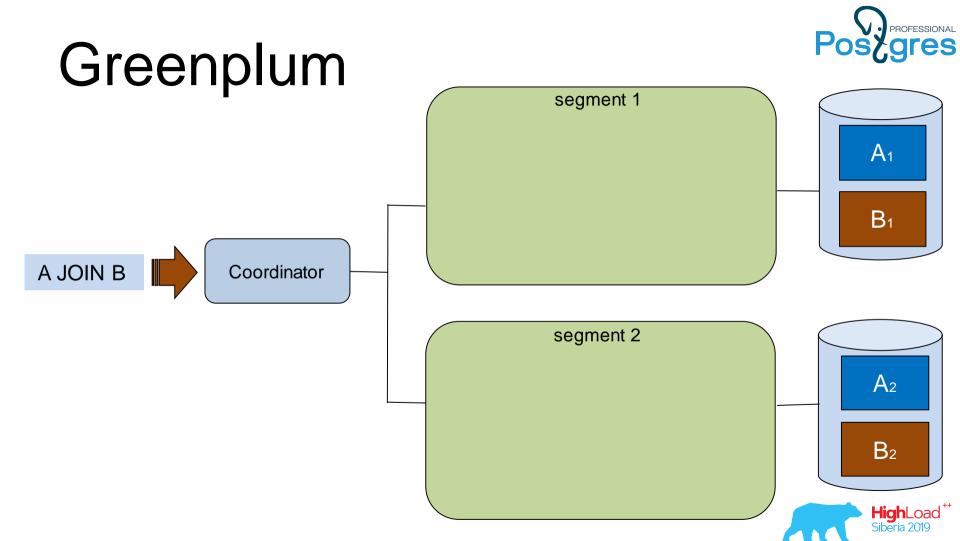


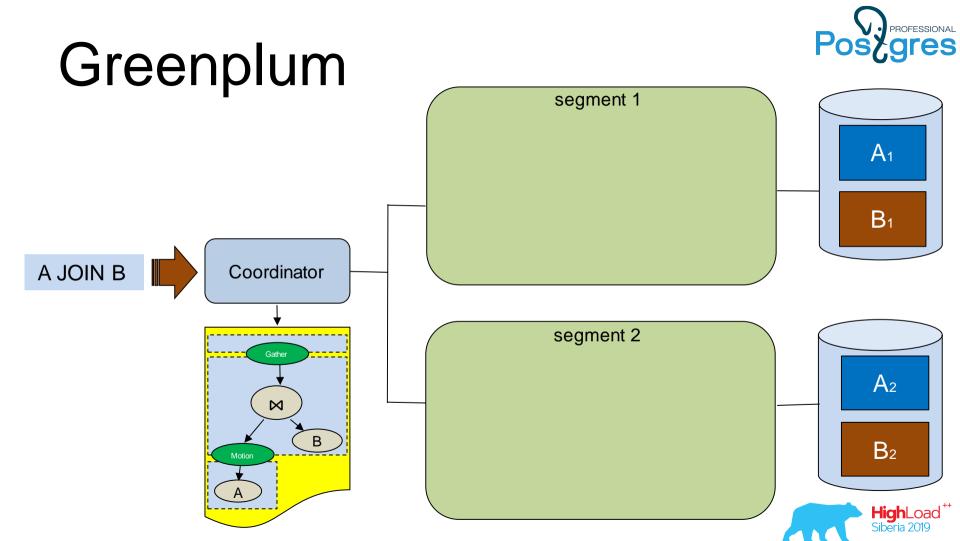
#### Greenplum

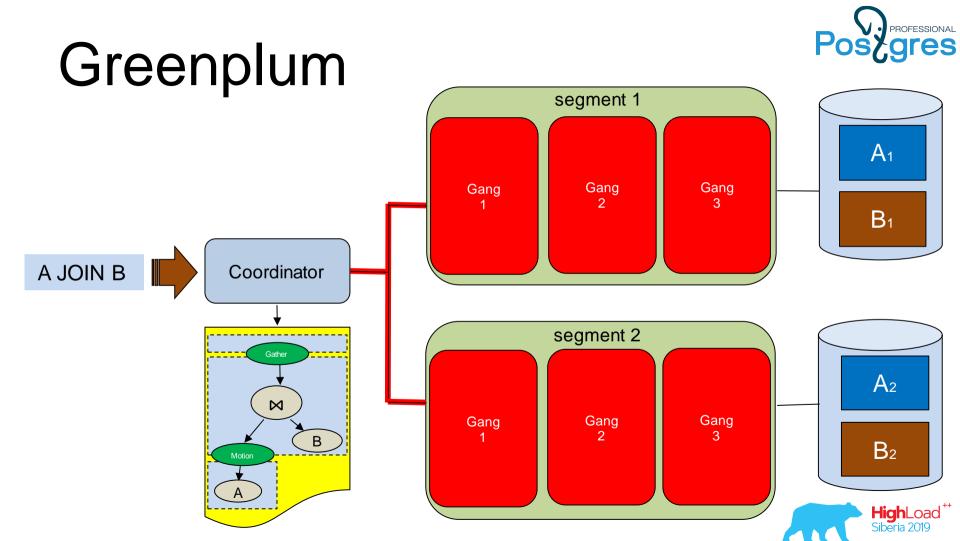
- Строит план один раз
- Делит план горизонтально (в целях перераспределения данных) и вертикально (подпланы, которые могут быть выполнены параллельно)
- Собственный планнер
- Собственный механизм шардирования и прунинга

