

GreenPlum

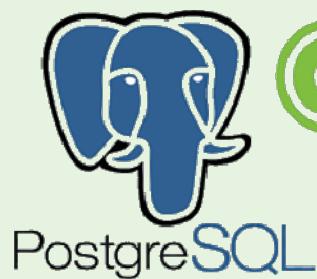
Просто о сложном!



Что это такое?

— □ ×

GreenPlum – это высоконагруженная, транзакционная, аналитическая система построенная на MPP архитектуре на основе PostgreSQL.



GREENPLUM
DATABASE®



A C I D



01

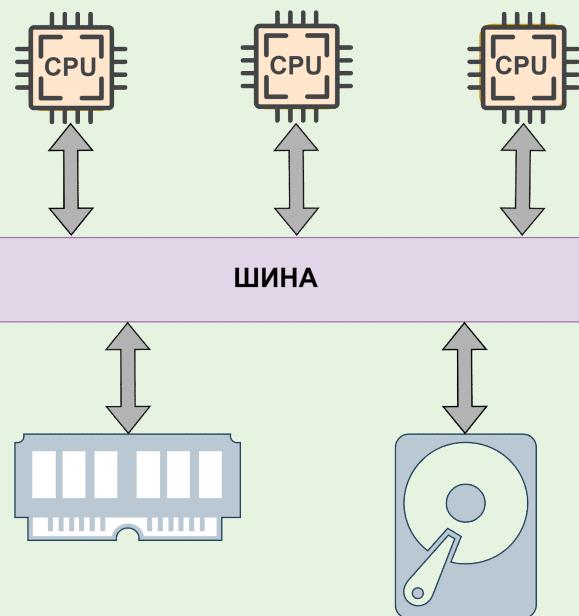
Архитектура GreenPlum

- SMP и MPP архитектуры
- Архитектура GreenPlum
- Операции SQL

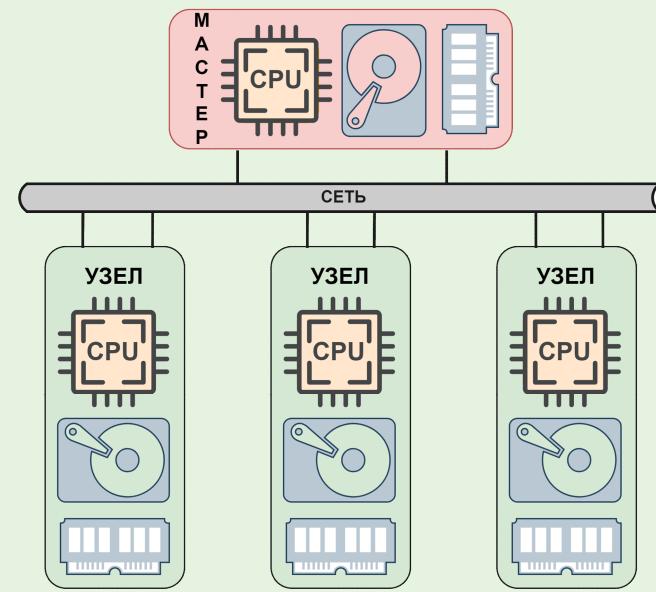


SMP vs MPP

SMP

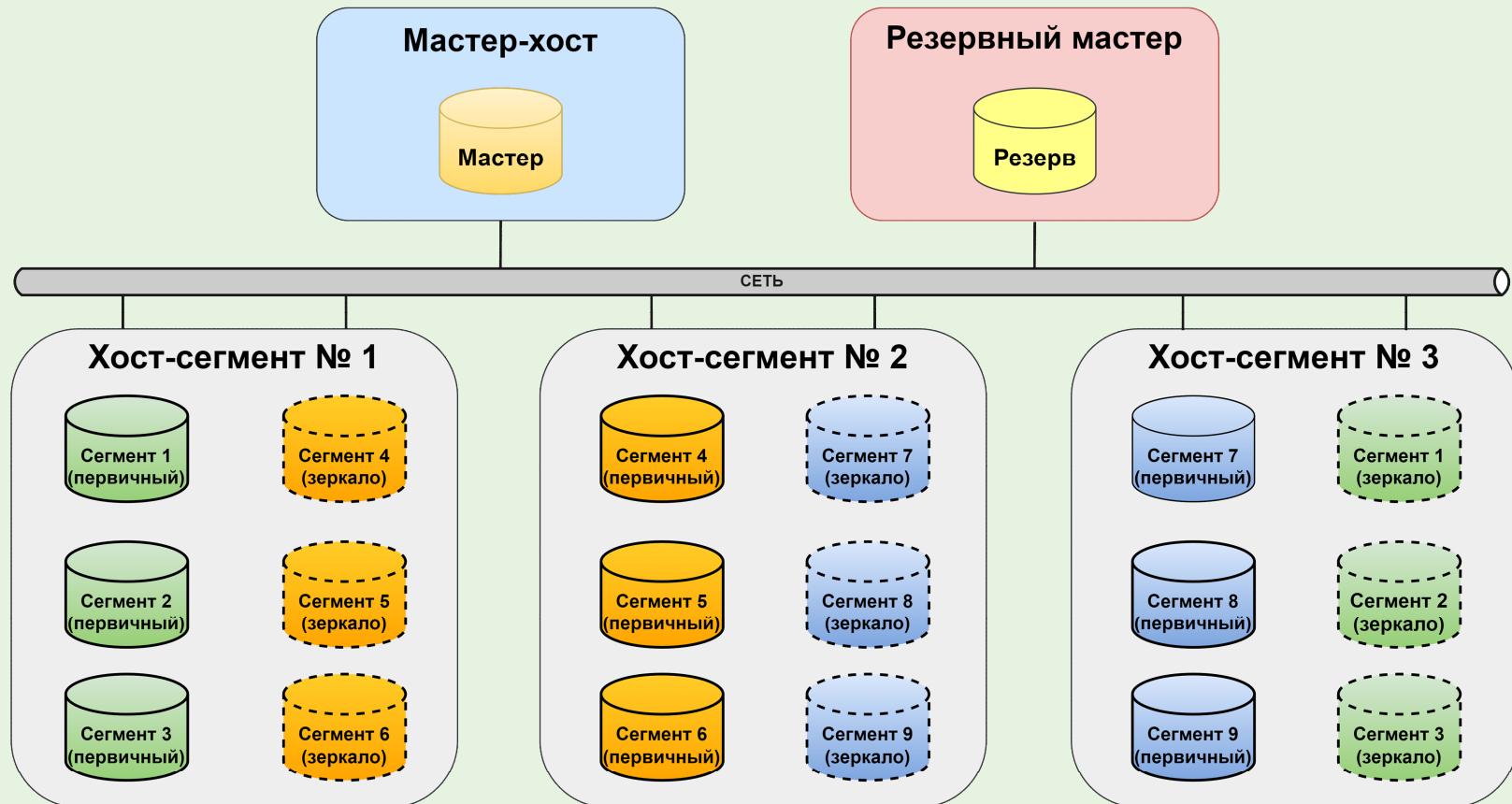


MPP



Архитектура кластера GreenPlum

— □ ×



Архитектура кластера GreenPlum

— □ ×

Мастер-хост

Резервный мастер

Операции SQL

— □ ×

SQL

DDL

CREATE
DROP
ALTER
TRUNCATE

DML

INSERT
UPDATE
DELETE

DCL

GRANT
REVOKE

TCL

COMIT
ROLLBACK
SAVE POINT

DQL

SELECT

X

Сервер 1 (первичный)

Сервер 2 (реплика)

Сервер 3 (реплика)

Х

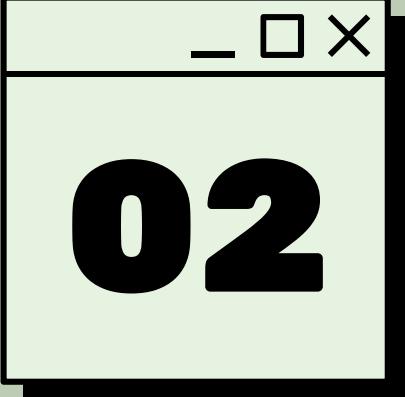
Сервер № 3

Сегмент 1 (зеркало)

Сегмент 2 (зеркало)

Сегмент 3 (зеркало)





02

.....

>>>



Реплицирование(зеркалирование), сегментов

- Репликация мастер-сегмента
- Журнал предзаписи (WAL-журнал)
- Репликация первичного сегмента

Архитектура кластера GreenPlum

— □ ×

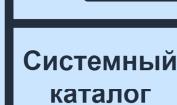
Мастер-хост

Резервный мастер

Зеркалирование мастер-сегмента GreenPlum

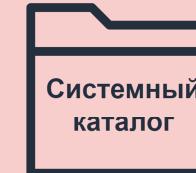
— □ ×

Мастер-хост



WAL-журнал

Резервный мастер



WAL-журнал



Хост № 1

Сегмент 1 (первоначальный)

Сегмент 2 (первоначальный)

Сегмент 3 (первоначальный)



Хост № 3

Сегмент 1 (зеркало)

Сегмент 2 (зеркало)

Сегмент 3 (зеркало)

поток SQL-команд

Сегменты 1, 2 и 3 зеркально отражены на хосте № 3.

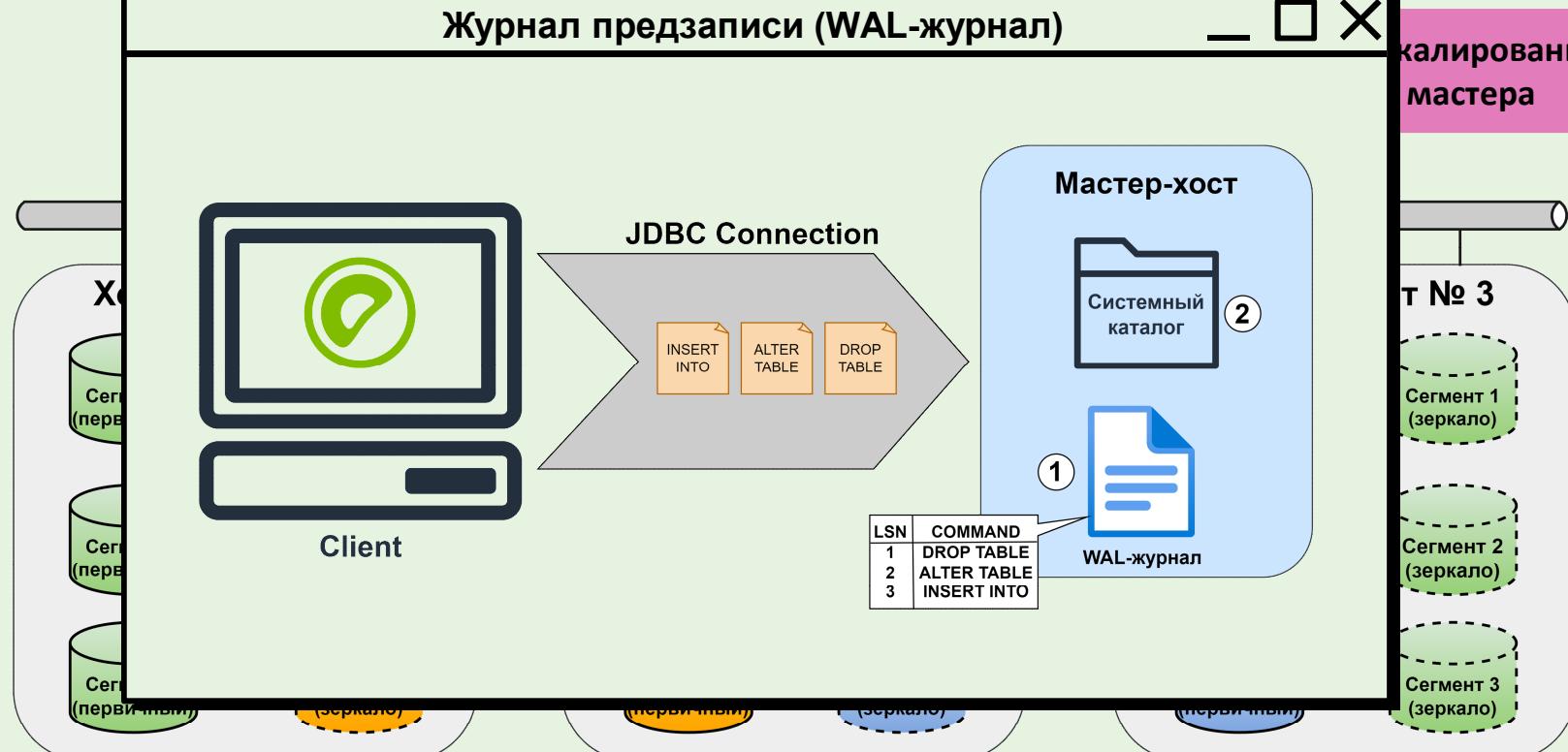
Архитектура кластера GreenPlum

Мастер-хост

Резервный мастер

Журнал предзаписи (WAL-журнал)

