

BI Partner

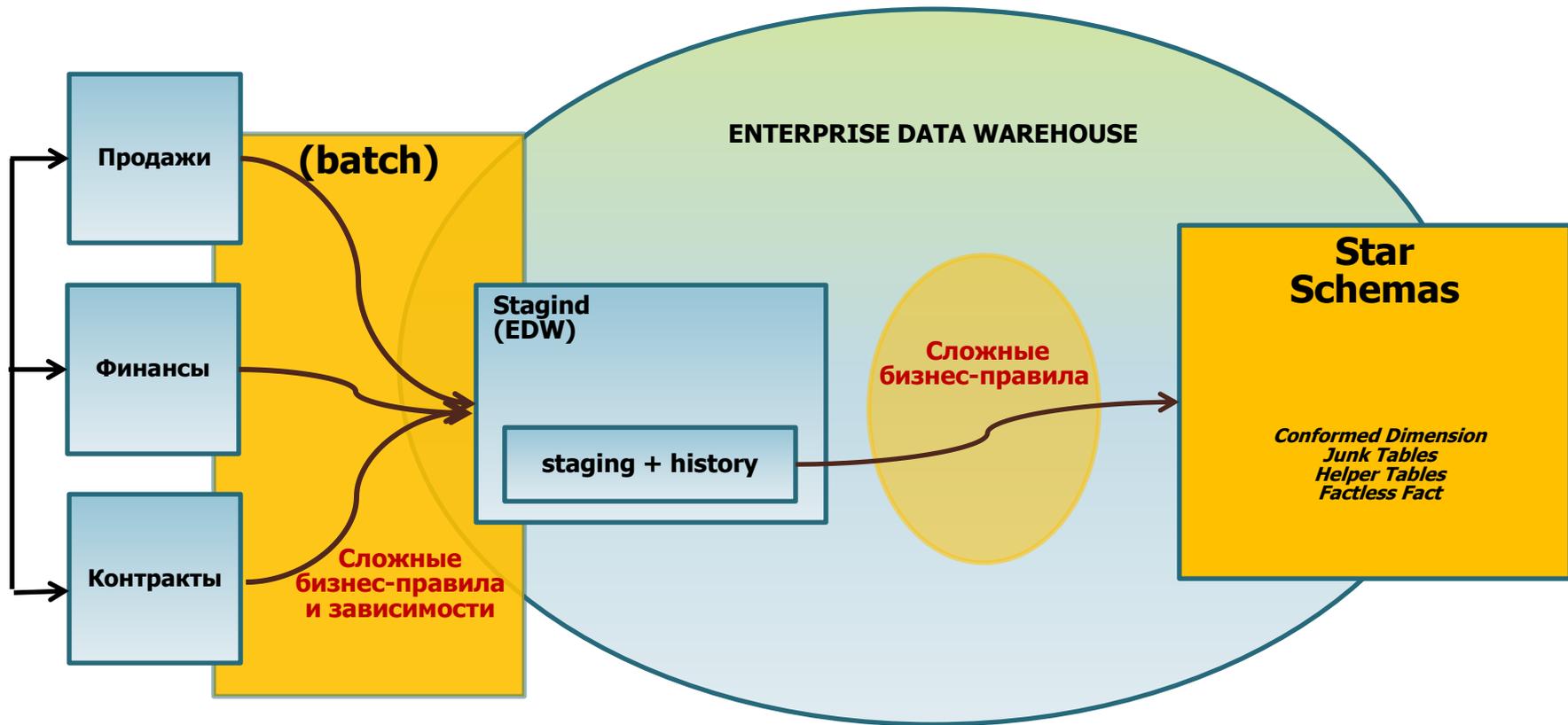
Аналитические Системы для Бизнеса

Data Vault –
новая парадигма хранилищ
данных для Big Data

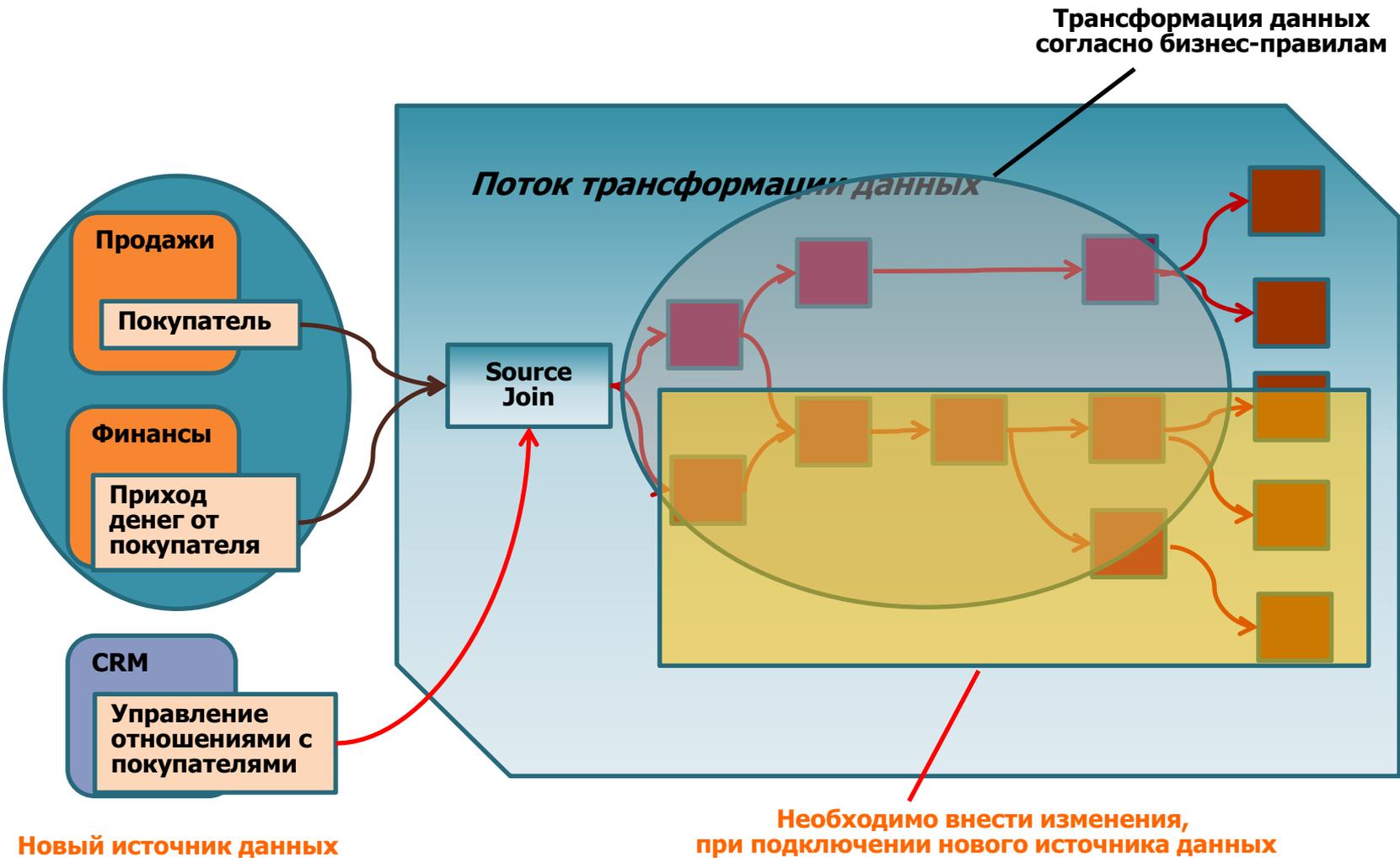
Сергей Сухарев
Руководитель BI-практики

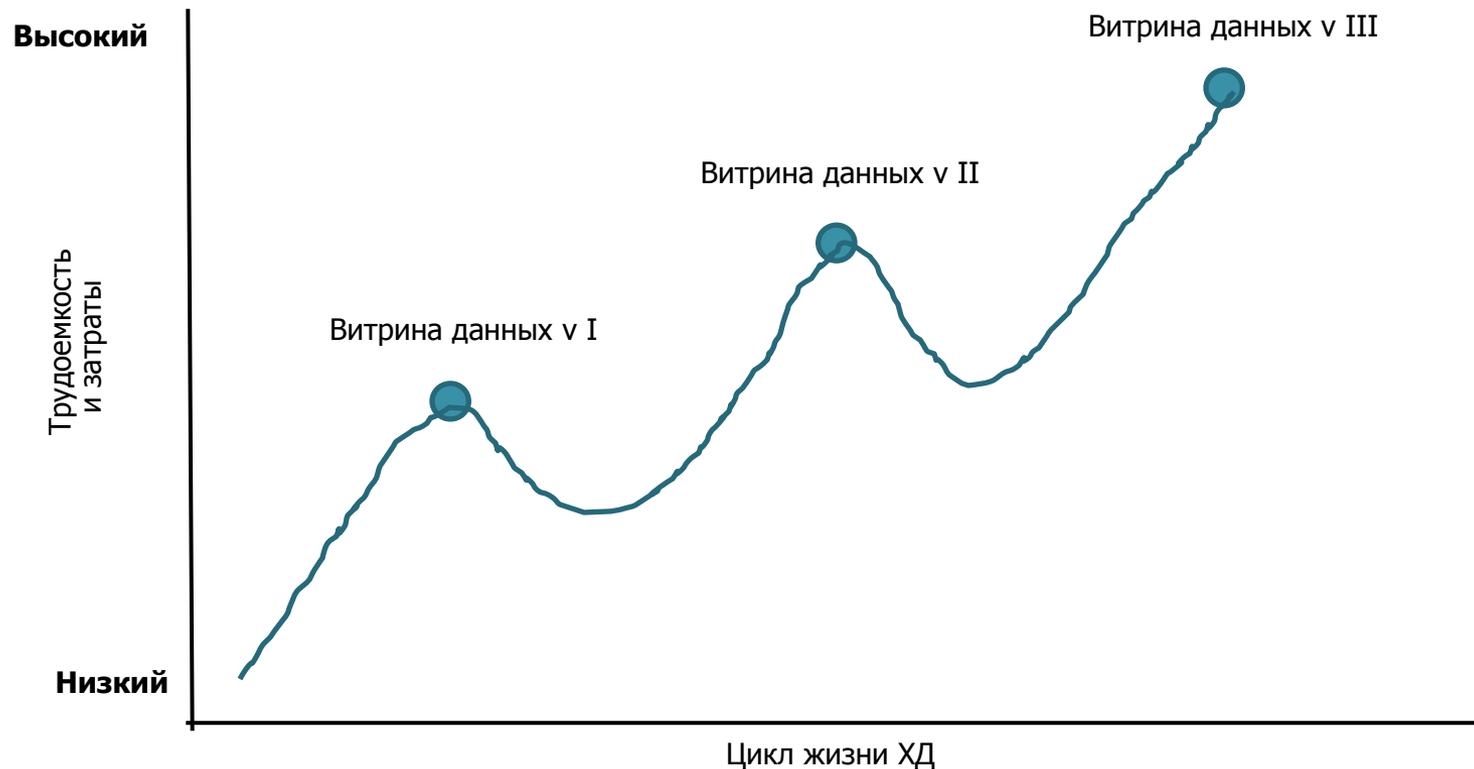
Содержание

- **Проблемы классических корпоративных хранилищ данных**
- Корпоративные хранилища данных 2.0 (DWH 2.0)
- Хранилища данных в HADOOP
- Витрины данных на продуктах Oracle и SAP
- Витрины данных HIVE + HBase (альтернатива)
- BI Partner, предложение к сотрудничеству



- Зависимость от качества данных
- Не все данные из источника представлены в ХД
- В процессе обработки в данные вносятся изменения
- Высокий риск некорректного расчета показателей
- История изменения первичных данных не может быть восстановлена по данным ХД
- Трудности при проектировании больших систем (не все показатели одинаково трактуются бизнес-пользователями)