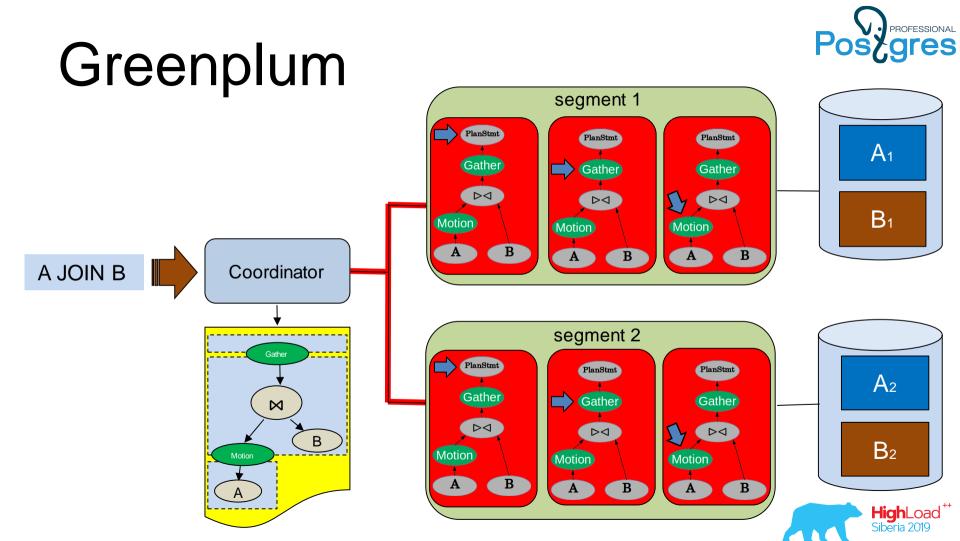


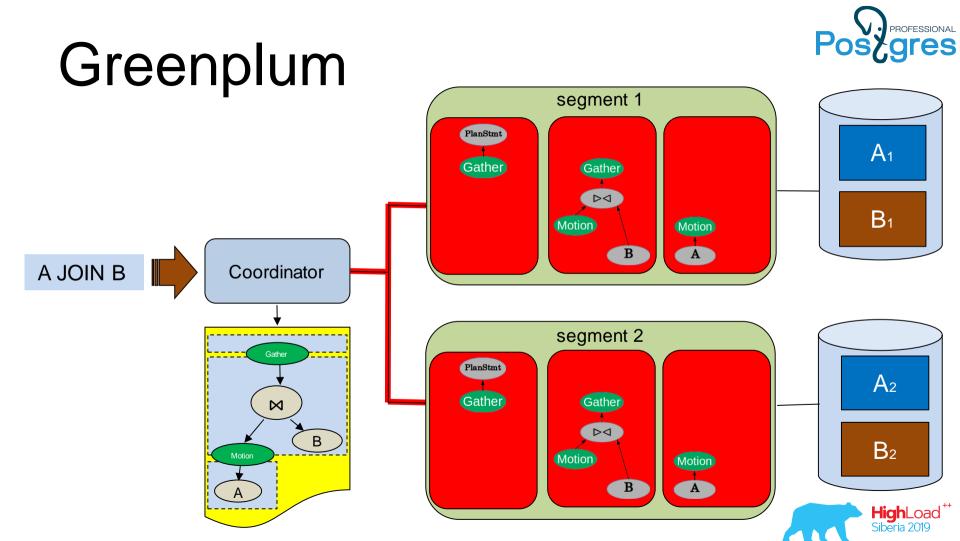












Greenplum segment 1 PlanStmt  $A_1$ Gather ---Gather  $B_1$ Motion A JOIN B Coordinator **4...** segment 2 PlanStmt  $A_2$ Gather <---Gather M В  $B_2$ Motion Motion **High**Load \*\* Siberia 2019





- Акцент на OLTP либо OLAP
- Высокая трудоёмкость апгрейда до новой версии ванильного постгреса
- Невысокая вероятность поддержки сообществом PostgreSQL в будущем