Масштабирование
мастера

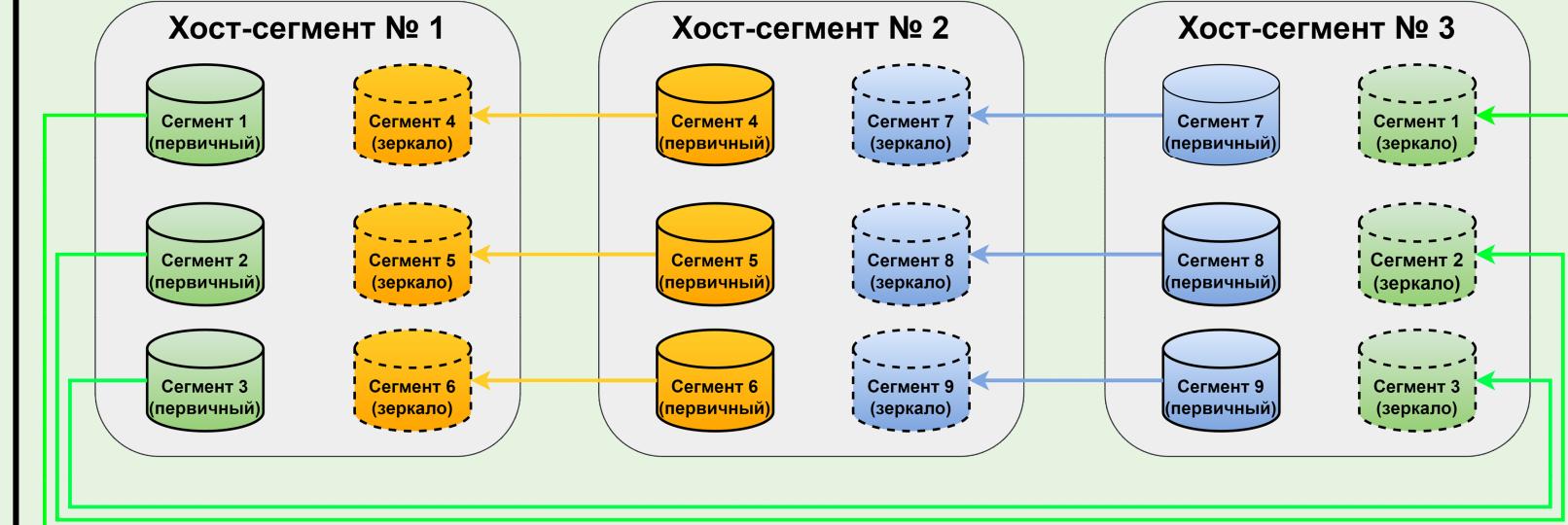


Мастер-хост кластера GreenPlum вышел из строя



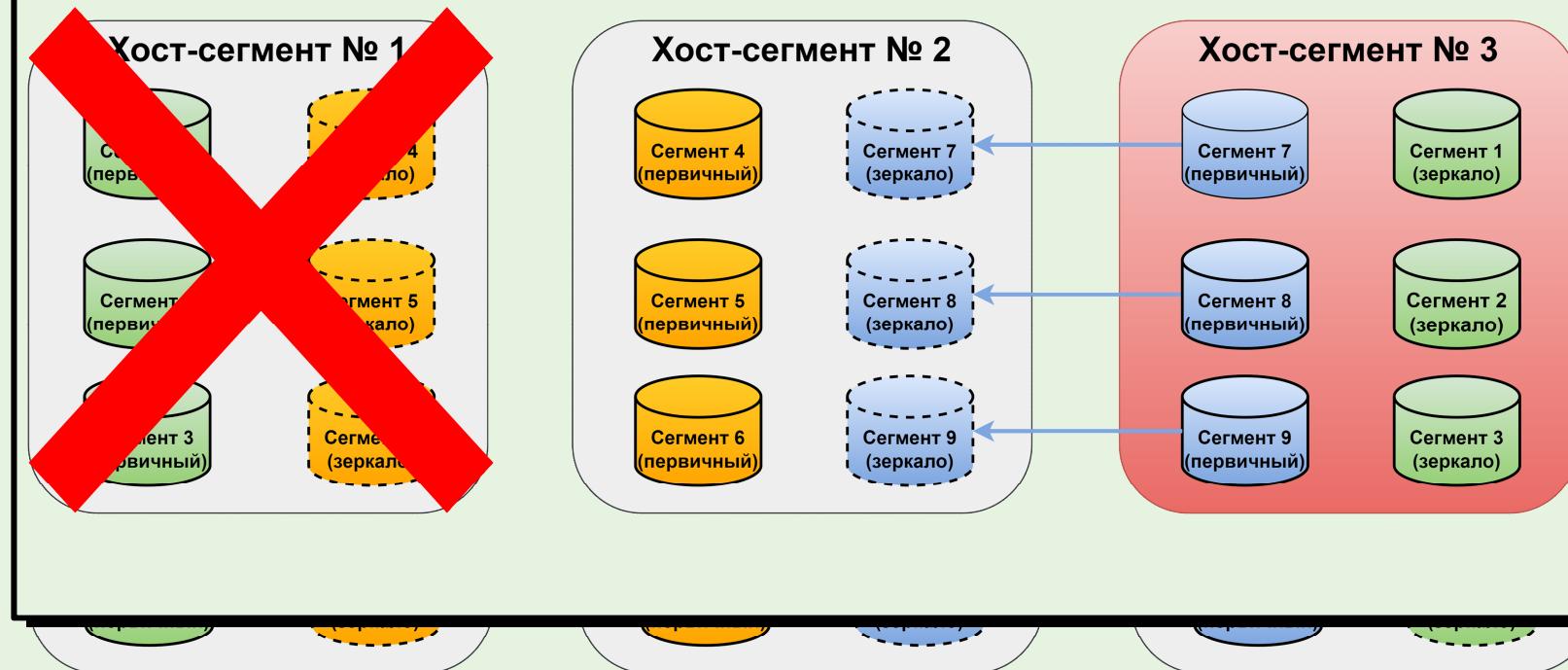
Архитектура кластера GreenPlum

Групповое зеркалирование первичных сегментов



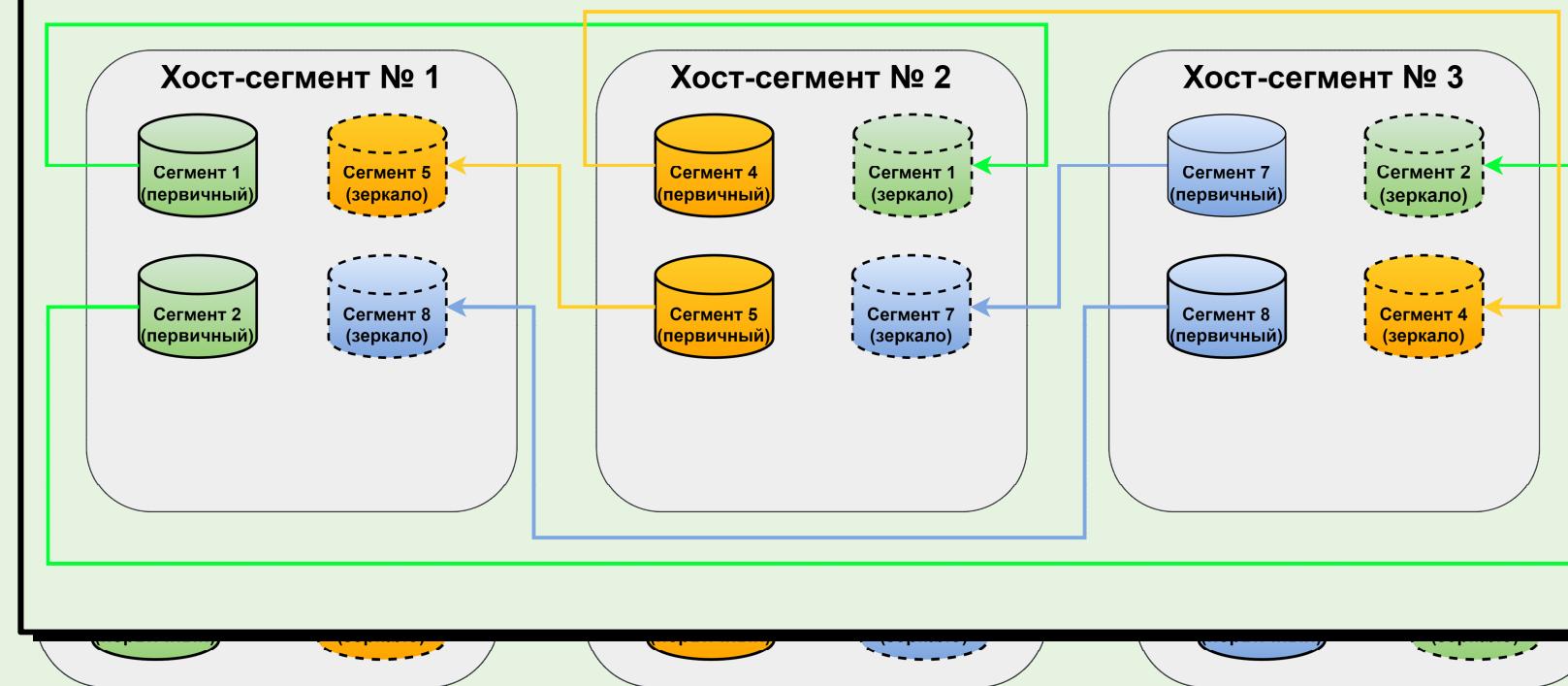
Архитектура кластера GreenPlum

Сломался хост-сегмент при групповом зеркалировании сегментов



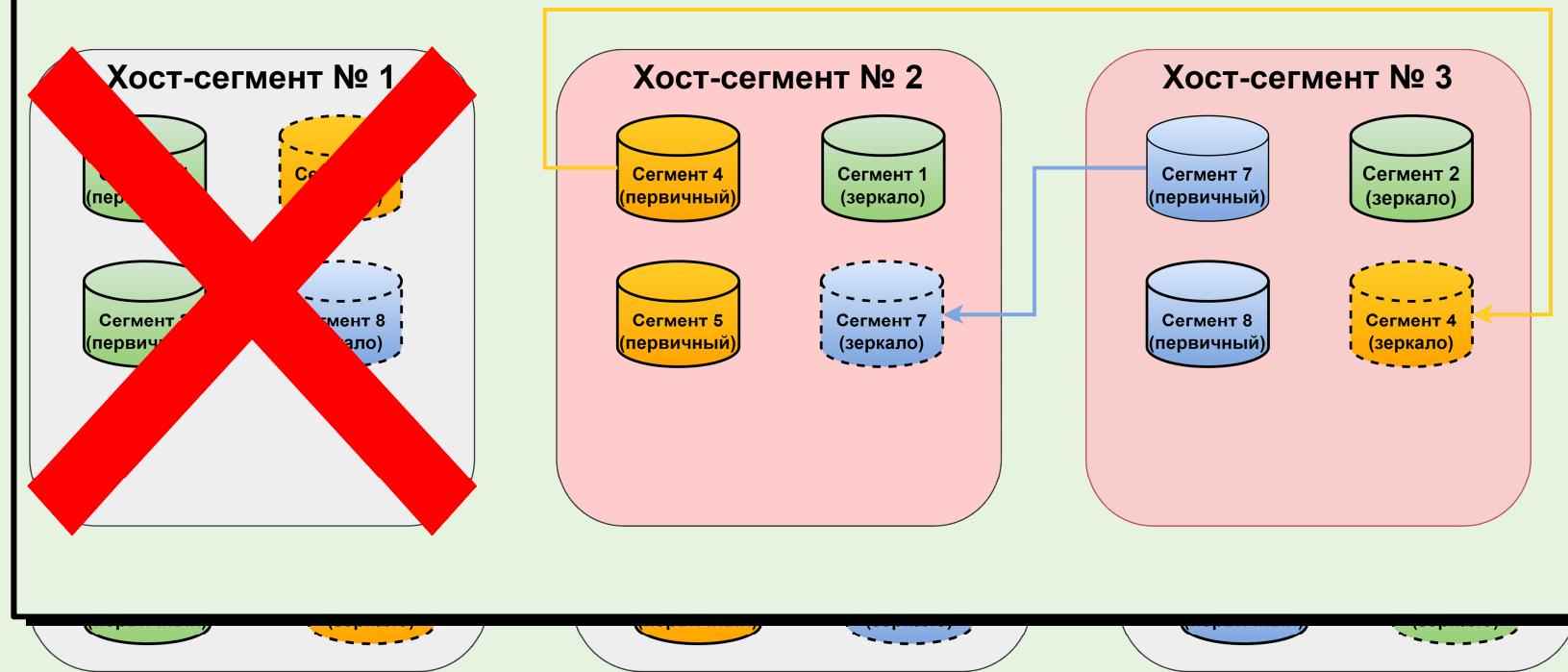
Архитектура кластера GreenPlum

Расширенное зеркалирование первичных сегментов



Архитектура кластера GreenPlum

Расширенное зеркалирование первичных сегментов



03

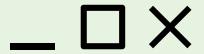
.....
>>>



Дистрибуция

- Дистрибуция и их виды в GreenPlum
- Виды таблиц, применяемые в GreenPlum

Дистрибуция



Дистрибуция — это принцип проектирования базы данных, при котором данные разбиваются на части и размещаются на разных шардах\сегментах.

Сегмент 1

Table_name		
id	date	num
1	2024-03-01	200
5	2024-03-14	7000
7	2024-05-12	1700
8	2024-05-19	3200
10	2024-06-01	30
12	2024-06-28	3270

Сегмент 2

Table_name		
id	date	num
2	2024-03-05	700
4	2024-03-12	3400
9	2024-05-30	2800

Сегмент 3

Table_name		
id	date	num
3	2024-03-06	6400
6	2024-04-17	1300
11	2024-06-17	400
13	2024-07-11	1450

При создании таблицы для обозначения дистрибуции используется ключевое слово
«DISTRIBUTED»

Дистрибуция \ Шардирование

Дистрибуция
разбиение

Дистрибуция по ключу

```
CREATE TABLE Table_name(  
    id UUID,  
    date TIMESTAMP,  
    num INT  
)  
DISTRIBUTED BY id;
```

Сегмент 1

Table_name		
id	date	num
1	2024-03-01	200
1	2024-03-14	7000
1	2024-05-12	1700
2	2024-05-19	3200
2	2024-06-01	30
2	2024-06-28	3270

Сегмент 2

Table_name		
id	date	num
3	2024-03-05	700
4	2024-03-12	3400
4	2024-05-30	2800

При создании таблицы для ее основания дистрибуции используется ключевое слово
«DISTRIBUTED»

Дистрибуция \ Шардирование