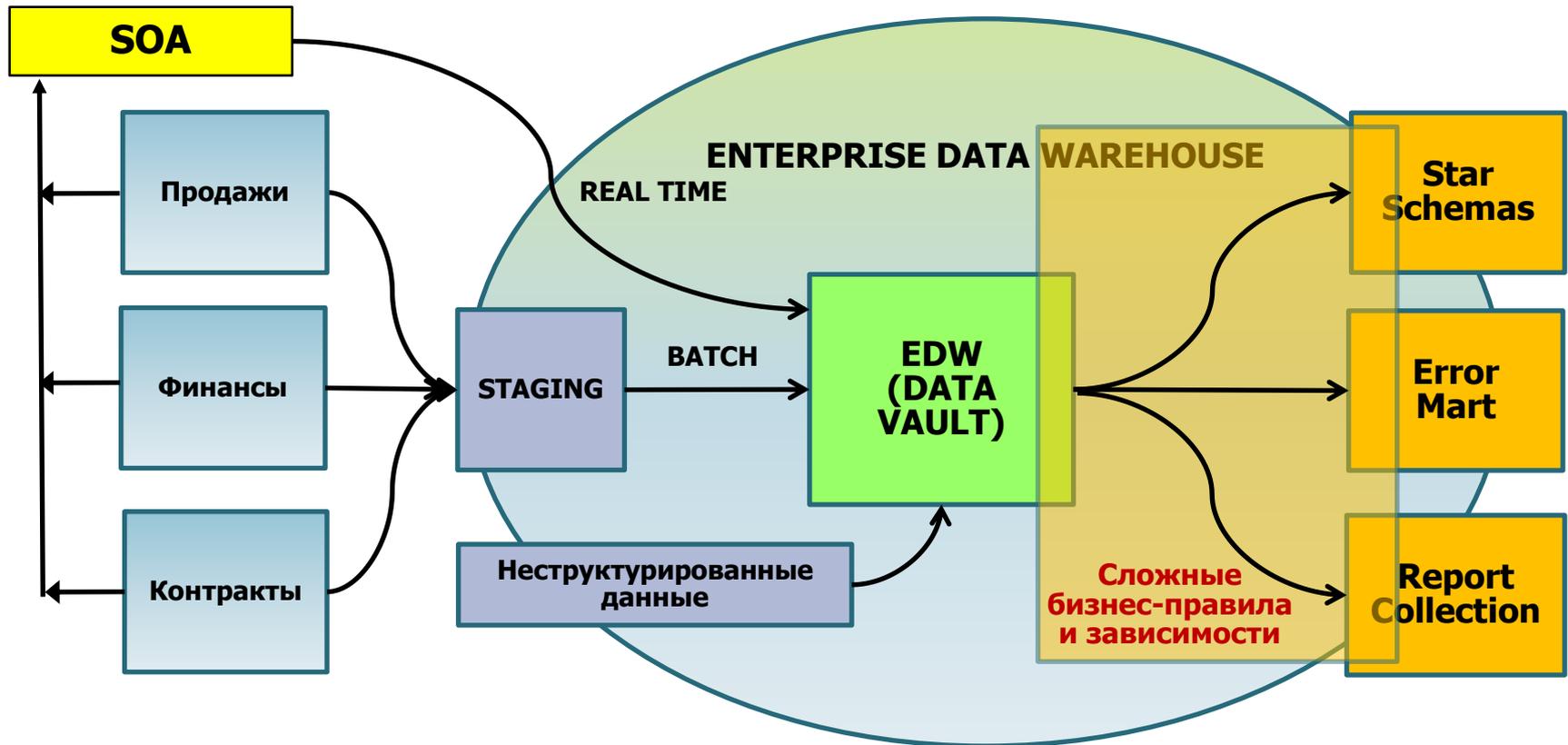




- Изменения и корректировки **согласованных витрин данных** вызывают экспоненциальный рост кривой затрат в течение всего времени жизни хранилища данных
- **Результат: каждое бизнес-подразделение строит собственную витрину данных**