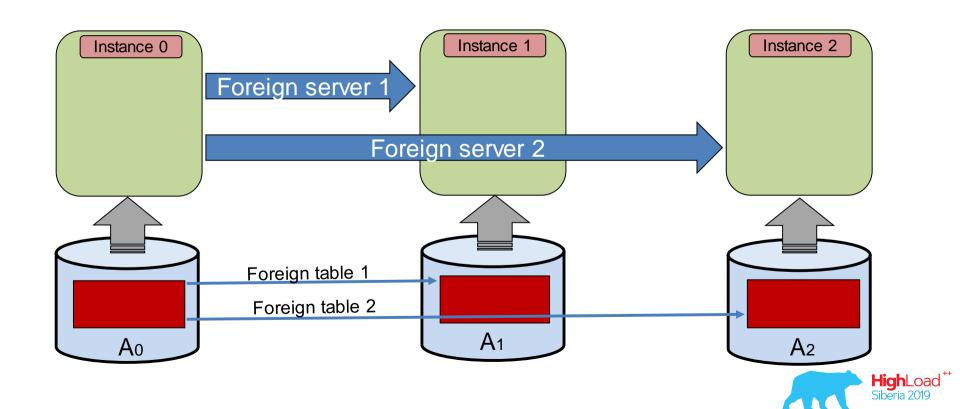
#### Шардман: принципы

- EXTENSION везде, где возможно
- Partitioning для шардирования и FDW для коннектов между инстансами
- Каждый инстанс может быть координатором
- Можно конфигурировать для OLTP и OLAP
- Ориентир на интеграцию в ванильный постгрес