Дистрибуция
разбиение

Дистрибуция случайным образом

```
CREATE TABLE Table_name(  
    id UUID,  
    date TIMESTAMP,  
    num INT  
)  
DISTRIBUTED RANDOMLY;
```

Сегмент 1

Table_name		
id	date	num
1	2024-03-01	200
2	2024-03-14	7000
2	2024-05-12	1700
3	2024-05-19	3200
4	2024-06-01	30
4	2024-06-28	3270

Сегмент 2

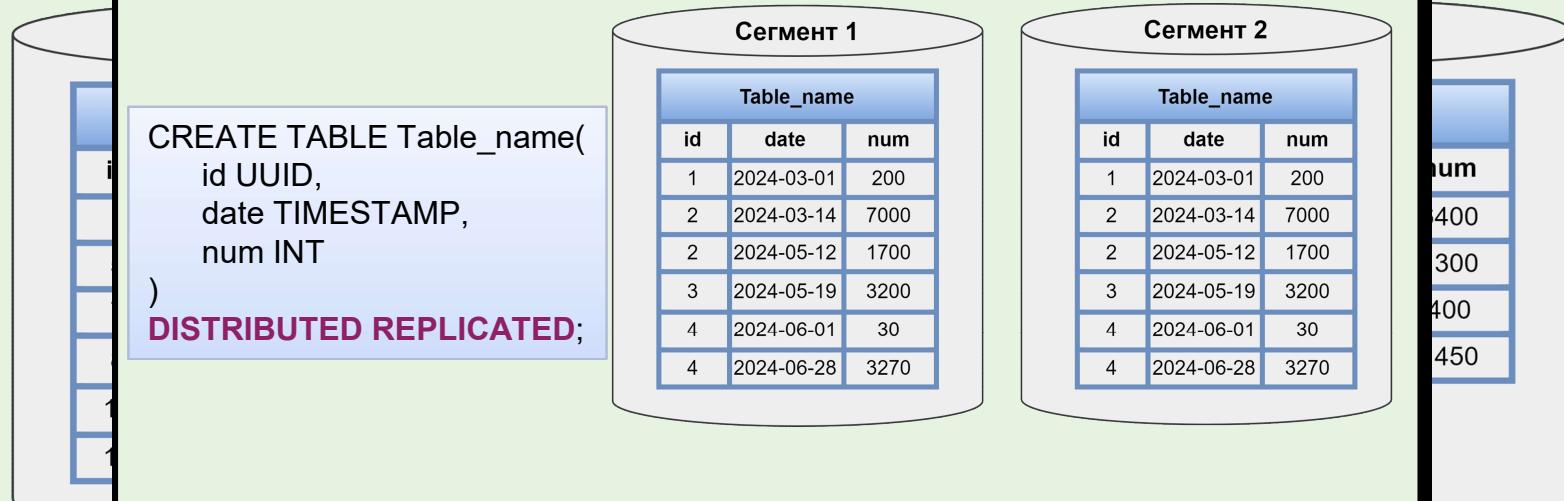
Table_name		
id	date	num
1	2024-03-05	700
2	2024-03-12	3400
3	2024-05-30	2800
4	2024-05-30	2800
4	2024-05-30	2800

При создании таблицы для ее назначения дистрибуции используется ключевое слово
«DISTRIBUTED»

Дистрибуция \ Шардирование

Дистрибуция
разбивка

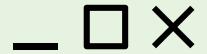
Дистрибуция таблицы по всем сегментам



```
CREATE TABLE Table_name(  
    id UUID,  
    date TIMESTAMP,  
    num INT  
)  
DISTRIBUTED REPLICATED;
```

При создании таблицы для ее назначения дистрибуции используется ключевое слово
«DISTRIBUTED»

Дистрибуция \ Шардирование



Дистрибуция — это принцип проектирования базы данных, при котором данные разбиваются на части и размещаются на разных шардах\сегментах.

Сегмент 1

Table_name		
id	date	num
1	2024-03-01	200
5	2024-03-14	7000
7	2024-05-12	1700
8	2024-05-19	3200
10	2024-06-01	30
12	2024-06-28	3270

Сегмент 2

Table_name		
id	date	num
2	2024-03-05	700
4	2024-03-12	3400
9	2024-05-30	2800

Сегмент 3

Table_name		
id	date	num
3	2024-03-06	6400
6	2024-04-17	1300
11	2024-06-17	400
13	2024-07-11	1450

При создании таблицы для д обозначения дистрибуции используется ключевое слово
«DISTRIBUTED»

Виды таблиц применяемые в GP

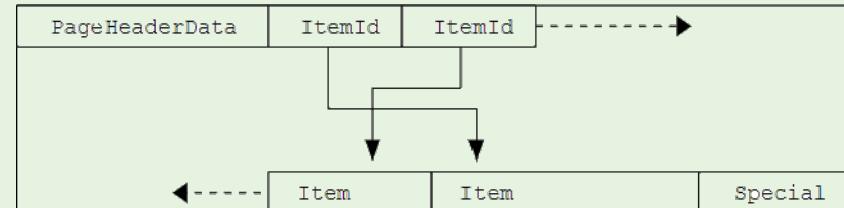


Основные виды таблиц:

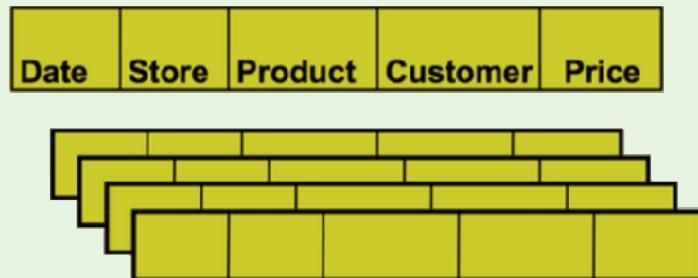
1. HEAP-таблица
2. APPEND-OPTIMIZED таблица

Дополнительные виды:

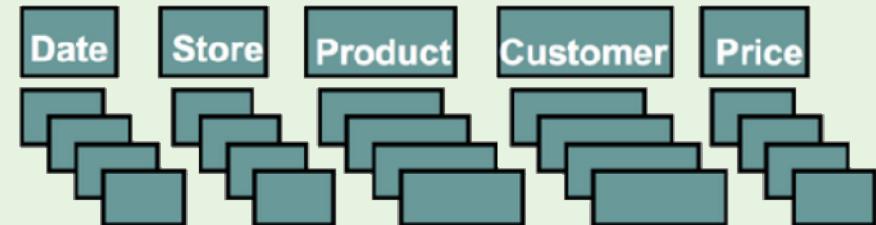
1. Временная таблица
2. Нежурналируемая таблица
3. Внешняя таблица



