

Pivotal®

Решение средствами Pivotal HAWQ задачи СОРМ



*“It does look similar—but this one
is powered by Hadoop”*

А.Ермаков – Pivotal

Pivotal™

Платформа Pivotal

Advanced Analytics



Pivotal
Greenplum
Database



Pivotal
HAWQ

Apps at Scale



Pivotal GemFire



Redis



Rabbit MQ



Pivotal Labs &
Data Science Labs



Pivotal Cloud Foundry

Data Processing



Spring XD



Spark



Pivotal HD & Open Data Platform

Commodity Hardware

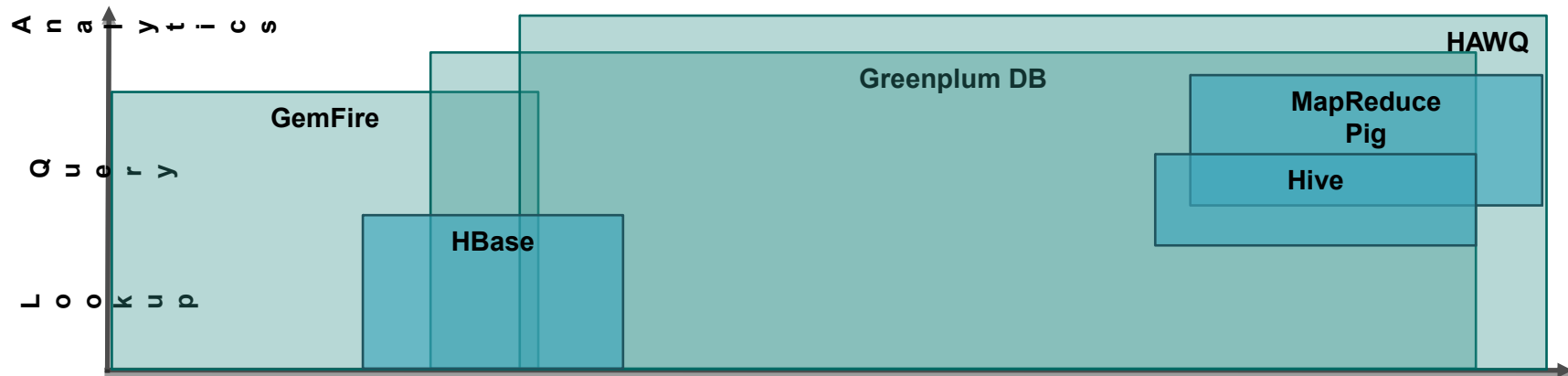
Appliance

Hybrid Cloud

Pivotal Cloud Foundry

Pivotal™


Data Access



Real Time

- NoSQL Object Store
- Transaction Support


GemFire



Interactive

- MPP DB
- Polymorphic Data Storage


Greenplum



Interactive

- SQL on Hadoop
- Query Optimizer


HAWQ



Batch

- Real-time looking
- Simple data queries

Hbase



CORM (Постановка задачи)

CORM-3 — система для обеспечения долгосрочного хранения и оперативного доступа к данным об абонентах оператора и оказанных им услугах связи.

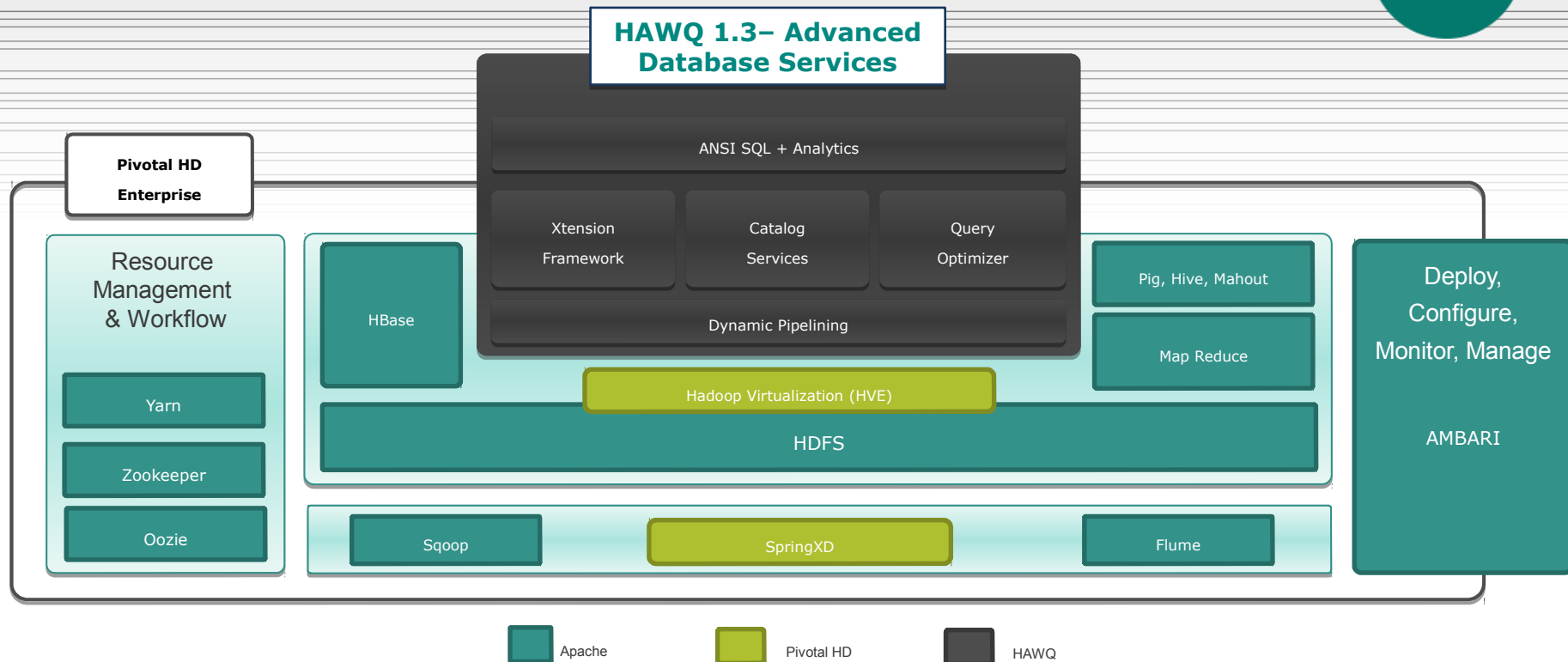
Результатом работы системы является структурированный ответ, содержащий информацию об абоненте, номере телефона, звонках, посещенных сайтах, сессиях, прокси и др. по одному или нескольким параметрам.