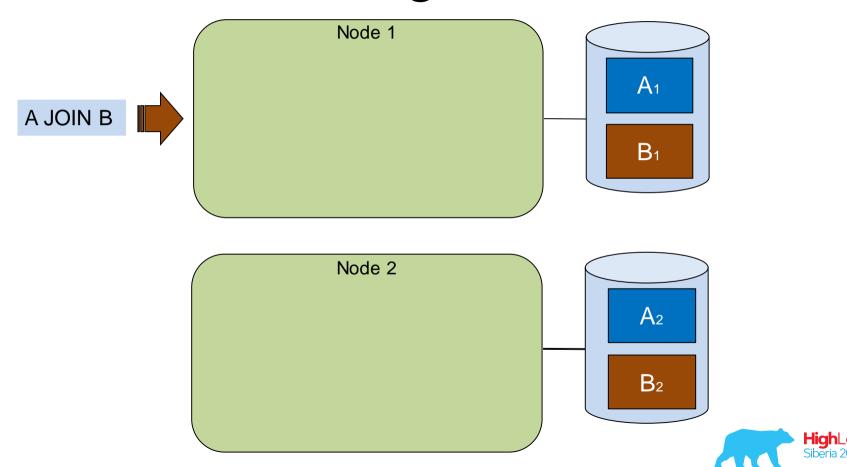


#### FDW + Partitioning

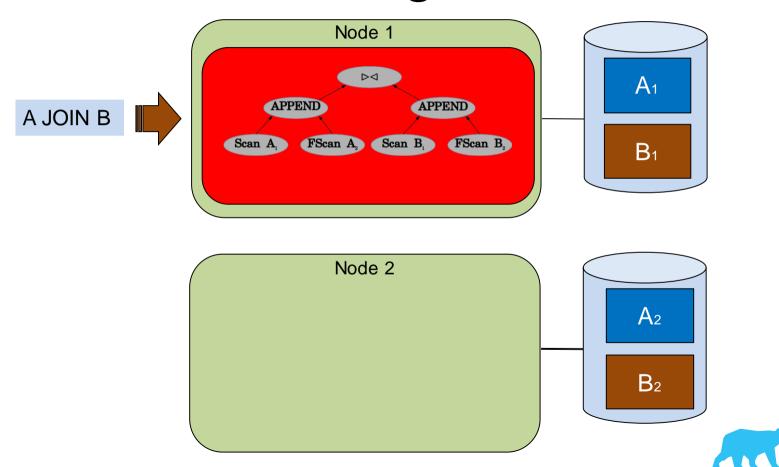






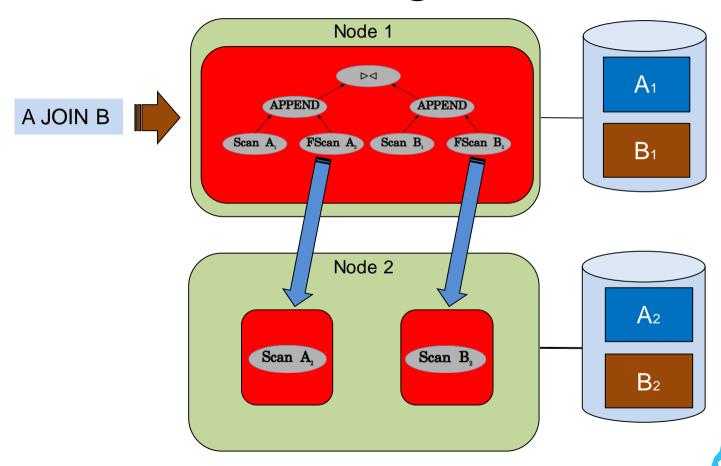






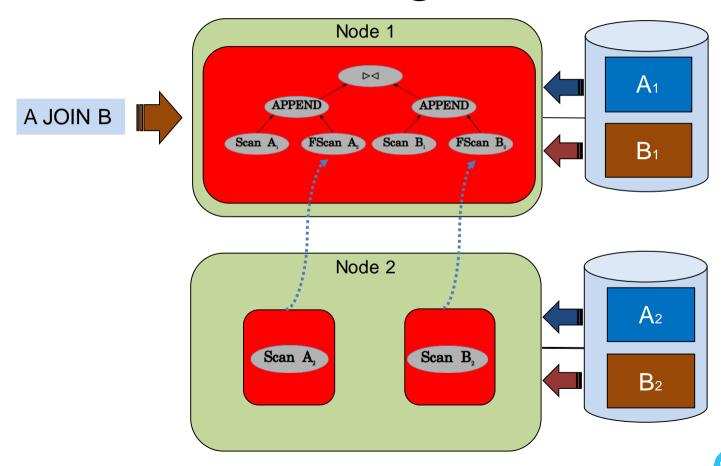


**High**Load \*\*\* Siberia 2019





**High**Load \*\*\* Siberia 2019

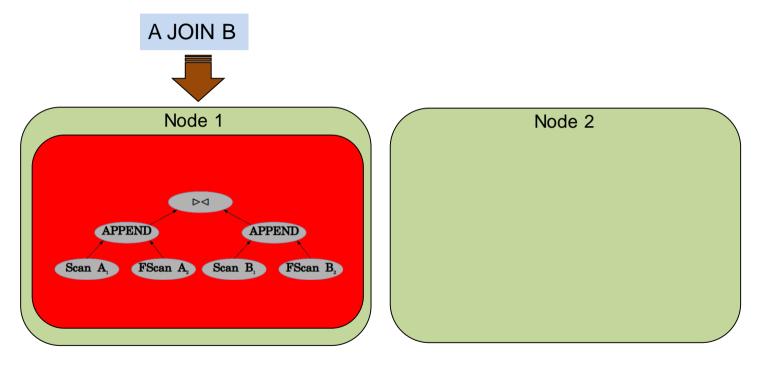


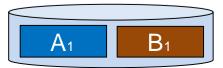


- Нет механизмов динамического решардирования
- Передает не план, а запрос
- Последовательный APPEND