row-store



column-store



Виды таблиц применяемые в GP

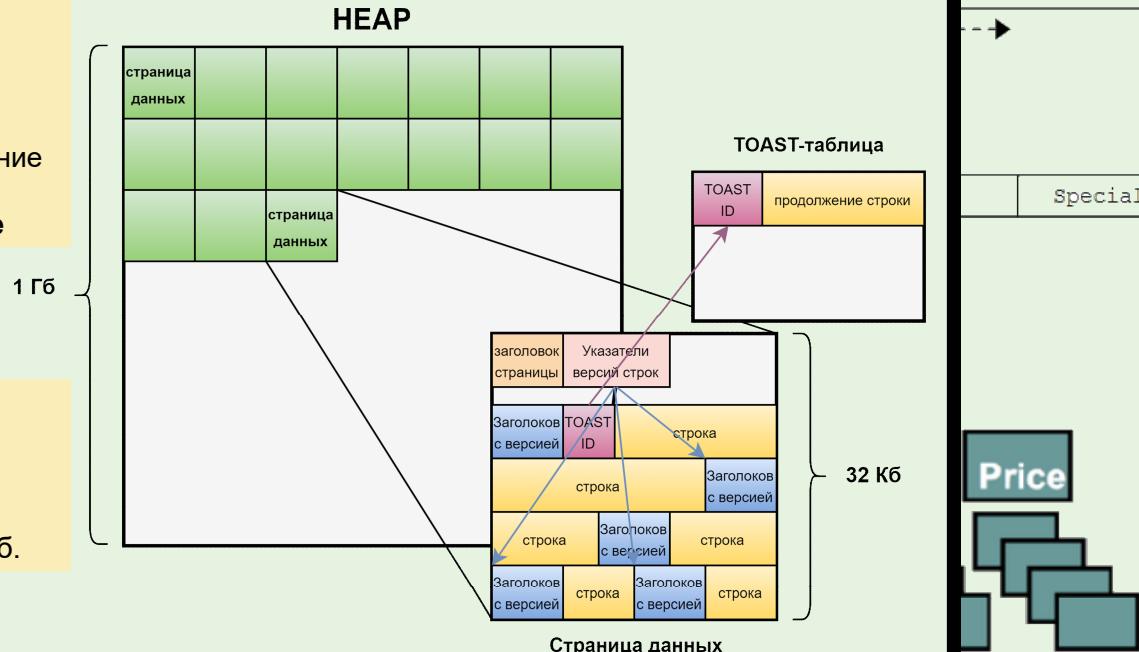
HEAP таблица

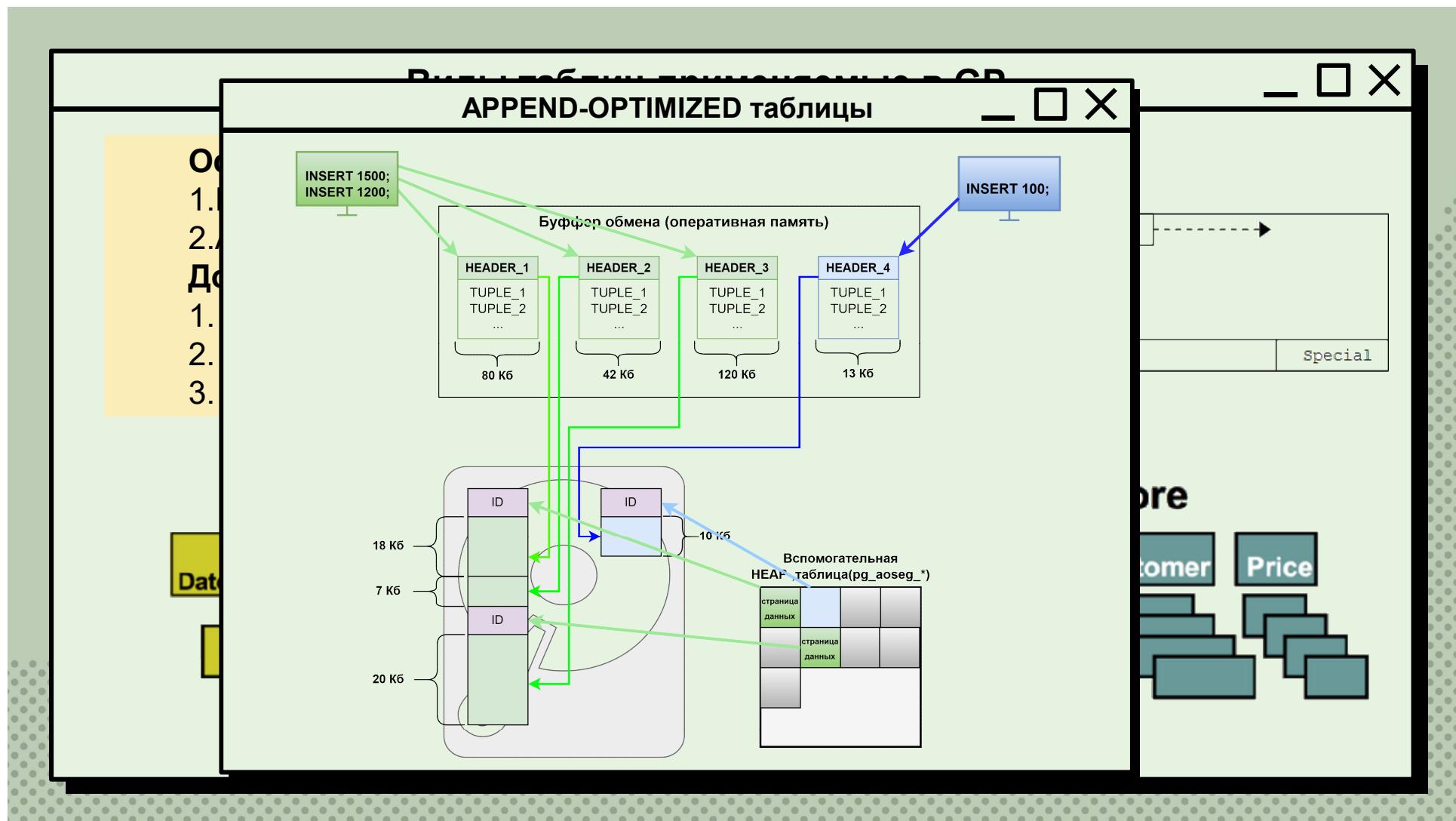
Ограничения:

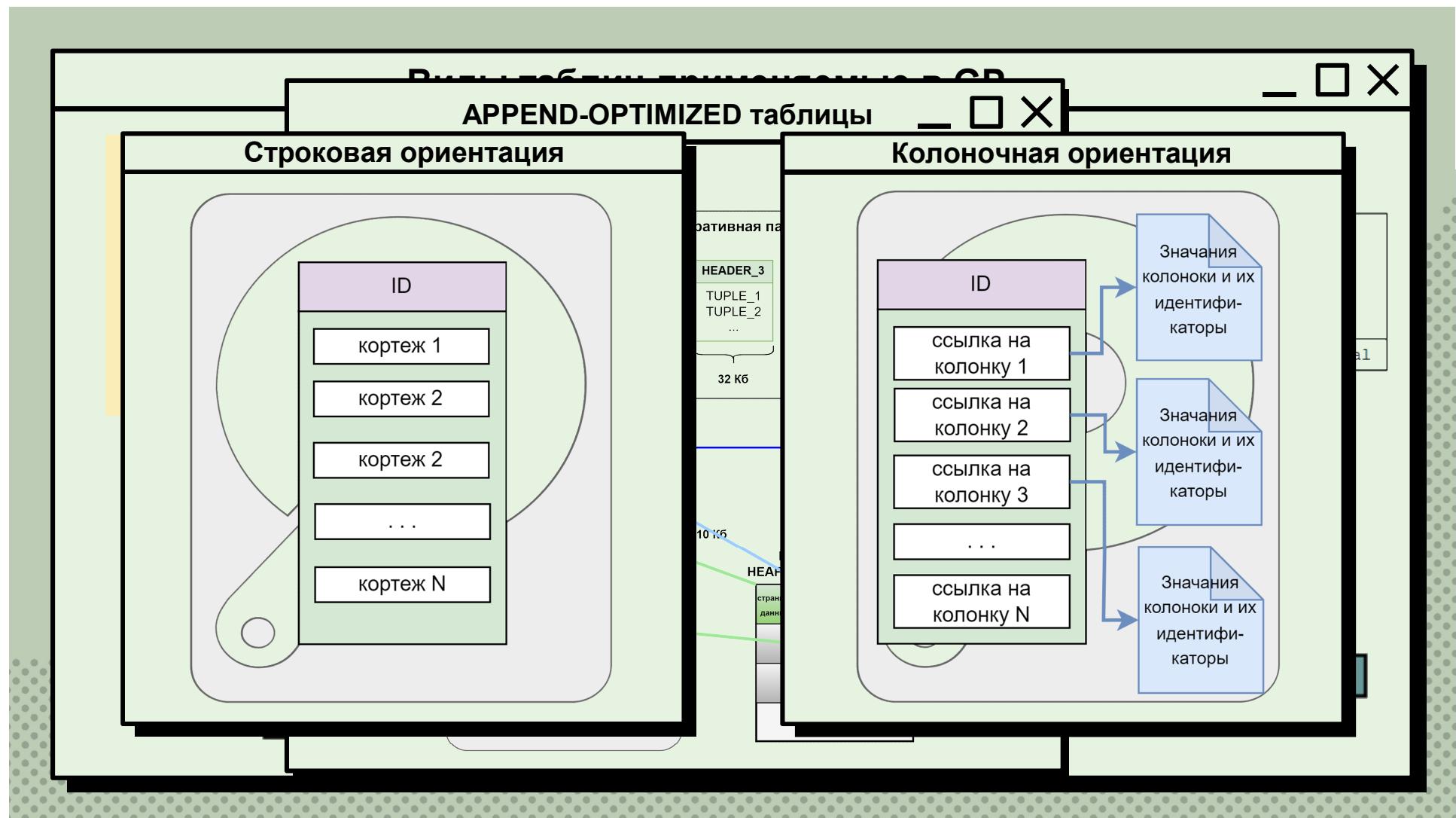
- HEAP-таблицы поддерживают только **построчно-ориентированное хранение данных**.
- Не поддерживают **сжатие**

Использование:

- Если в таблице часто меняются данные;
- Для таблиц **малого размера**, например до Гб.







Виды таблиц применяемые в GP

Временные таблицы

Временные таблицы — живут в рамках одной сессии и автоматически уничтожаются при завершении сессии.

```
CREATE TEMPORARY TABLE temp_table(  
    id_col SERIAL,  
    col_2 TEXT)  
WITH (  
    appendoptimized=TRUE,  
    orientation=column,  
    compressstype=zstd,  
    compresslevel=1)  
DISTRIBUTED BY (id);
```

Применение:

Сохранение данных промежуточного расчёта и использования их в дальнейшем.

Виды таблиц применяемые в GP

Нежурналируемая таблица

Нежурналируемые таблицы — не записывают свои транзакции в WAL-журнал, поэтому не зеркалируются.

```
CREATE UNLOGGED TABLE unlogged_table (
    id_col SERIAL,
    col_2 TEXT)
WITH (
    appendoptimized=true,
    orientation=row)
DISTRIBUTED BY (id);
```

Преимущества и недостатки:

- ✓ Скорость записи в такую таблицу увеличивается.
- ✗ Уменьшается надёжность хранения данных.

Виды таблиц применяемые в GP

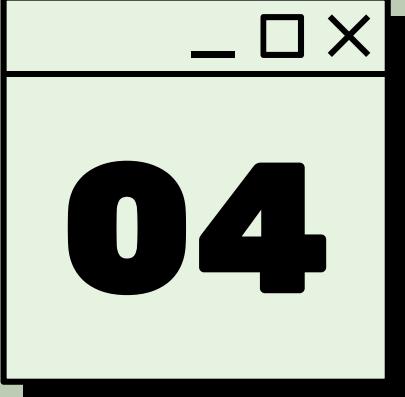


Внешние таблицы

Внешние таблицы —
предназначены для
параллельной загрузки
данных в Greenplum или
выгрузки из него.

```
CREATE EXTERNAL TABLE ext_expenses (
    id_col int,
    col_2 TEXT)
LOCATION (
    'gpfdist://localhost:2107/*.csv',
    'gpfdist://localhost:2108/*.csv ')
FORMAT 'CSV' ( DELIMITER ',' );
```