Почему Hadoop?

- Экономическая эффективность (TCO, ROI)
- Производительность и масштабируемость
- Большие объёмы данных
- Хранение данных в исходном формате
- Возможность сжатия данных (zlib, quicklz, gzip)
- Устойчивость к отказам приложений
- Прозрачная масштабируемость приложений

Pivotal HD 3.0

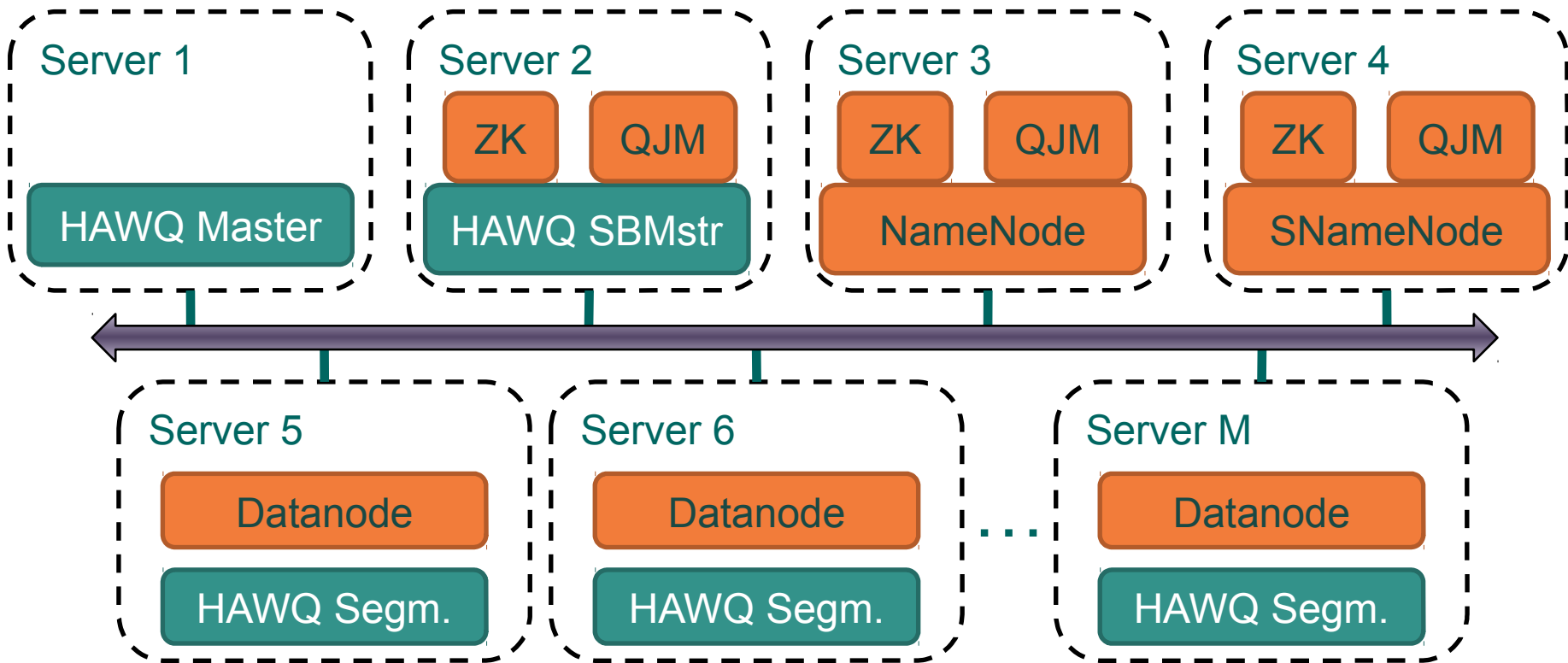


HAWQ



- Полная поддержка ANSI SQL
- Параллельная обработка запросов
- Лучший в классе оптимизатор
- Горизонтальное масштабирование
- Поддержка неструктурированных файлов
- Углубленная аналитика (MADLib)
- Поддержка шифрования данных и аутентификации пользователей