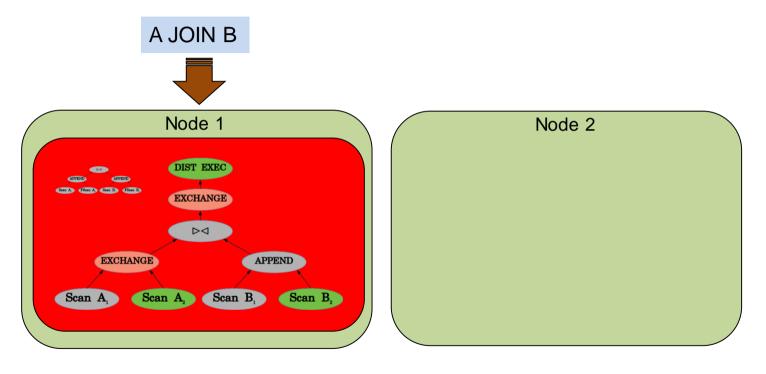


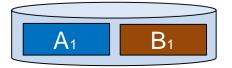


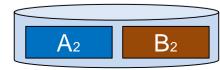






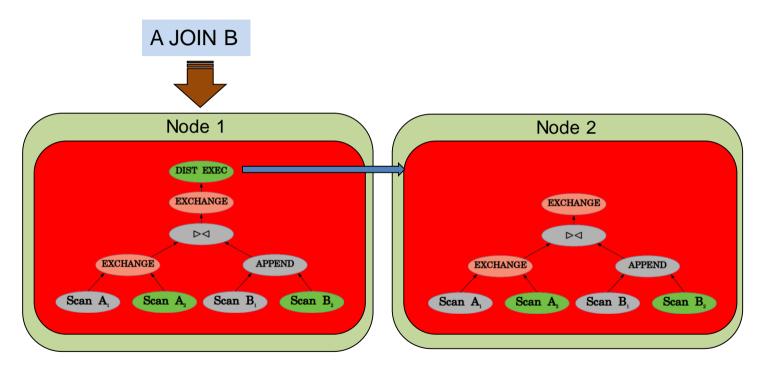


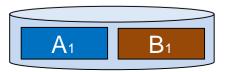








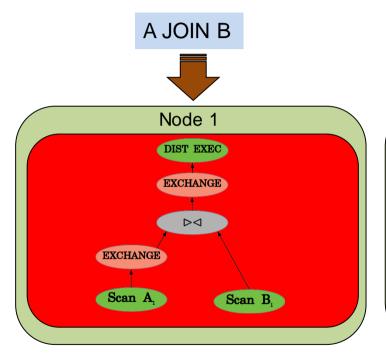


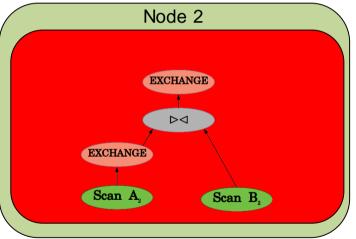


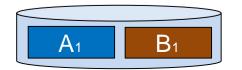








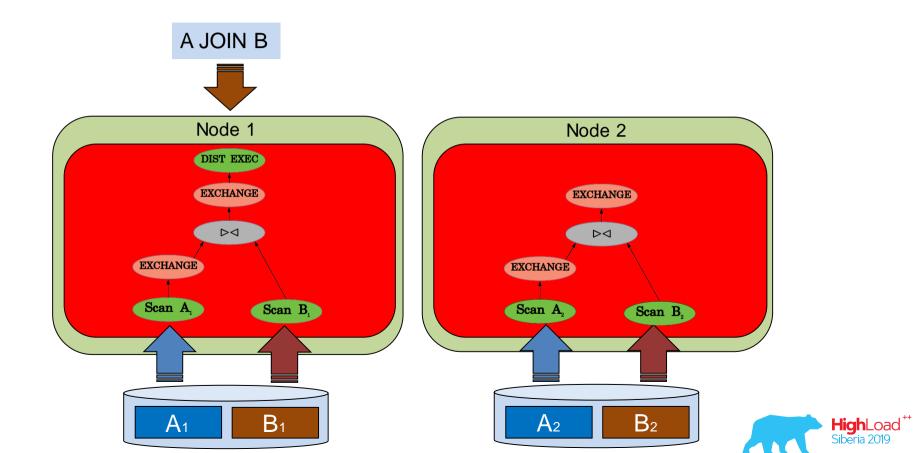




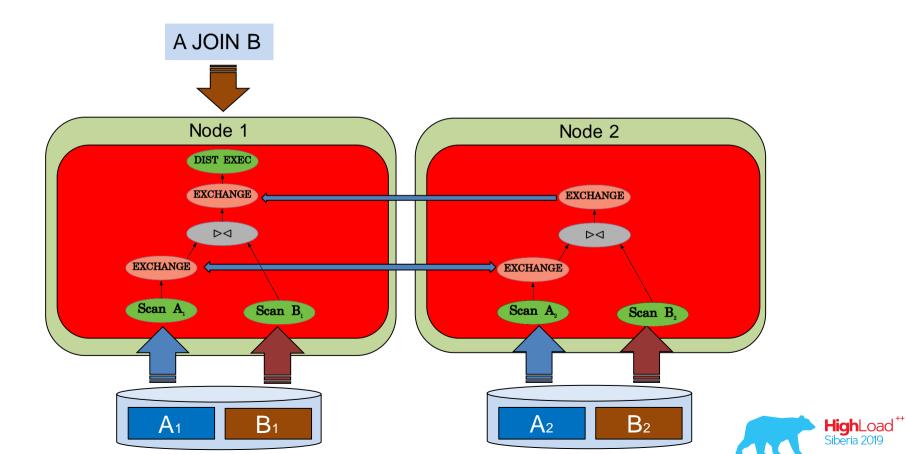






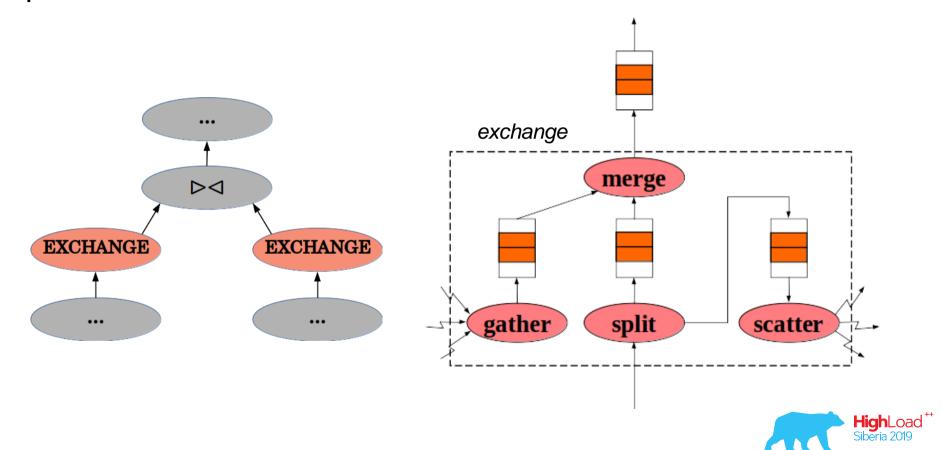