Special



04

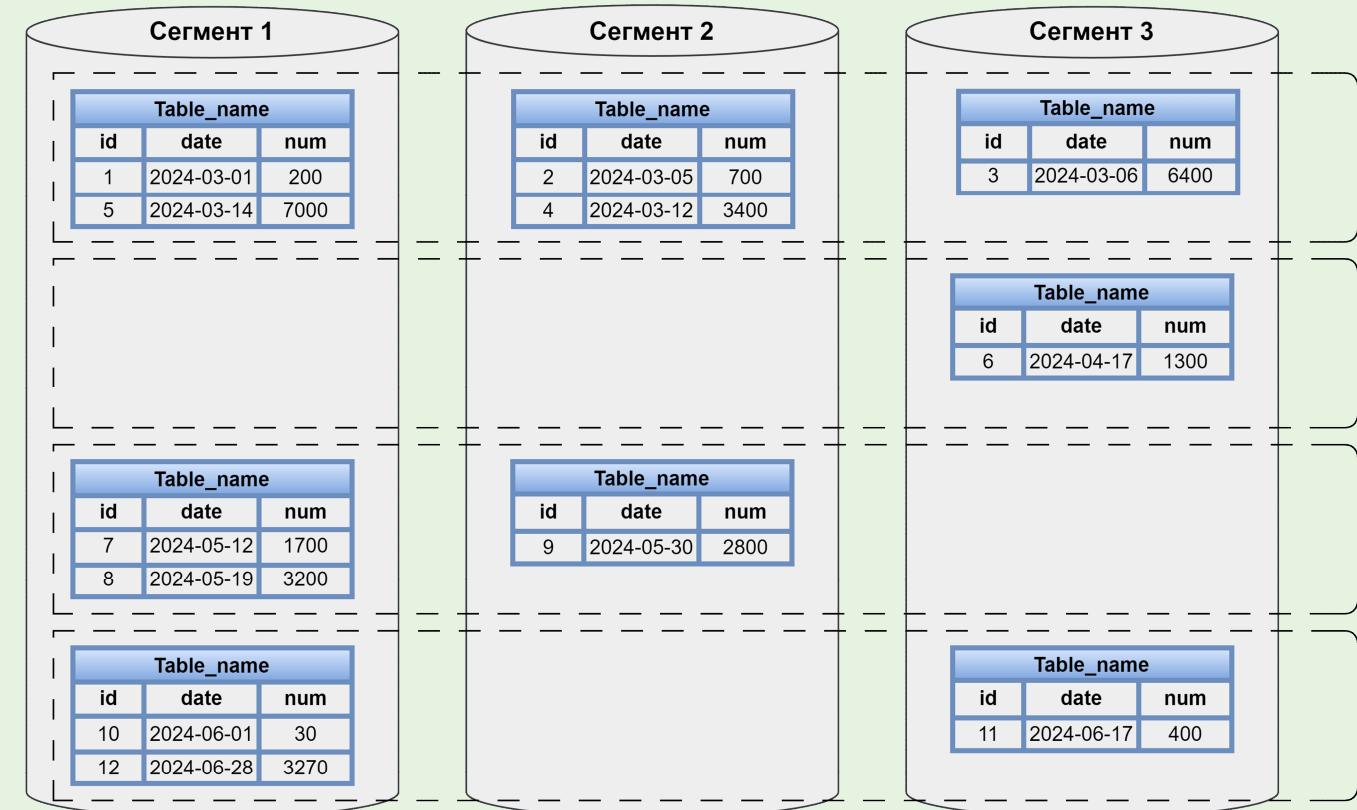
.....
>>>



Партицирование

- Партицирование и их виды в GreenPlum
- Групповое партицирование в GreenPlum

Партицирование



партиция
2024-03

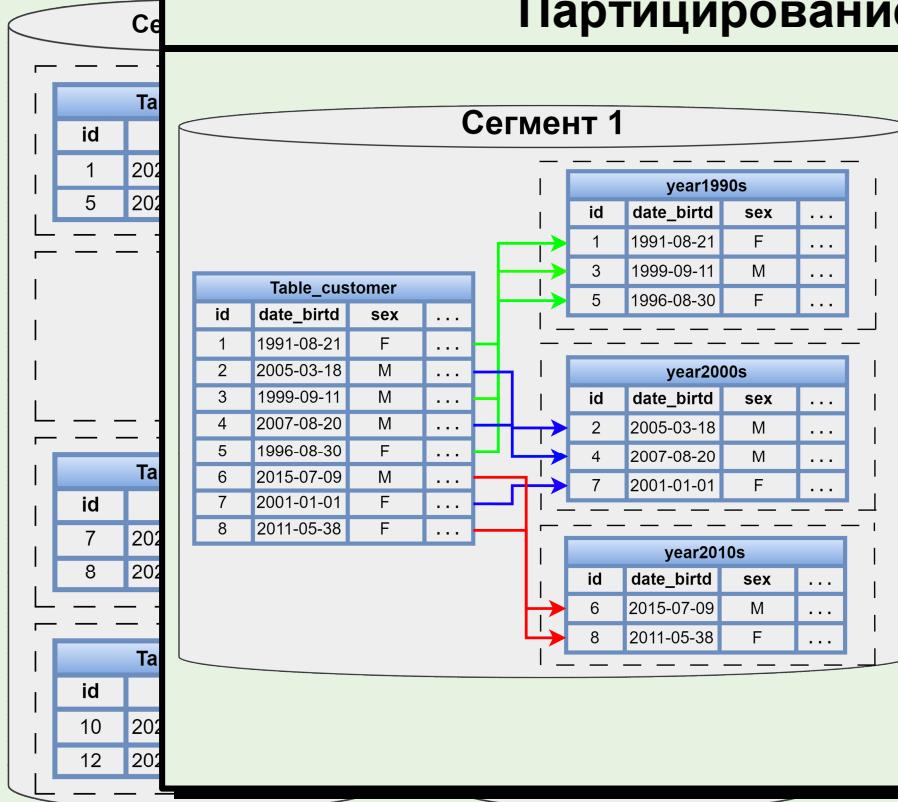
партиция
2024-04

партиция
2024-05

партиция
2024-06

Партицирование

Партицирование по диапазону



CREATE TABLE ...

WITH ...

DISTRIBUTED BY ...

PARTITION BY RANGE (date_birth)

(

PARTITION year1990s START ('1990-01-01')
END ('2000-01-01')

WITH (appendoptimized=TRUE,
compressstype=zstd),

PARTITION year2000s START ('2000-01-01')
END ('2010-01-01'),

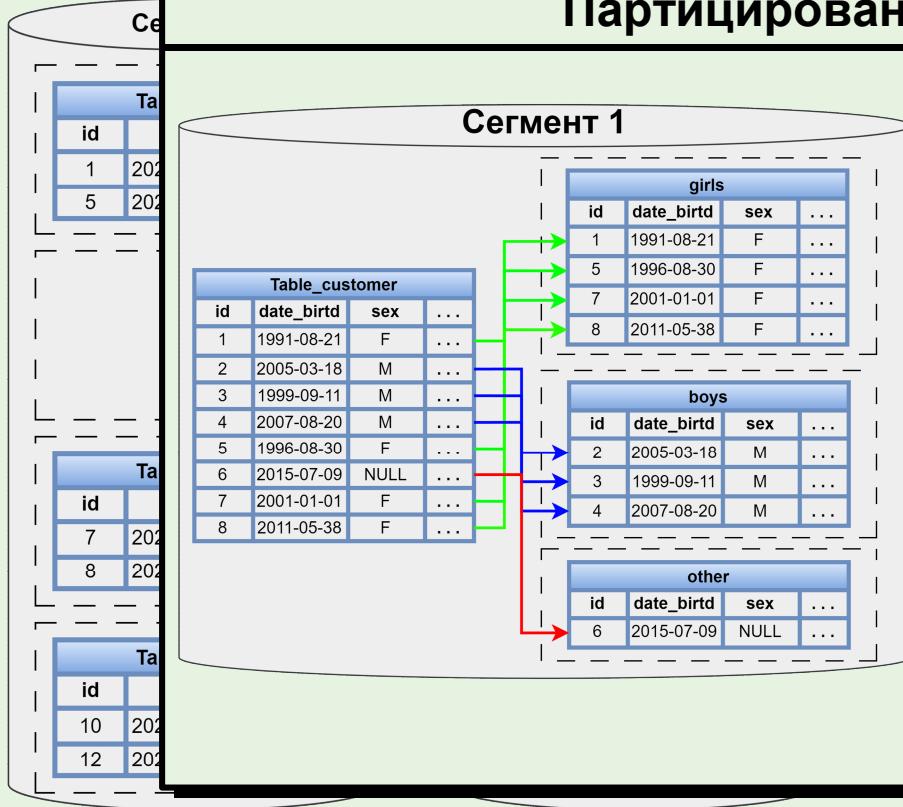
PARTITION year2010s START ('2010-01-01')
END ('2020-01-01'));

(START ('1990-01-01 00:00:00')
END ('2020-01-01 00:00:00'))

EVERY (INTERVAL '10 YEAR'));

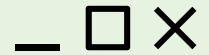
Партицирование

Партицирование по списку

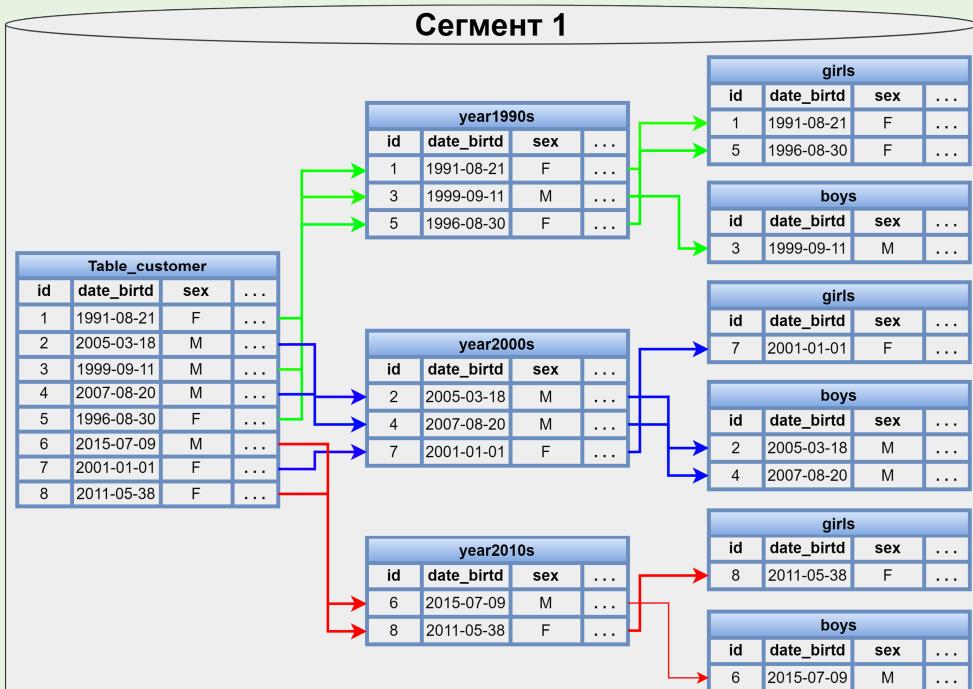


```
CREATE TABLE ...  
WITH...  
DISTRIBUTED BY ...  
PARTITION BY LIST (sex)  
( PARTITION girls VALUES ('F'),  
PARTITION boys VALUES ('M'),  
DEFAULT PARTITION other );
```

Партицирование



Многоуровневое партицирование



```
CREATE TABLE ...  
WITH ...  
DISTRIBUTED BY ...  
PARTITION BY RANGE (date_birth)  
SUBPARTITION BY LIST (sex)  
SUBPARTITION TEMPLATE (  
    SUBPARTITION girls VALUES ('F'),  
    SUBPARTITION boys VALUES ('M')  
)  
(START ('1990-01-01')  
END ('2025-01-01')  
EVERY (INTERVAL '10 YEAR'));
```

05

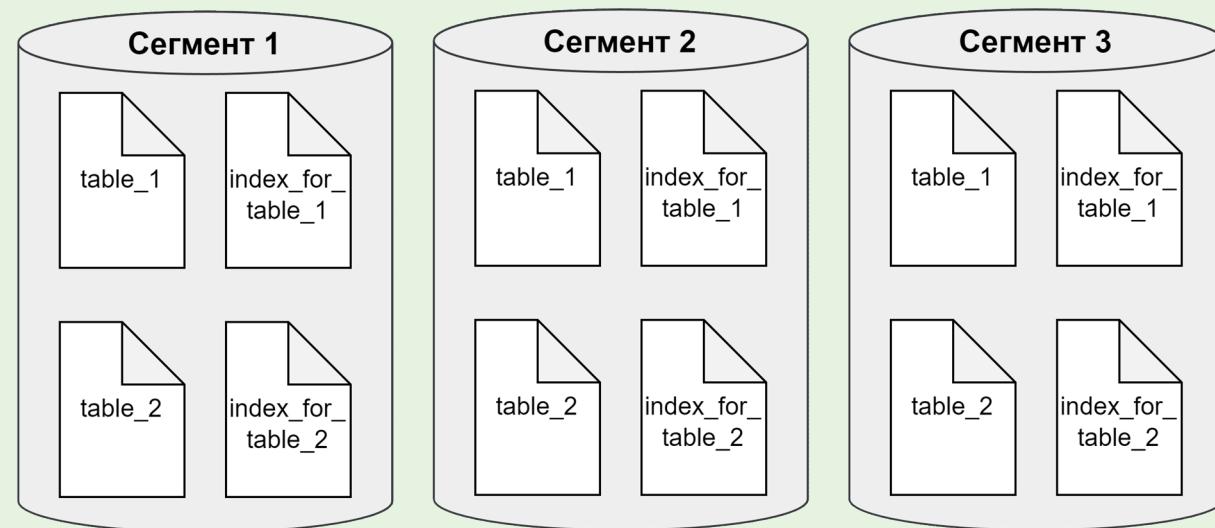
Индексы

- Индексы и их виды в GreenPlum
- B-tree индекс
- Bitmap индекс

Партицирование



Индекс базы данных – это отдельная структура данных, которая создается на основе с данных в таблице, которая ускоряют процесс извлечения этих данных, и тем самым повышая производительность запросов.