Компоненты Pivotal HAWQ



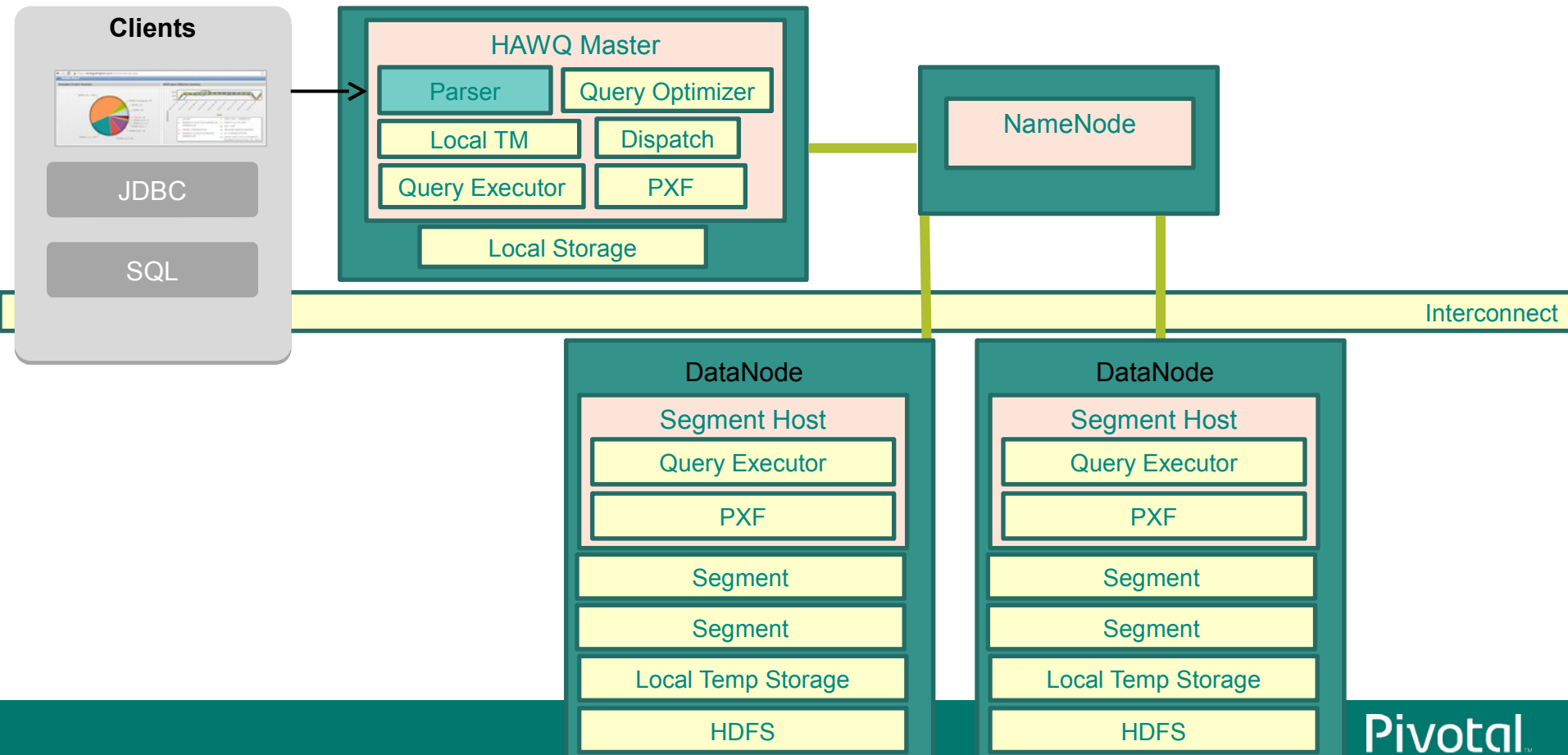
Компоненты Pivotal HAWQ

- И Master, и Segment'ы являются модифицированными процессами postgres
- При открытии подключения к Master-серверу – форк postmaser для работы с сессией
- При отправке процесса на выполнение – форк postmaster на сегментах
- План выполнения запроса разделяется на независимые блоки, slice, каждый из которых является отдельным процессом на сегменте