## Инкапсуляция параллелизма Post gres

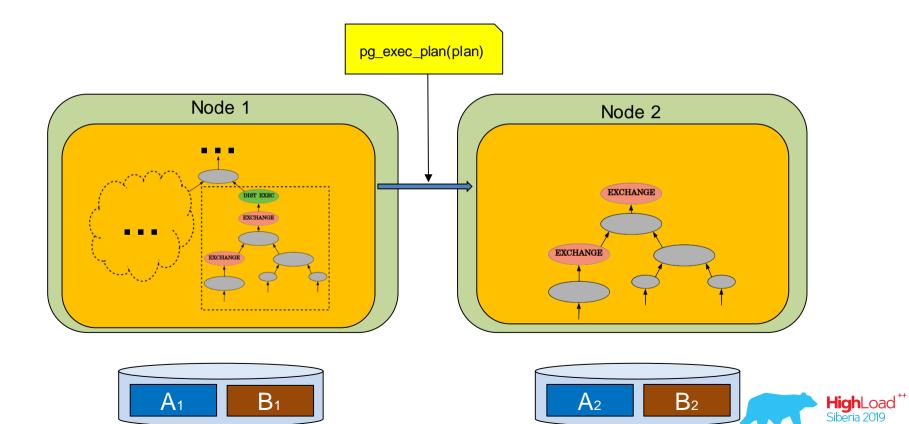
параллельное выполнение JOIN



#### Рассылка планов

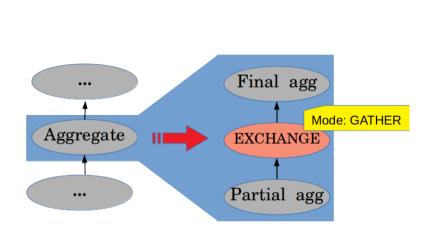


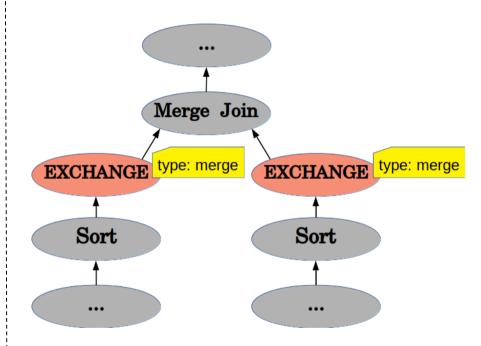
Dist Exec custom-node



#### To Do: агрегаты и сортировки









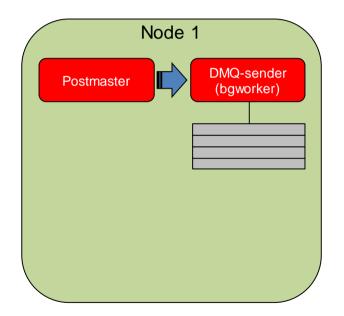
## Message passing

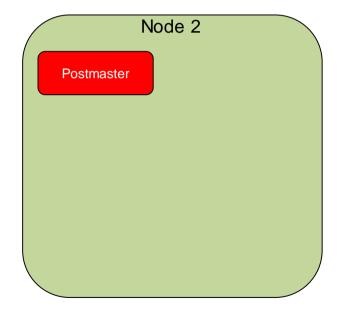


- Citus PostgreSQL receiver
- Postgres-XL PostgreSQL receiver
- Greenplum socket-based subsystem
- Shardman DMQ