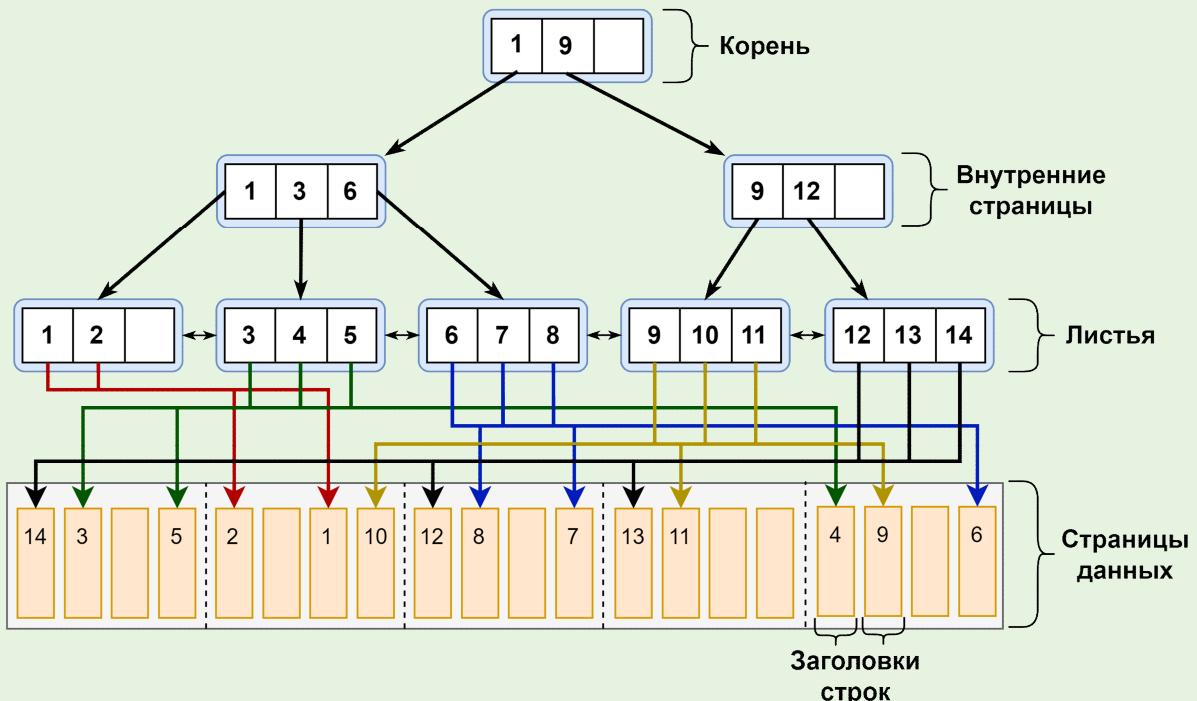
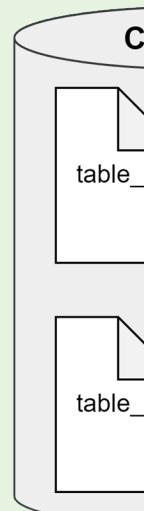
Виды индексов:

- B-tree
- Bitmap
- GIST
- SP-GIST
- GIN

Партицирование

B-tree индекс

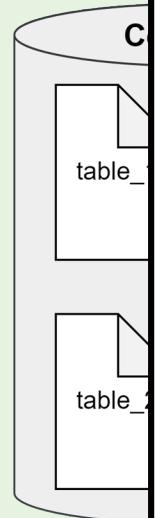
Индекс
в таблице
произво



Партицирование

Индекс
в таблице
представляет

B-tree индекс



Table_customer			
id	date_birtd	sex	age
1	1991-08-21	F	33
2	2005-03-18	M	19
3	1999-09-11	M	25
4	2007-08-20	M	17
5	1996-08-30	F	28
6	2015-07-09	M	9
7	2001-01-01	F	23
8	2011-05-38	F	13

	10001011
F	10001011

	10000000
33	10000000
19	01000000
25	00100000
17	00010000
28	00001000
9	00000100
23	00000010
13	00000001

SELECT *
FROM Table_customers
WHERE age > 20 and sex = 'F'

10001011
&&
(10000000 ||
00100000 ||
00001000 ||
00000010)

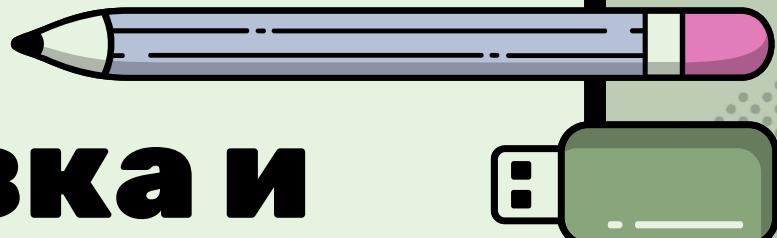
10001011
&&
10101010

Result			
id	date_birtd	sex	age
1	1991-08-21	F	33
5	1996-08-30	F	28
7	2001-01-01	F	23

05

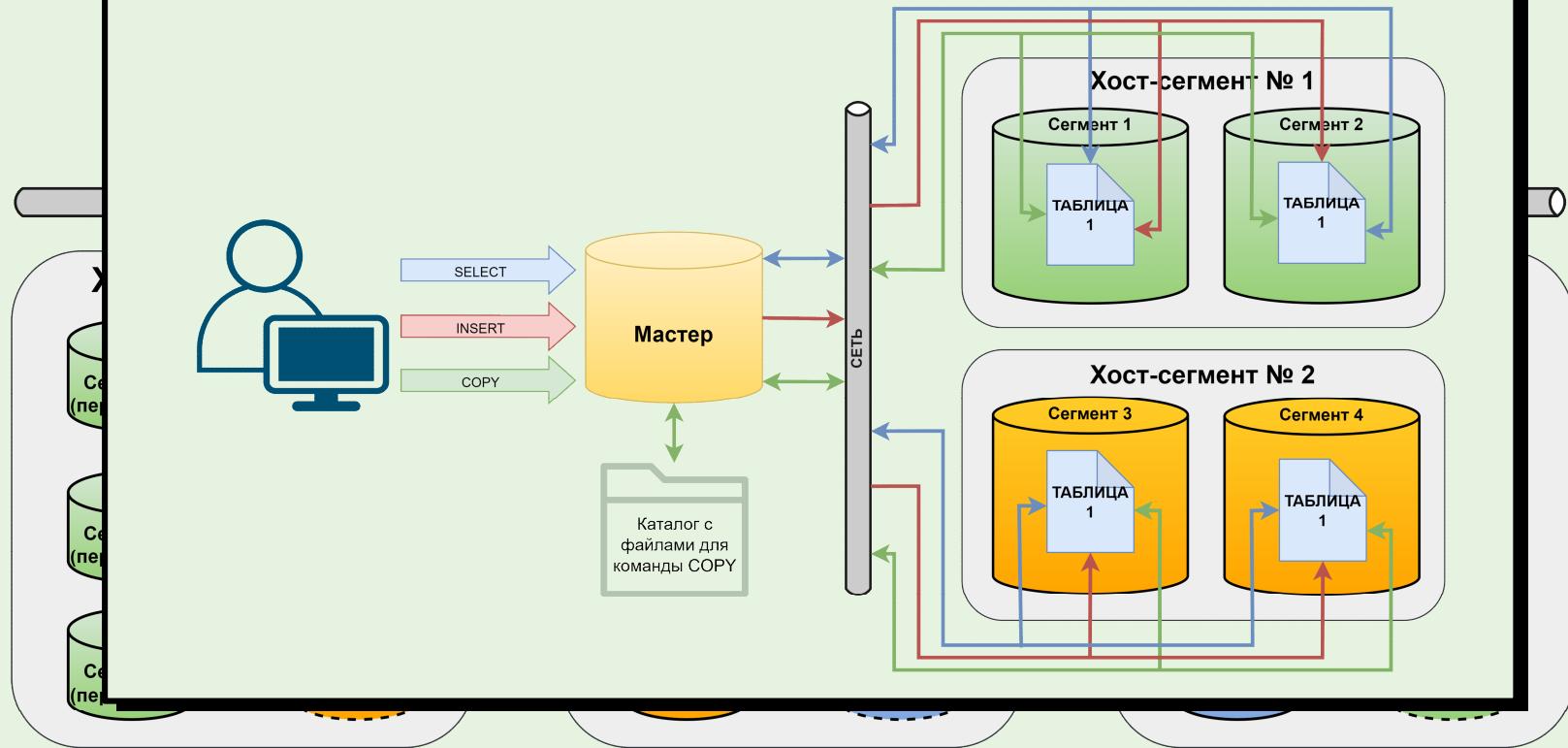
Загрузка и выгрузка данных

- SQL-команды для загрузки и выгрузки данных
- Внешние таблицы Greenplum
- Протоколы доступа к внешним данным