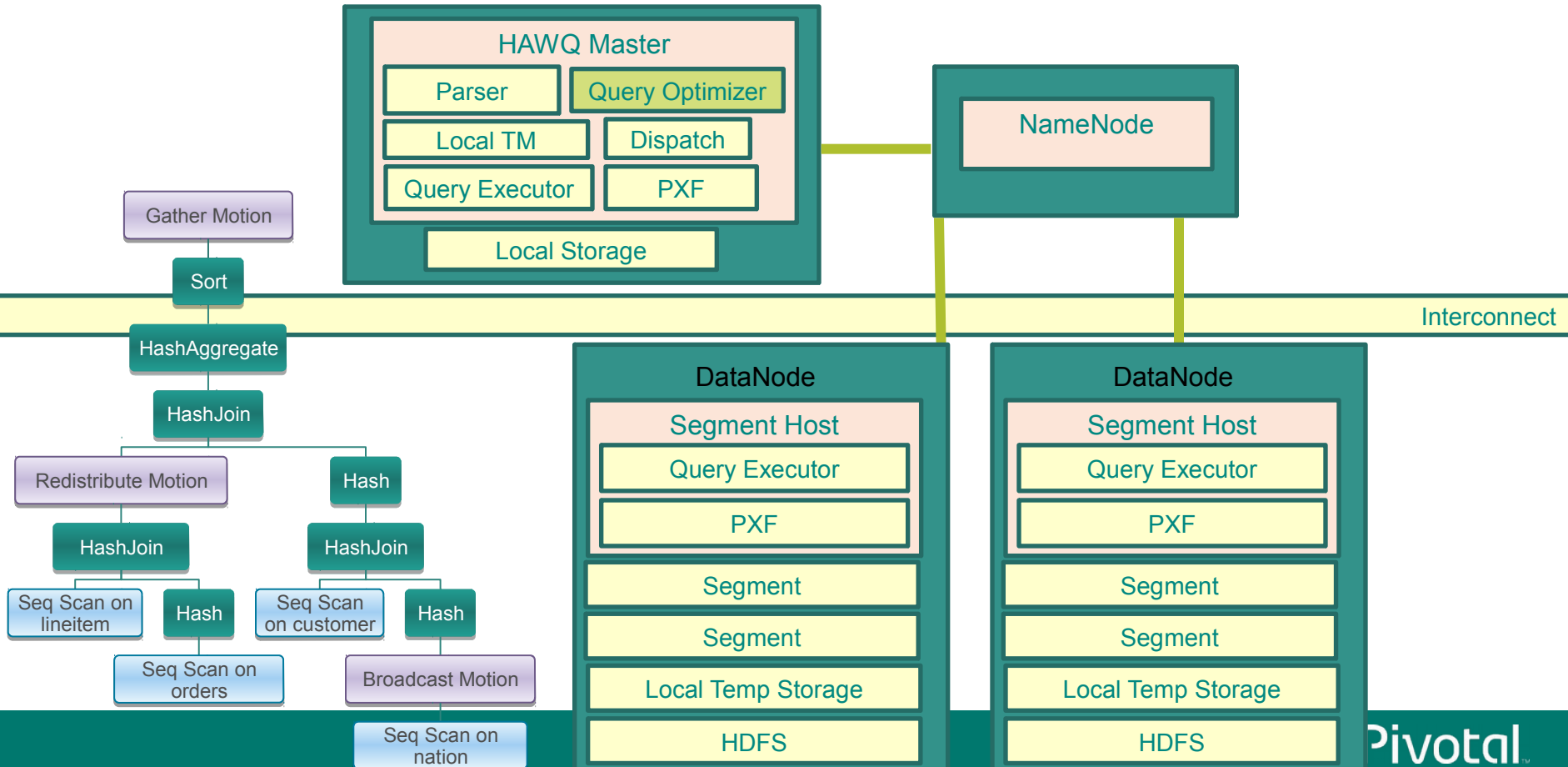
Компоненты Pivotal HAWQ

- Таблицы могут храниться как:
 - Row-oriented таблицы (сжатие quicklz, zlib)
 - Column-oriented таблицы (сжатие quicklz, zlib, rle)
 - Parquet таблицы
- Каждый сегмент имеет отдельную директорию на HDFS и отдельно хранит свой шард данных
- При поколонном хранении каждый столбец представляет собой 1 файл
- Parquet позволяет хранить таблицу по столбцам и при этом не нагружать NameNode множеством файлов

