

Greenplum для разработчиков и архитекторов баз данных

Оптимизация сложных запросов в MPP-кластерах: Greenplum, Arenadata DB, Cloudberry Database



**Меня хорошо видно
& слышно?**



Защита проекта

Тема: Оптимизация сложных запросов в MPP-кластерах: Greenplum, Arenadata DB, Cloudberry Database



Игорь Щербаков

Разработчик баз данных

План защиты

Цель и задачи проекта

Какие технологии использовались

Что получилось

Выводы

Вопросы и рекомендации

Цель и задачи проекта

Цель проекта:

Разбор методов оптимизации сложных запросов в MPP-кластерах

1. Собрать методы оптимизации запросов из разных источников
2. Создать MPP-кластеры: Greenplum, Arenadata DB, Cloudberry Database
3. Проанализировать планы выполнения различных запросов в MPP-кластерах
4. Выявить преимущества и недостатки MPP-кластеров



Какие технологии использовались

1. MPP - Massively Parallel Processing
2. Virtualization, Docker
3. Greenplum, Arenadata DB, Cloudberry Database
4. PXF - Platform Extension Framework
5. Dbeaver, PostgreSQL



Использованные MPP-кластеры

Кластер	Версия кластера	Версия PostgreSQL	Время на первую установку
Greenplum	6.23, 6.27	9.4	1 месяц
Arenadata DB	7.2	12.12	1 день
Cloudberry Database	1.5.1, 1.6	14.4	1 час

Варианты оптимизаторов (планировщиков): GPORCA, Postgres.



База данных «Авиаперевозки»

1. Курс “QPT. Оптимизация запросов” фирмы “PostgresPro”
(редакция 27.12.2024):

<https://postgrespro.ru/education/demodb>

2. Книга: Домбровская Г., Новиков Б., Бейликова А.
Оптимизация запросов в PostgreSQL. - М.: ДМК Пресс, 2022.
(PostgreSQL Query Optimization. The Ultimate Guide to Building
Efficient Queries. - Apress, 2021.)

https://drive.google.com/drive/folders/13F7M80Kf_somnjb-mTYAnh1hW1Y_g4kJ?usp=sharing

8 таблиц, ~22 млн. строк, ~1.5 Gb. Выбрано 33 запроса.



ER-диаграмма «Авиаперевозки»

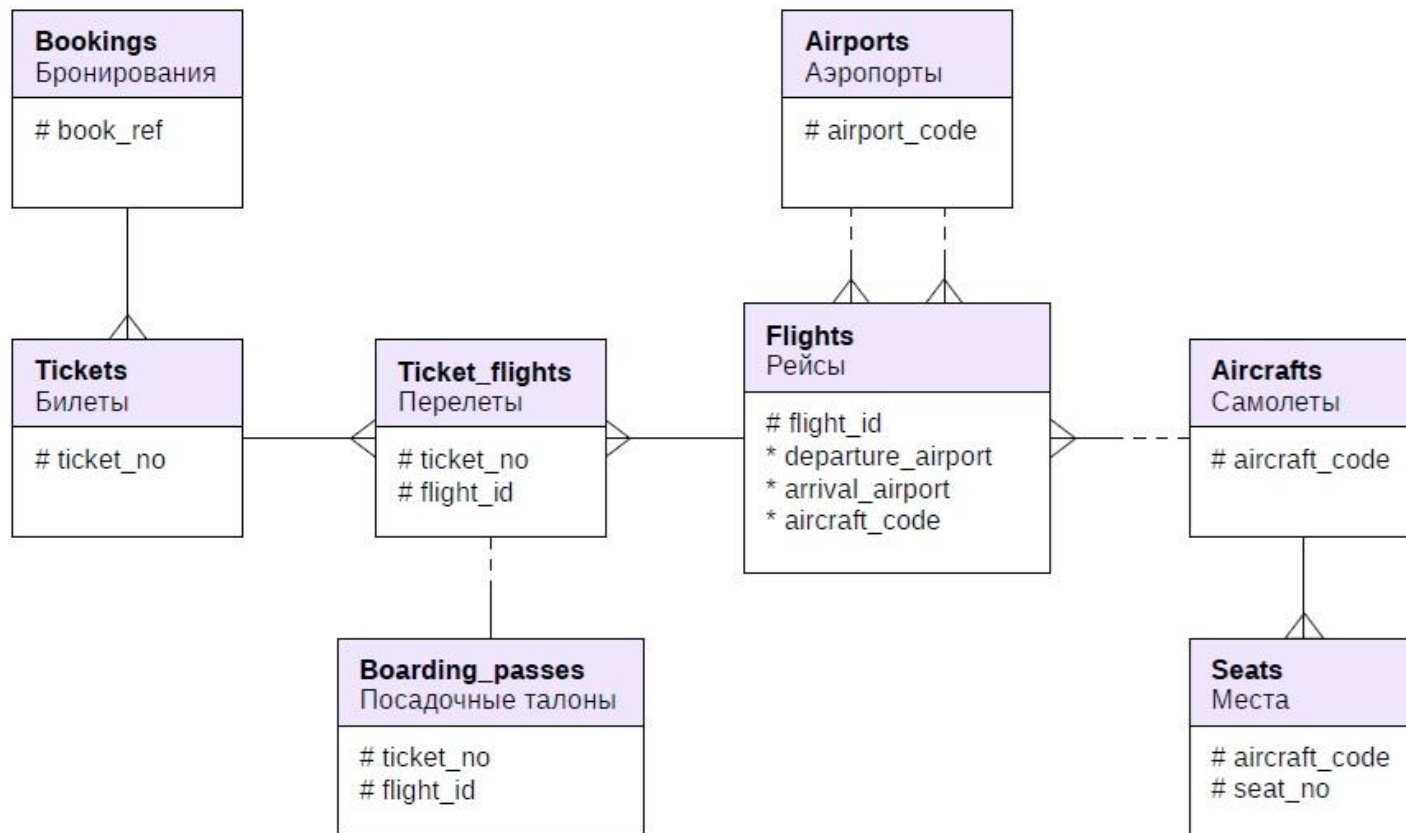
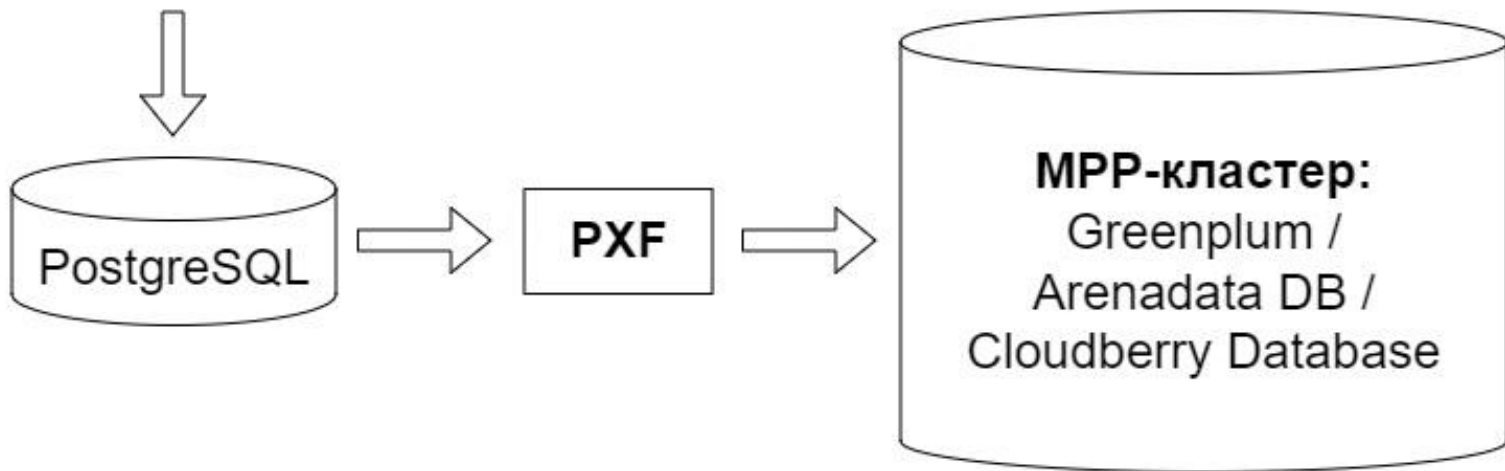