Содержание

- Проблемы классических корпоративных хранилищ данных
- **Корпоративные хранилища данных 2.0 (DWH 2.0)**
- Хранилища данных в HADOOP
- Витрины данных на продуктах Oracle и SAP
- Витрины данных HIVE + HBase (альтернатива)
- BI Partner, предложение к сотрудничеству



Основные преимущества

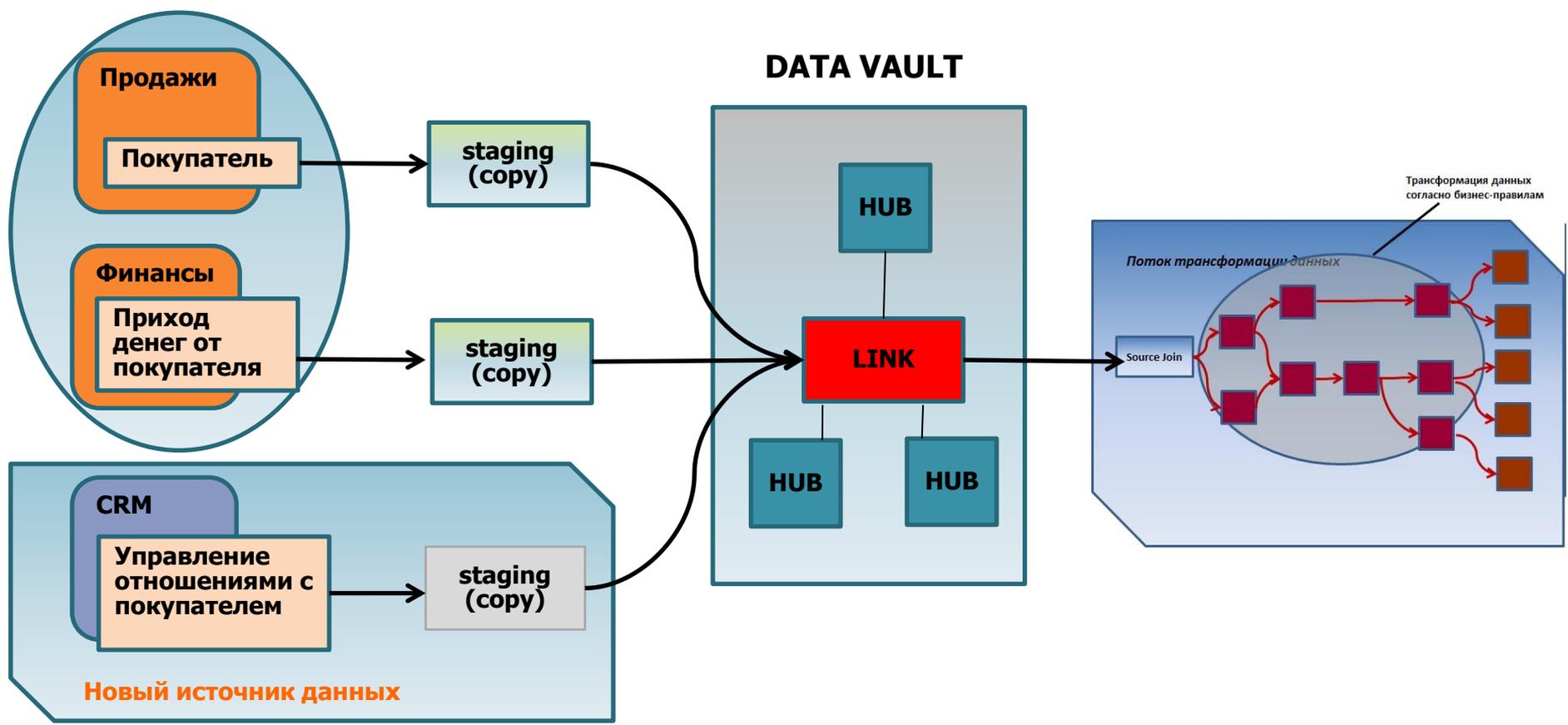
- Гибкость
- Производительность
- Отказоустойчивость
- Масштабируемость
- Легкость в поддержке и развитии
- Легкость аудита данных