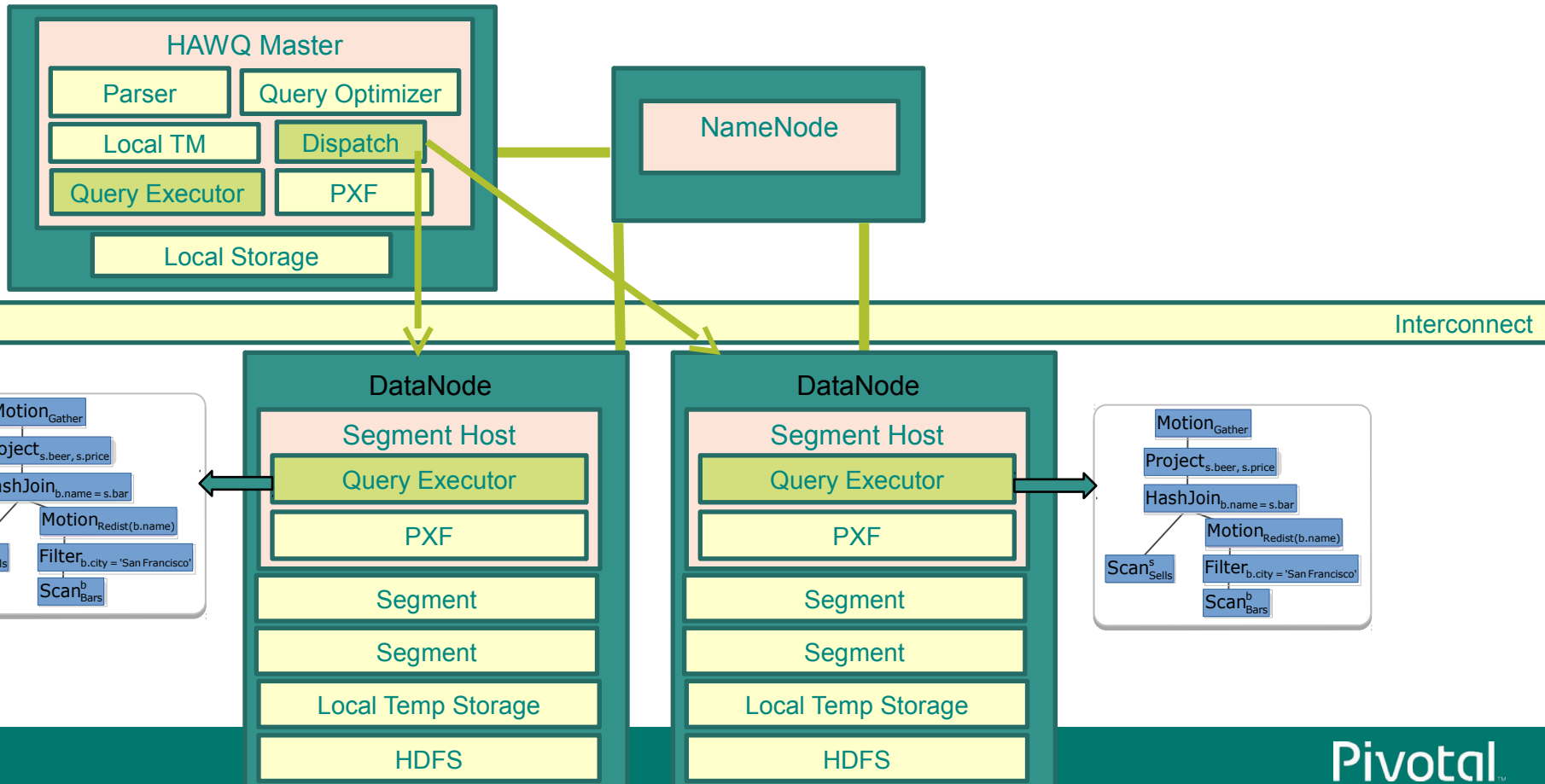
Выполнение запросов в HAWQ



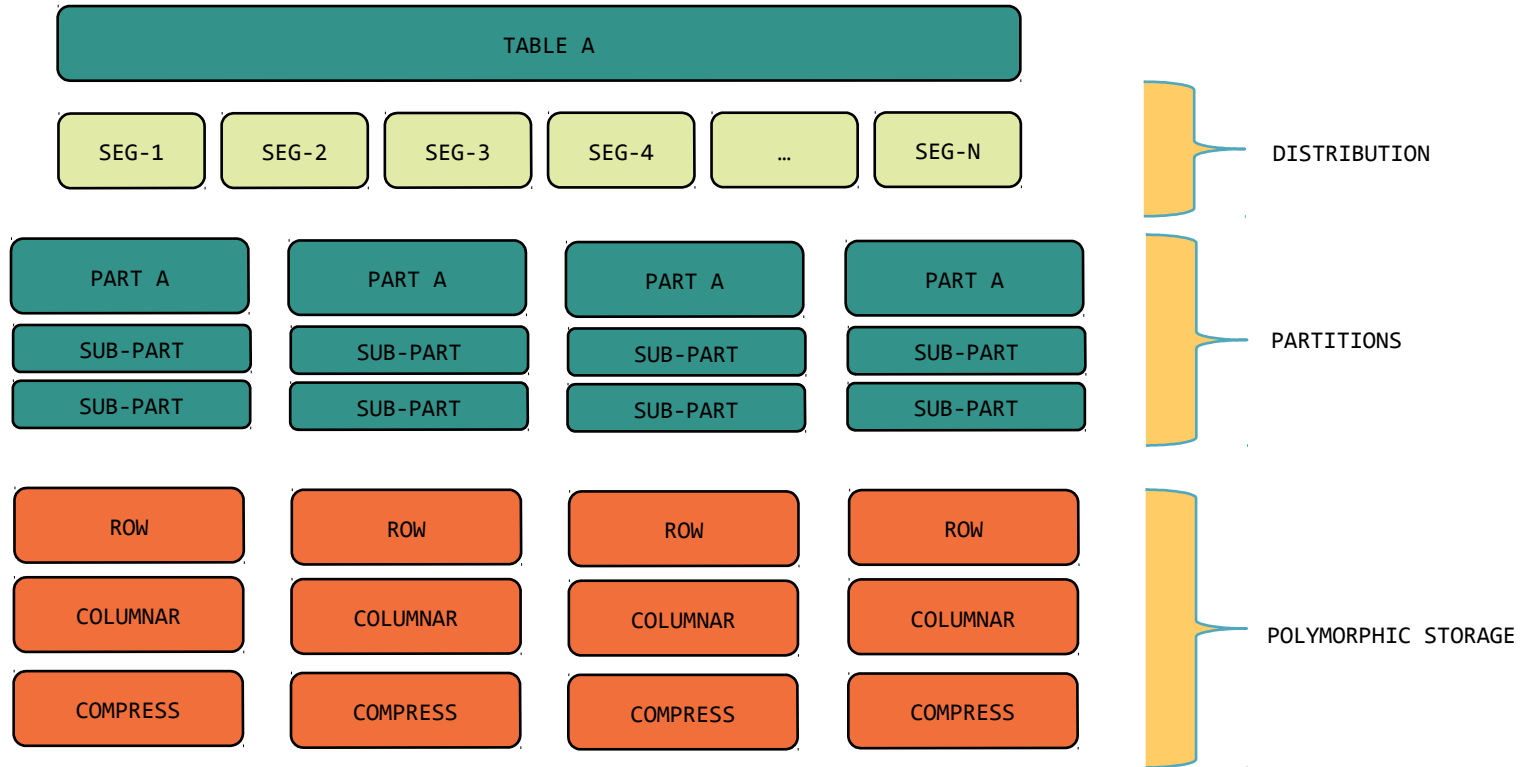
Выполнение запросов в HAWQ



Выполнение запросов в HAWQ



HAWQ Хранение данных



CORM (Исходные данные)

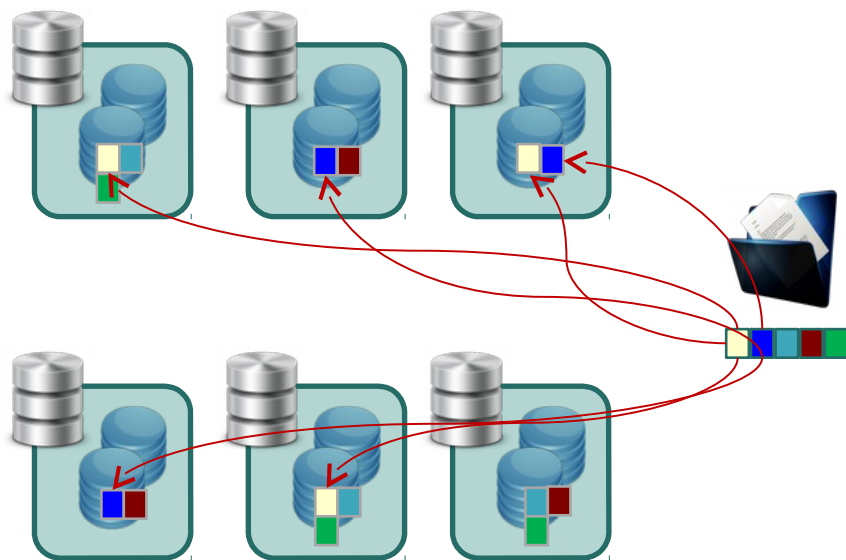
Предположим что в год требуется хранить и анализировать 200 ТВ информации (без учета сжатия);

Если принять продолжительность квартала 92 дня, в день прирост будет составлять ~ 540 GB.

Квартал	Объем данных, ТВ
Квартал 1	40
Квартал 2	45
Квартал 3	55
Квартал 4	60
Итого	200

Отказоустойчивость

3-х кратная репликация данных * (+ учет снижения производительности при наполнении более 70%)



$$200 \text{ TB} * 3 / 0.7 = 857 \text{ TB}$$

* - данный коэффициент является настраиваемым

Хранение данных #1

- Рекомендуемый объем данных на каждом их узлов не более *48 TB (4 TB диск x 12 сокетов)*.
- Для хранения 857 TB данных потребуется:
 - $857 TB / 48 TB = 18$ узлов.
- Если рассматривать текущую ситуацию, при которой объем данных за один день составляет 543 GB, то на каждом из узлов будет находиться по:
 - $540 / 18 = 30.24 GB$ (Primary Data)
- Пусть на каждом из узлов функционирует по 4 сегмента HAWQ, таким образом каждый из сегментов будет оперировать объемом:
 - $30,24 GB / 4 = 7.6 GB$ данных