Схема кластеров

БД "Авиаперевозки"



Источники по методам оптимизации

1. Лекция «Оптимизация запросов» и другие лекции
2. Документация по Greenplum: <https://techdocs.broadcom.com/>
3. Документация по Arenadata: <https://docs.arenadata.io>
4. 5 лайфхаков оптимизации SQL-запросов в Greenplum
<https://habr.com/ru/companies/rostelecom/articles/442758/>
5. Другие материалы из интернета



Методы оптимизации запросов MPP (1)

1. Распределение данных, соединения по ключам дистрибуции. Для ключей дистрибуции, по которым будут соединения, использовать одинаковые типы данных.
2. Партиционирование, partition elimination
3. Использование append optimized таблиц, использование колоночной ориентации
4. Использование distributed replicated для маленьких таблиц (справочников)
5. Использование unlogged таблиц
6. Использование temporary таблиц для хранения промежуточных результатов вычислений

Методы оптимизации запросов MPP (2)

7. Добиваться равномерного распределения данных между сегментами. Избегать skew – перекосов.
8. По возможности, избегать переноса недостающих данных с одного сегмента на другой: broadcast motion и redistribute motion.
9. По возможности, избегать операций сортировки (order by в запросе, sort в плане запроса)
10. Поддержание статистики в актуальном состоянии, регулярный сбор
11. Использование индексов
12. Управление оптимизацией с помощью параметров

Что получилось

1. Созданы кластеры: Greenplum, Arenadata DB, Cloudberry Database
2. Базы данных заполнены специальными данными для воспроизведения сложных случаев оптимизации запросов
3. Из разных источников собраны методы оптимизации запросов
4. Получены планы выполнения большого числа сложных и интересных запросов
5. Сделаны выводы



Выводы (1)

1. MPP-кластеры отстают от БД PostgreSQL в оптимизации запросов
2. Чем больше версия MPP-кластера и БД, входящих в его состав, тем выше производительность запросов
3. Все запросы, которые можно выполнить с помощью Sec Scan + Hash Join, MPP-кластеры выполняют с помощью этих операций и обеспечивают хорошее время выполнения за счет распараллеливания по сегментам
4. Методы оптимизации, собранные при выполнении данной проектной работы, позволят ускорять запросы в предстоящих проектах
5. Полученные навыки создания MPP-кластеров позволят эффективно развертывать стенды для разработки и тестирования



Выводы (2)

6. Cloudberry Database 1.6 эффективнее других MPP-кластеров по некоторым видам запросов: поиск по диапазону, агрегирование, группировки
7. В Arenadata DB 7.2 действительно реализовано эффективное сканирование индексов
8. Производительность Arenadata DB 6.27 и 7.2 почти не различается
9. Переключение с оптимизатора GPORCA на оптимизатор PostgreSQL в большинстве случаев не дает выигрыша в производительности
10. ADCM – удобное средство для установки и администрирования Arenadata DB



Вопросы и рекомендации



если есть вопросы



если вопросов нет

Спасибо за внимание!