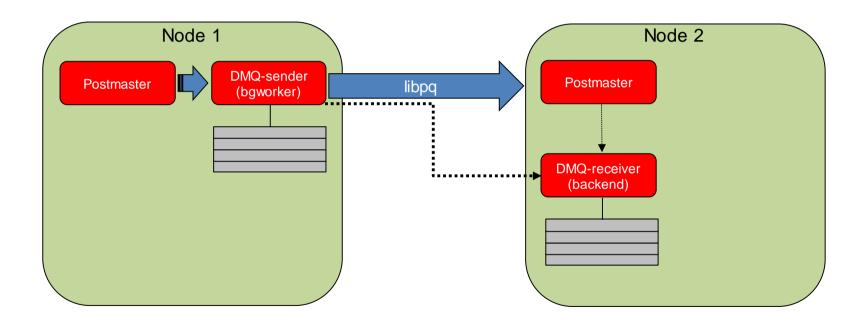






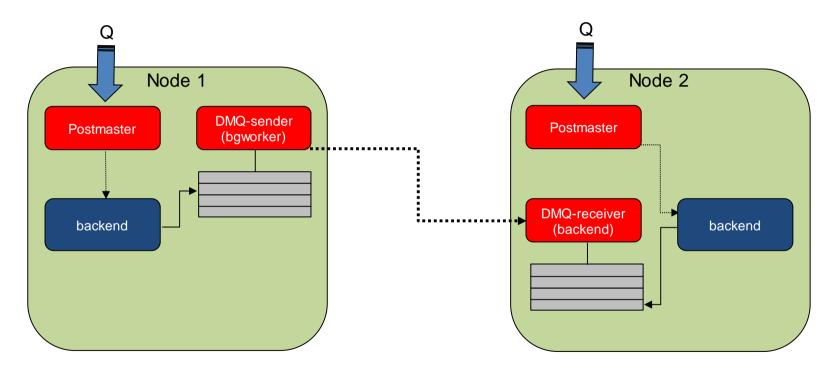






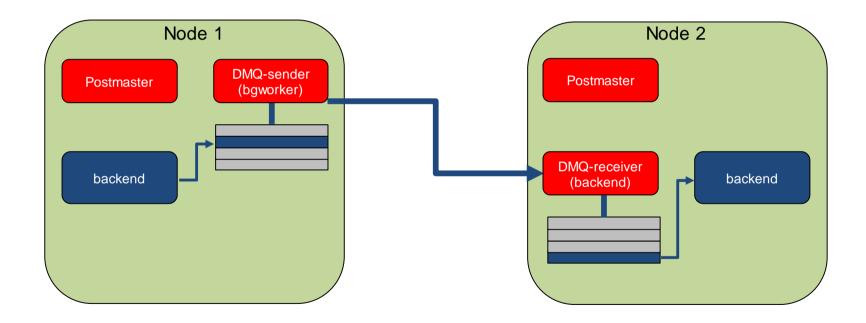






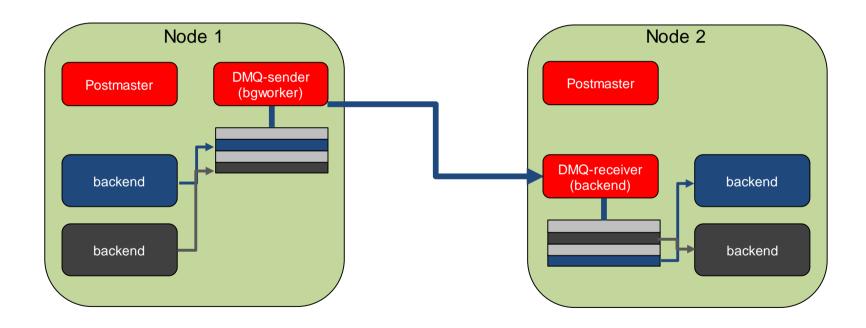
















#### Что получилось

- Bce PostgreSQL-инстансы равнозначны
- Только libpq соединения между инстансами
- Весь параллелизм инкапсулирован в custom-нодах планнера
- Параллелизм задействуется только для той части плана, где это необходимо
- Шардинг и прунинг выполняются средствами partitioning a и FDW





#### Шардман: OLTP

- Координатором может выступать любой узел
- План запроса создается однократно
- При использовании только локальных данных инстанса выполнение запроса не отличается от ванильного PostgreSQL
- Задействуется минимально необходимое количество инстансов, на которых выполняется минимально необходимый подплан



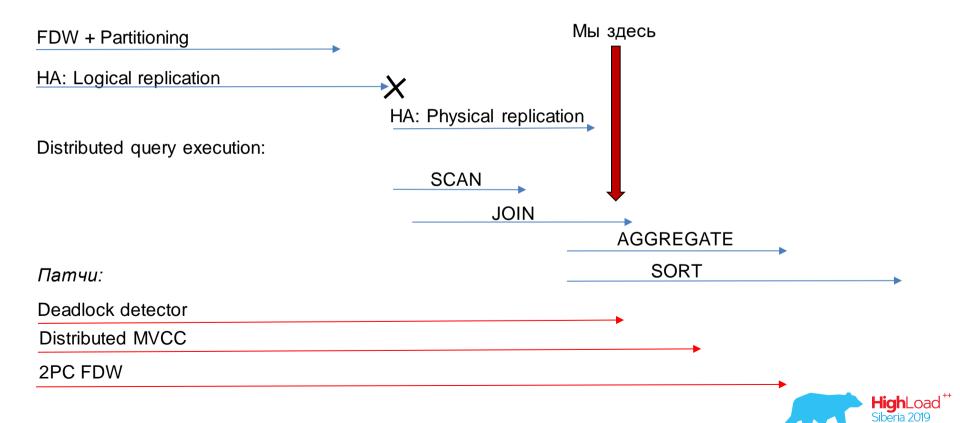
#### Шардман: OLAP

- Возможность параллельной загрузки данных
- SCAN, JOIN задействуют все инстансы (в перспективе агрегаты и сортировки)
- Доступны все планнеры постгреса