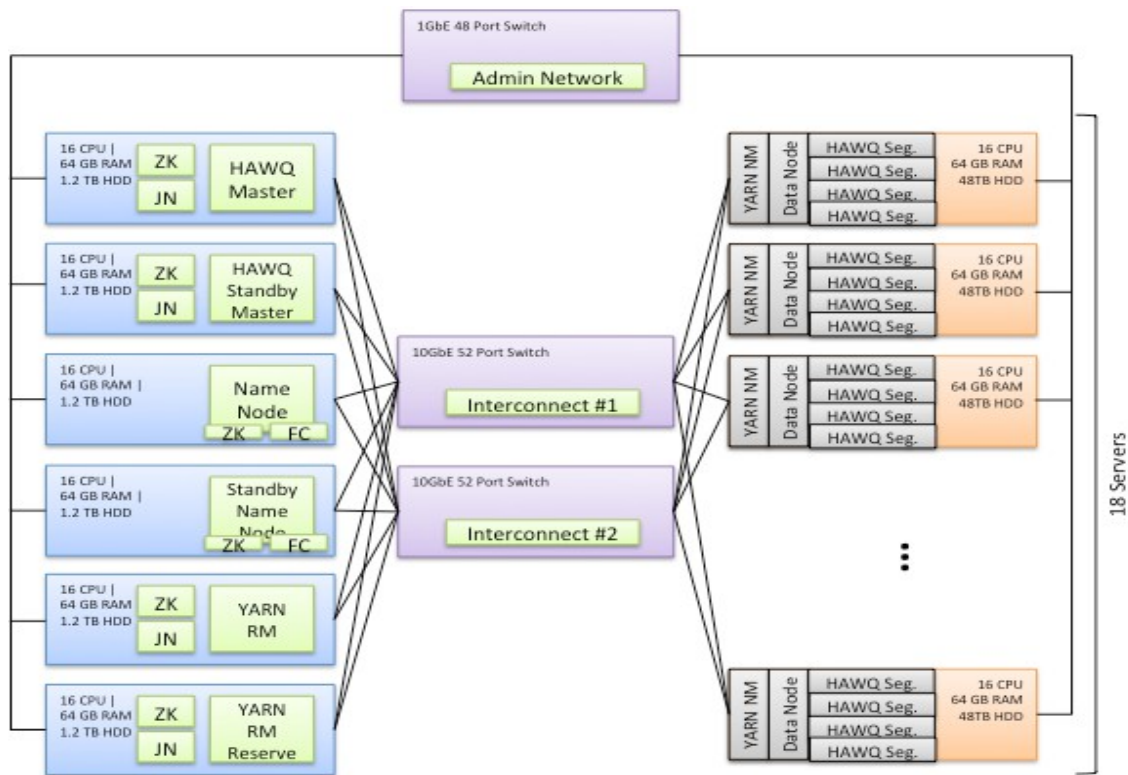
Хранение данных #2

- Предлагается разбить данные на интервалы по одному дню, а также разбить IMSI на 32 интервала по значению шести последних бит IMSI. Таким образом, каждый из дней будет разделен на 32 интервала, при этом каждому из интервалов будет соответствовать приблизительно:
 - $30.24 \text{ GB} / 32 = 900 \text{ MB}$ для узла в целом
 - $7.6 \text{ GB} / 32 = 200 \text{ MB}$ для сегмента HAWQ (файл HDFS)
- Таким образом, в файловой системе Hadoop будет находиться следующее количество файлов:
 - $18 \text{ узлов} * 4 \text{ сегм. на серв.} * 32 \text{ парт.} * 92 \text{ дня в кварт.} * 4 \text{ кварт.} = 847'872$ файлов