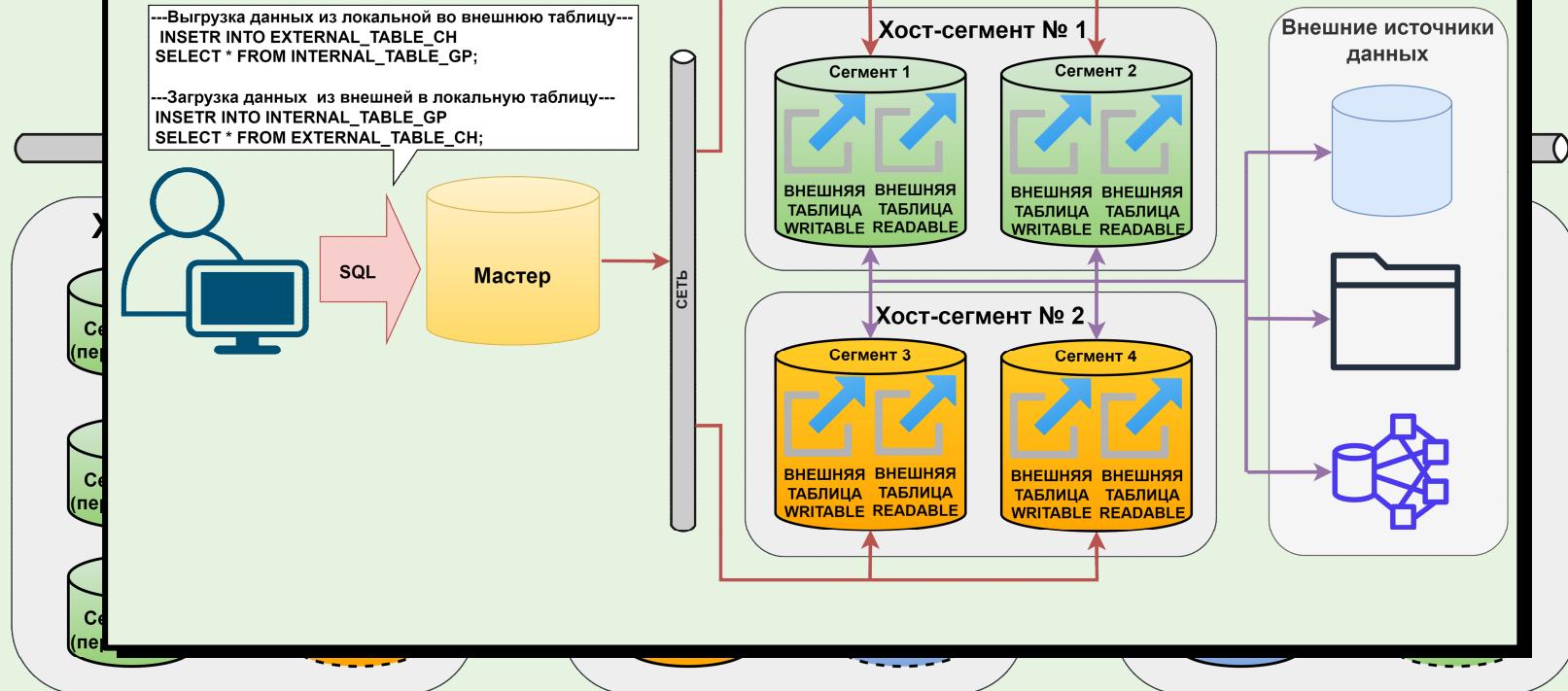
Загрузка и выгрузка данных в Greenplum

SQL-команды



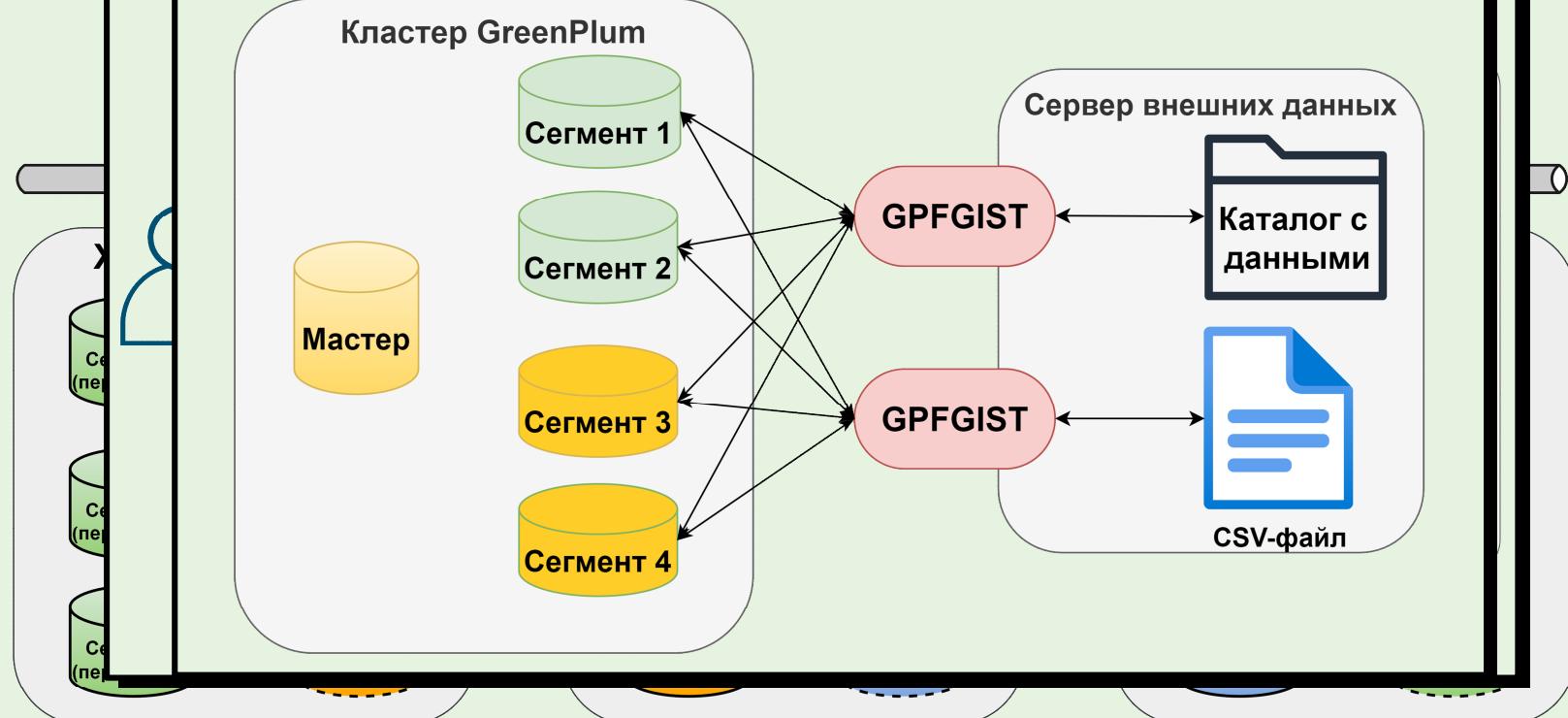
Загрузка и выгрузка данных в Greenplum

Внешние таблицы



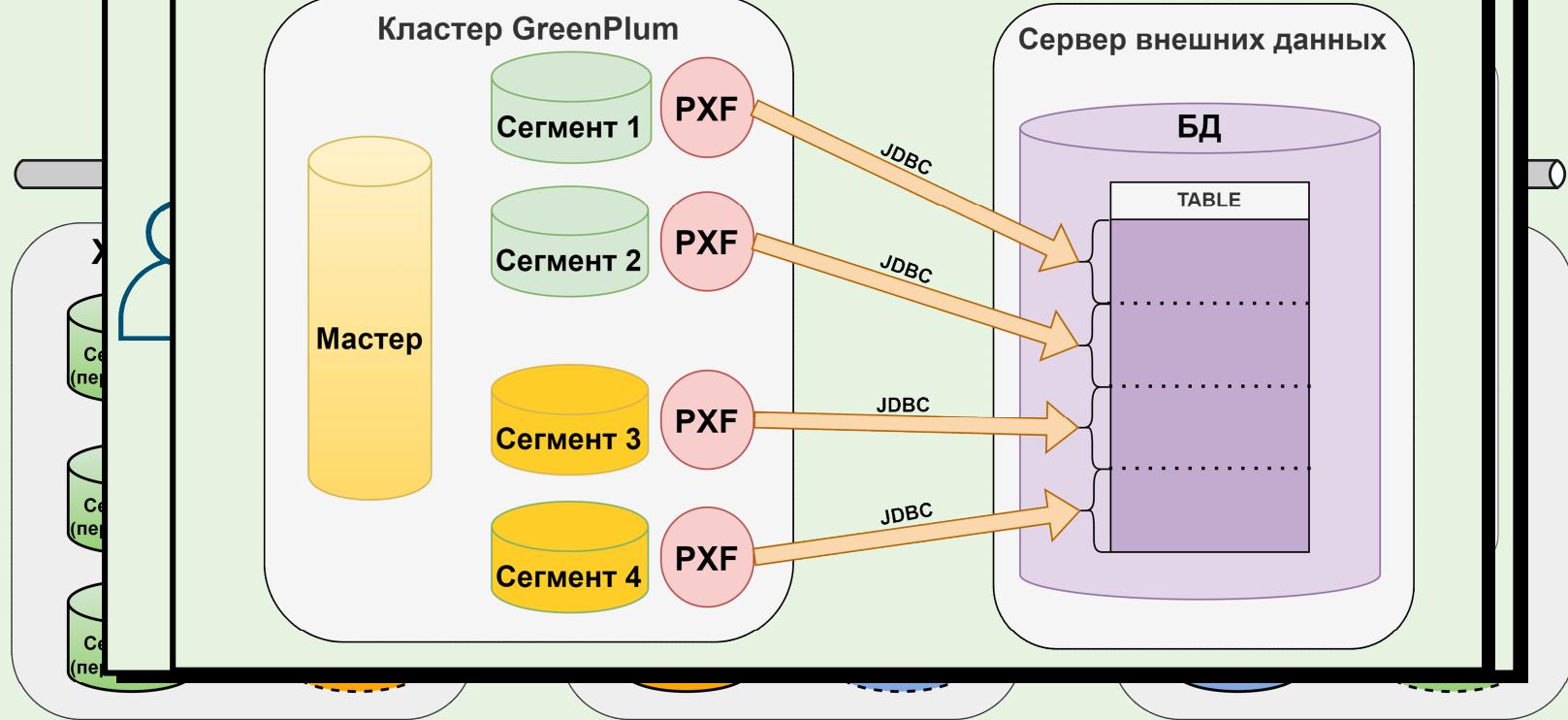
Загрузка и выгрузка данных в Greenplum

Протокол GPFDIST



Загрузка и выгрузка данных в Greenplum

Протокол PXF



06

Оптимизация

запросов

- Дерево выполнения запросов
- Виды операторов в плане запросов
- Hash Join, Merge Join, Nested Loop
- На что стоит обратить внимание в плане запросов

Дерево выполнения запроса



EXPLAIN

EXPLAIN ALALYZE

A-Z QUERY PLAN

Gather Motion 24:1 (slice1; segments: 24) (cost=0.00..1293.11 rows=190 width=98)

```
> Hash Left Join (cost=0.00..1293.06 rows=8 width=98)
  Hash Cond: (t1.user_sk = t2.user_sk)
  Join Filter: (t1.date >= (min(t2.date)))
-> Hash Left Join (cost=0.00..862.05 rows=8 width=102)
  Hash Cond: (t1.user_sk = t2.user_sk)
  Join Filter: (t1.date >= (min(t2.date)))
-> HashAggregate (cost=0.00..431.01 rows=8 width=90)
  Group Key:
    -> Seq Scan on t1 (cost=0.00..431.01 rows=8 width=90)
      Filter: issued_conversions
-> Hash (cost=431.01..431.01 rows=34 width=12)
  -> HashAggregate (cost=0.00..431.01 rows=34 width=12)
    Group Key:
      -> Seq Scan on t2 (cost=0.00..431.01 rows=34 width=12)
        Filter: conversions
-> Hash (cost=431.01..431.01 rows=8 width=12)
  -> GroupAggregate (cost=0.00..431.01 rows=8 width=12)
    Group Key:
      -> Sort (cost=0.00..431.01 rows=8 width=12)
        Sort Key:
          -> Seq Scan on t2 (cost=0.00..431.01 rows=8 width=12)
            Filter: issued_conversions
```