## Roadmap





#### Ссылки



- 1. https://github.com/postgrespro/shardman/
- 2. https://github.com/postgrespro/postgres\_cluster



# **High**Load Siberia 2019

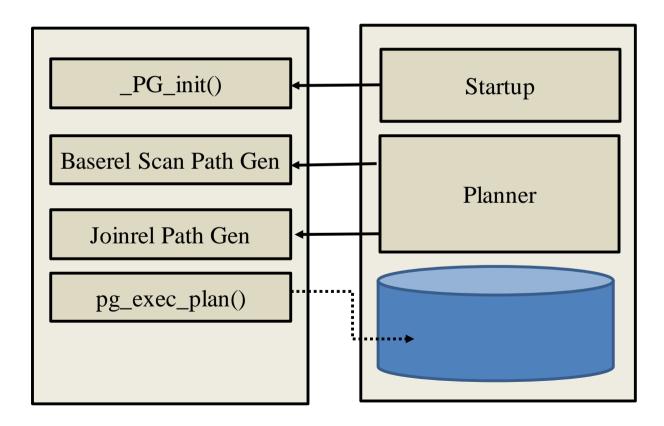
Профессиональная конференция для разработчиков высоконагруженных систем

# Вопросы?





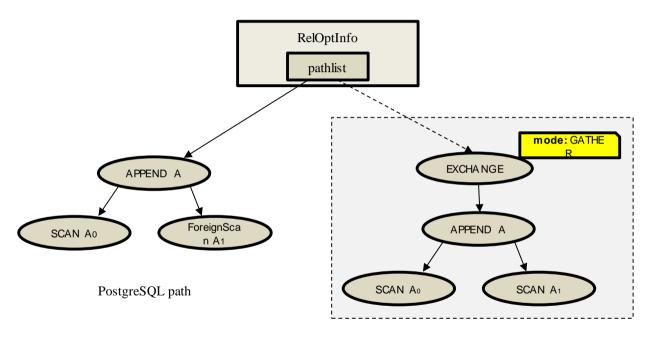
## Архитектура





## Baserel path generator



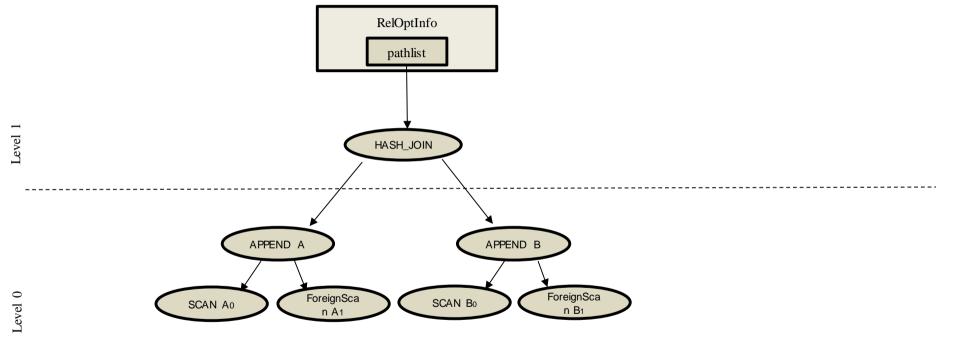


Shardman proposal



## JOIN Path generator

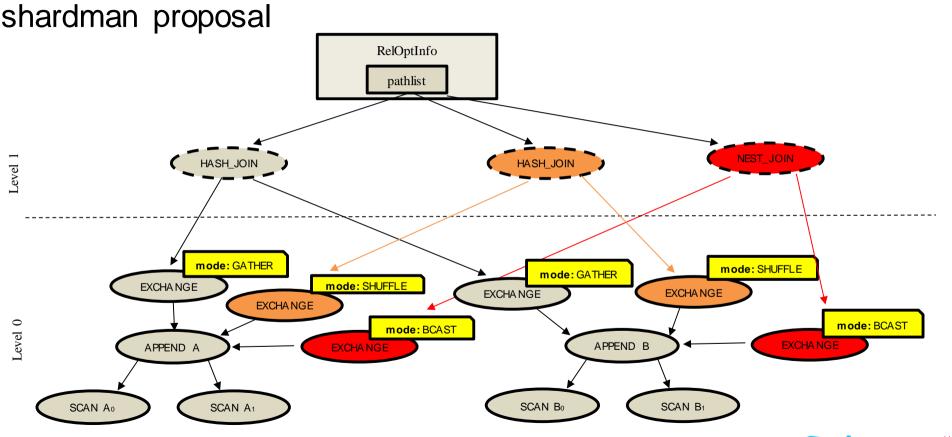






# JOIN Path generator







## Postgres-XL



SELECT \* FROM A,B WHERE a.payload=b.id;