Производительность

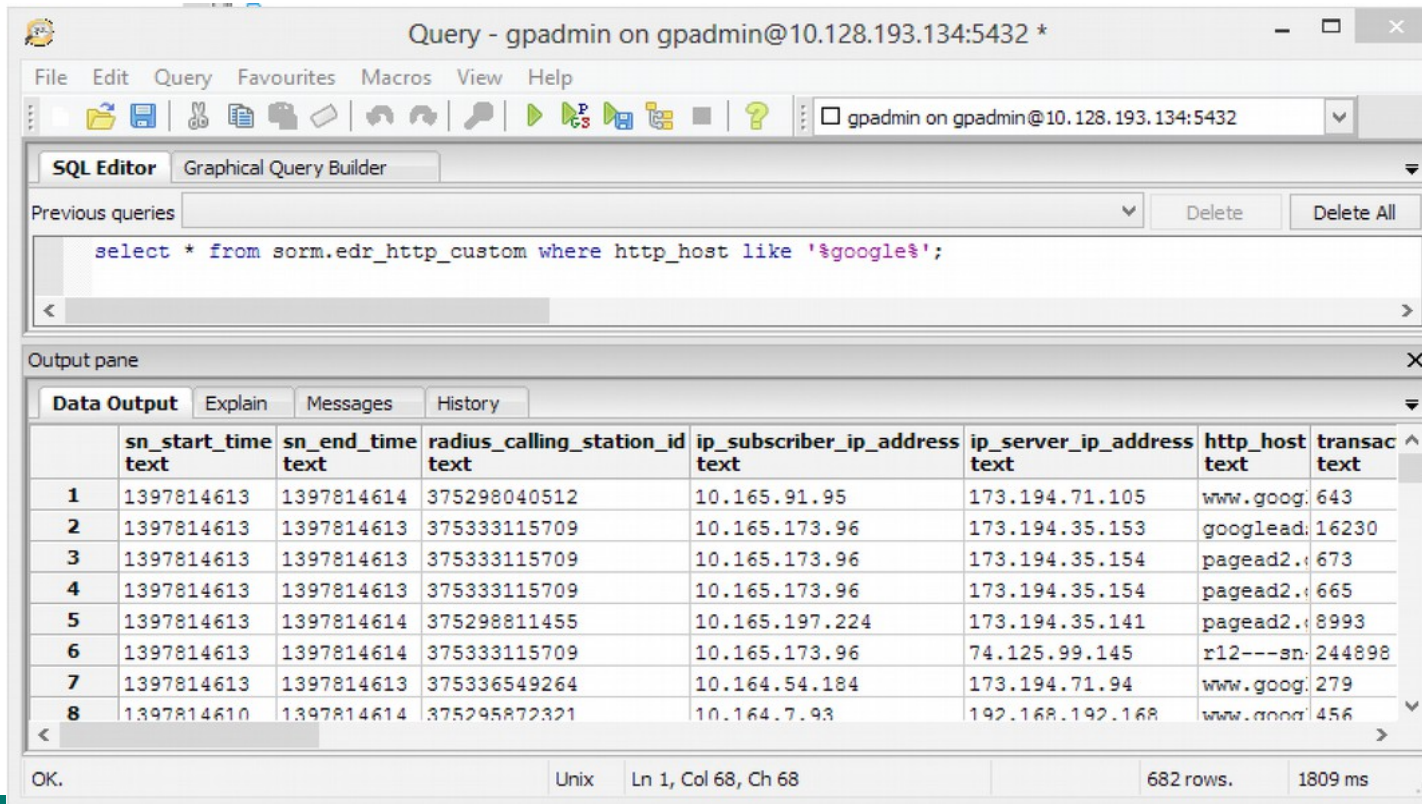
- Запрос по IMSI и периоду времени. Для анализа запроса за один день потребуется обработать 200 MB данных, что при размере блока в 64 MB составит 4 блока, то есть в среднем данные будут считываться с 4 дисков, каждый из которых сможет отдавать данные со скоростью 70 MB/сек, итого время отклика:
 - $200 \text{ MB} / 70 \text{ MBps} / 4 \text{ дисков} = 0.7 \text{ секунды}$
- При запросе данных за месяц работать будут все диски и время отклика должно составить:
 - $200 \text{ MB} * 30 \text{ дней} / 70 \text{ MBps} / 10 \text{ дисков} = 8.57 \text{ секунды}$
- По транслированному IP и периоду времени. Этим запросам свойственен короткий временной период запроса, обычно минуты. Для обработки такого запроса необходимо поднять с дисков 8.1 GB данных на каждом из узлов что приведет к задействованности всех дисков, то есть время работы составит:
 - $7600 \text{ MB} / 70 \text{ MBps} / 10 \text{ дисков} = 10.85 \text{ секунды}$
- По IP назначения + периоду времени. Обрабатывается аналогично предыдущему запросу, соответственно время отклика составит 10.85 секунды

Аппаратное обеспечение (Пример)



- Общий объем 860 ТБ (с учетом резервирования и 30% свободного пространства)

Пример реализации (анализ EDR)



The screenshot shows a SQL query tool window titled "Query - gpadmin on gpadmin@10.128.193.134:5432 *". The query editor contains the following SQL statement:

```
select * from sorm.edr_http_custom where http_host like '%$google$';
```

The output pane displays the results in a table with the following columns: **sn_start_time text**, **sn_end_time text**, **radius_calling_station_id text**, **ip_subscriber_ip_address text**, **ip_server_ip_address text**, **http_host text**, and **transac text**. The table contains 8 rows of data.

	sn_start_time text	sn_end_time text	radius_calling_station_id text	ip_subscriber_ip_address text	ip_server_ip_address text	http_host text	transac text
1	1397814613	1397814614	375298040512	10.165.91.95	173.194.71.105	www.goog	643
2	1397814613	1397814613	375333115709	10.165.173.96	173.194.35.153	googlead	16230
3	1397814613	1397814613	375333115709	10.165.173.96	173.194.35.154	pagead2.	673
4	1397814613	1397814613	375333115709	10.165.173.96	173.194.35.154	pagead2.	665
5	1397814613	1397814614	375298811455	10.165.197.224	173.194.35.141	pagead2.	8993
6	1397814613	1397814614	375333115709	10.165.173.96	74.125.99.145	r12---sn	244898
7	1397814613	1397814613	375336549264	10.164.54.184	173.194.71.94	www.goog	279
8	1397814610	1397814614	375295872321	10.164.7.93	192.168.192.168	www.goon	456

At the bottom of the window, the status bar shows "OK.", "Unix", "Ln 1, Col 68, Ch 68", "682 rows.", and "1809 ms".

Вопросы?



*“It does look similar—but this one
is powered by Hadoop”*

aermakov@pivotal.io

www.pivotal.io