A-Z QUERY PLAN

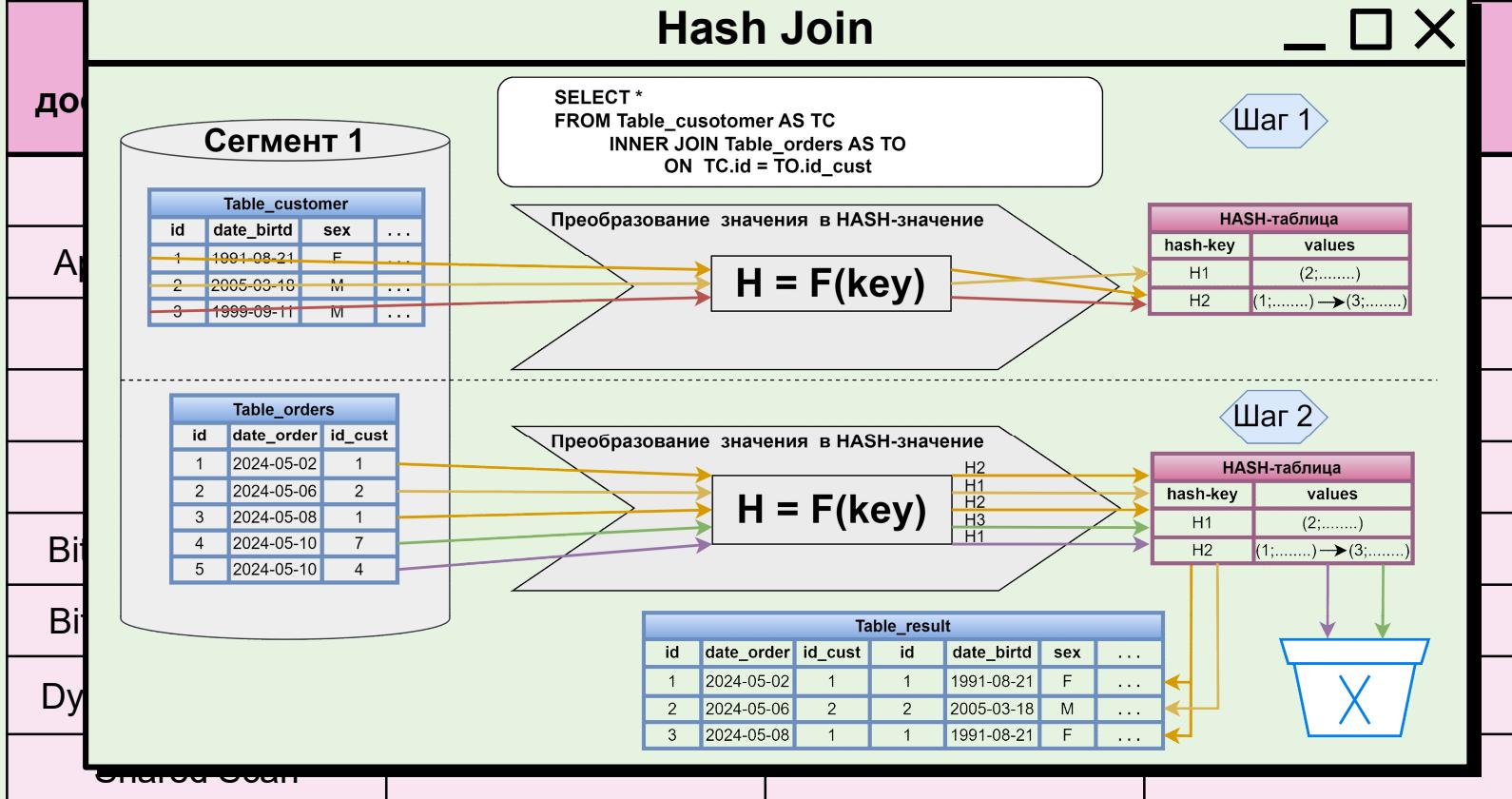
Виды операторов в плане запросов

— □ ×

Операторы доступа к данным	Операторы преобразований данных	Операторы перемещения данных	Служебные операторы
Seq Scan	Limit	Redistribute Motion	Result
AppendOnly Scan	Sort	Broadcast Motion	Append
Column Scan	Hash Join	Gather Motion	Sequence
External Scan	Merge Join		Subplan
Index Scan	Nested Loop Join		Materialize
Bitmap Index Scan	HashAggregate		Partition Selector
Bitmap Heap Scan	GroupAggregate		Hash
Dynamyc Seq Scan	WindowAgg		
Shared Scan			

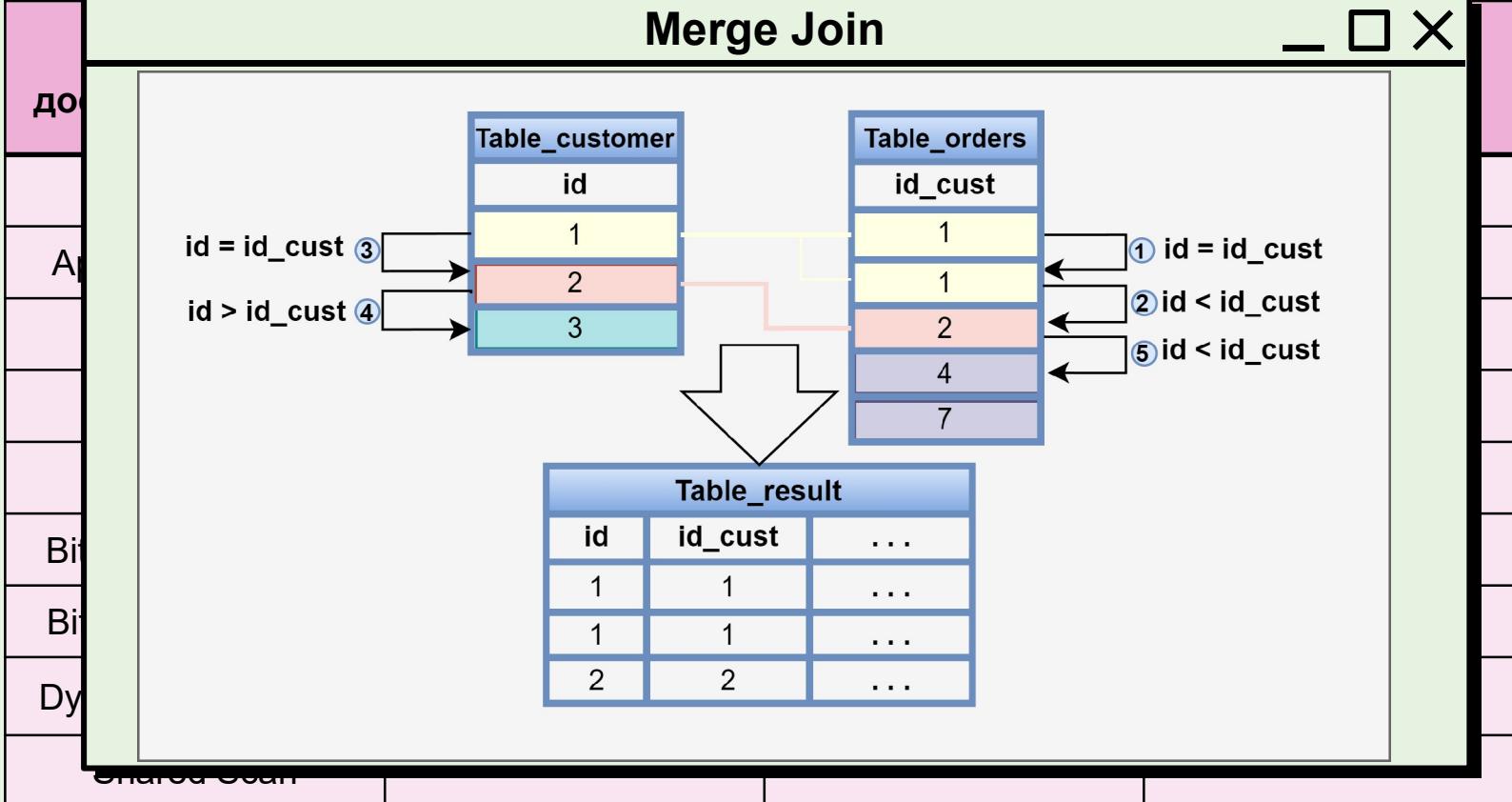
Виды операторов в плане запросов

Hash Join



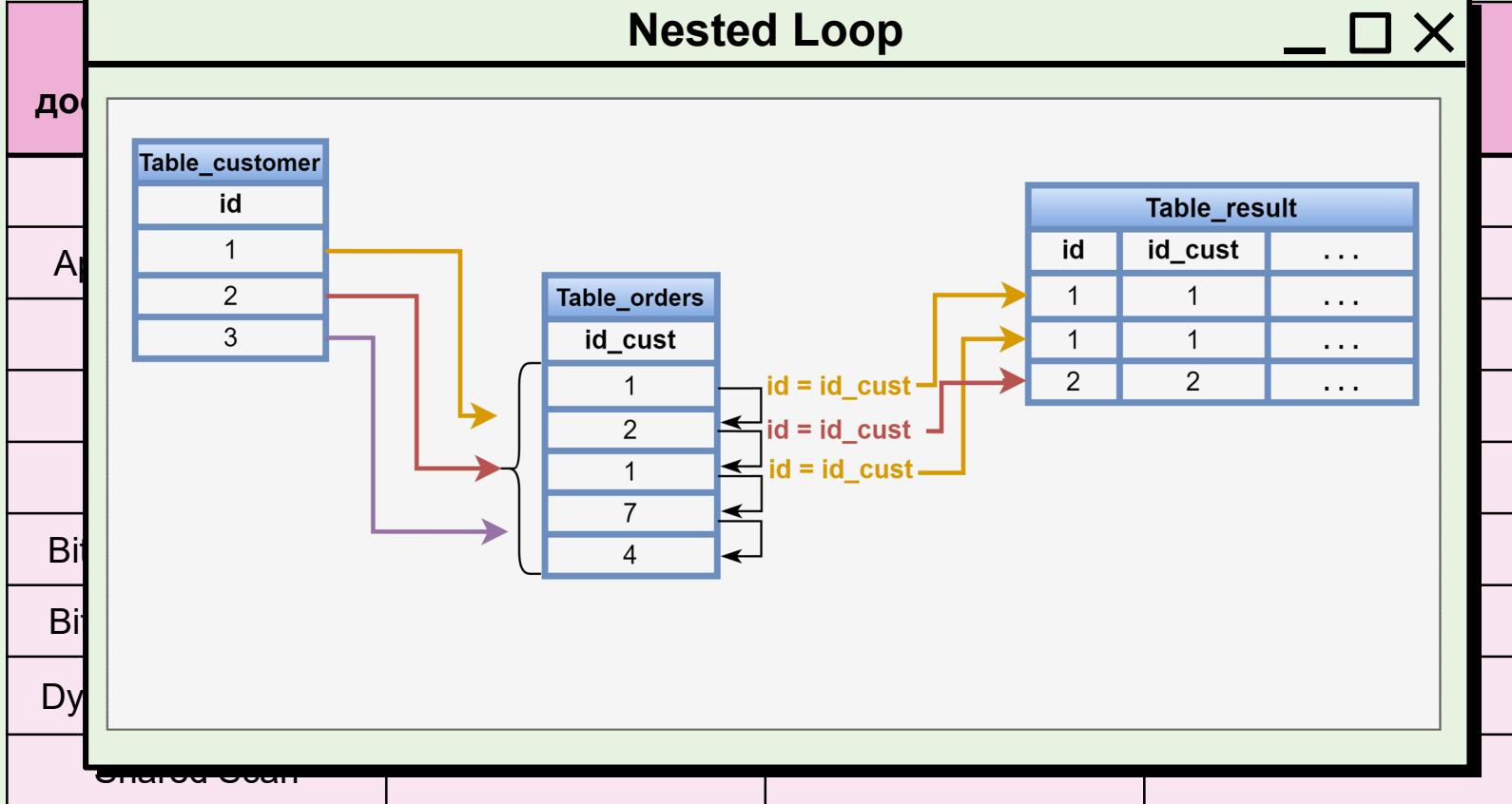
Виды операторов в плане запросов

Merge Join



Виды операторов в плане запросов

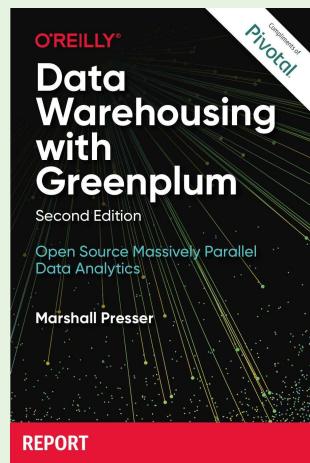
Nested Loop



Литература и ROADMAP-проект



Литература



VMware Greenplum 6 Documentation

VMware Greenplum 6



Roadmap

Files

shust/greenplum

Go to file

- DWH
- GREENPLUM
- Git
- HADOOP
- Linux
- NF
- PET_PROJECT
- files
- png
- .DS_Store
- .gitignore
- README.md

HalltapeRoadmapDE / GREENPLUM /

ShustGF · png center test · 6c1e470 · last week · History

This branch is 3 commits ahead of, 12 commits behind main.

Name	Last commit message	Last commit date
...		
README.md	png center test	last week

README.md

GREENPLUM

GreenPlum --- это массивно-параллельная, транзакционная аналитическая база данных построенная на основе БД PostgreSQL. В основу архитектуры GreenPlum заложена MPP-архитектура.

  **GREENPLUM DATABASE®**

 **PostgreSQL**