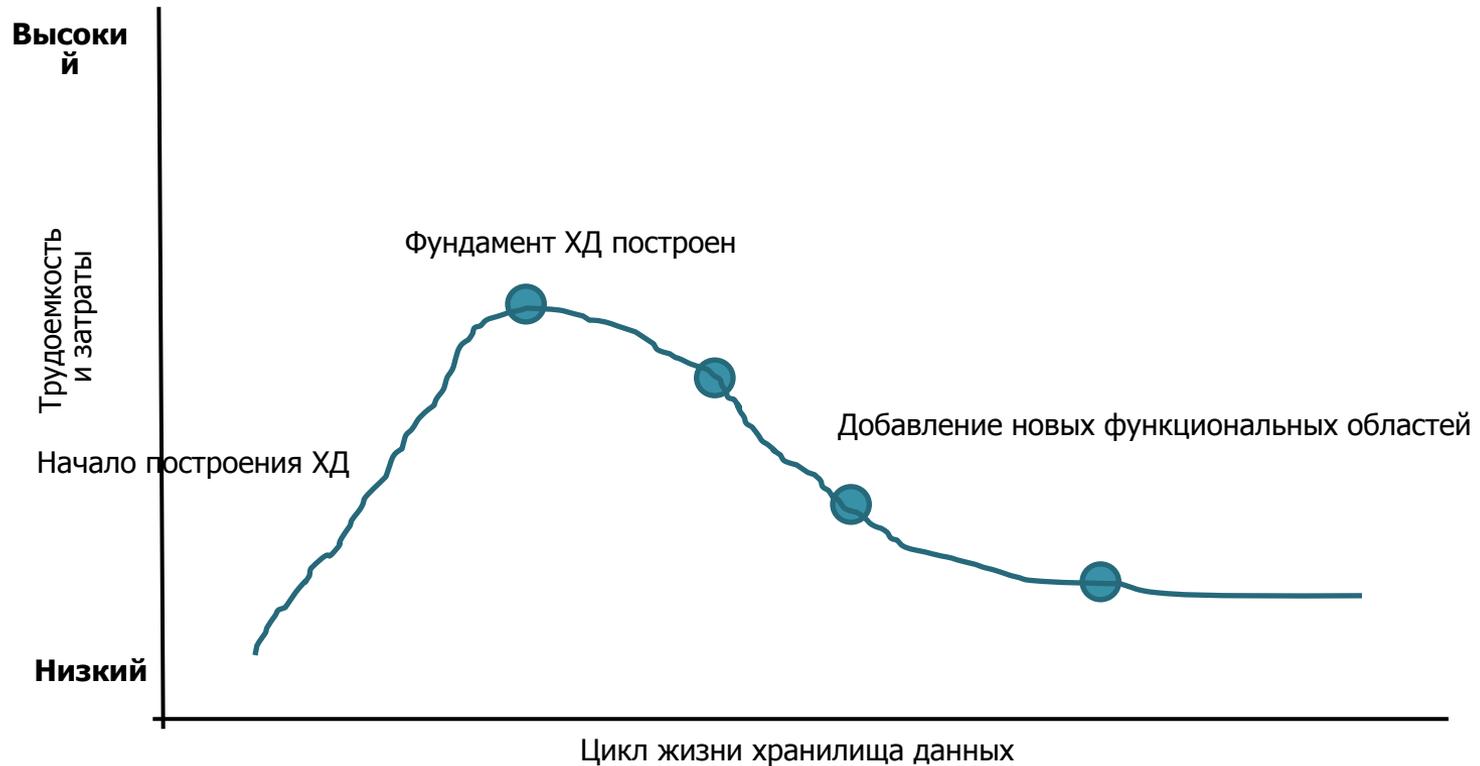
Трансформация данных, согласно бизнес-правилам перемещается ближе к бизнесу, это улучшает **время реакции, снижает стоимость и минимизирует последствия** изменений в корпоративном хранилище данных (EDW)

Согласно концепции DW 2.0 хранилища данных разделяются на два слоя:

- ◆ Слой хранения данных (собственно хранилище данных)
- ◆ Слой представления данных (витрины данных)

Методология Data Vault соответствует принципам построения DW 2.0





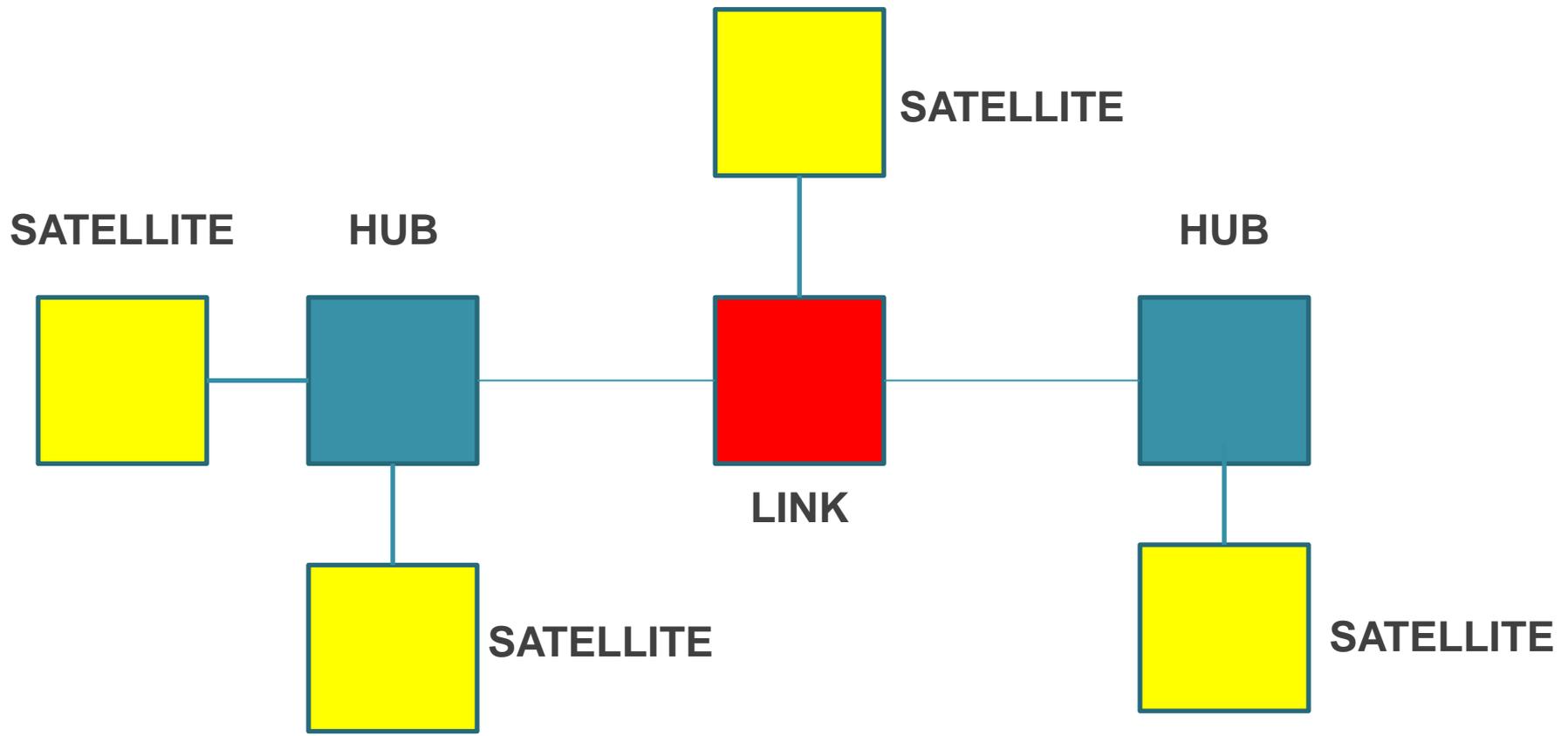
В Data Vault **не бывает** реинжиниринга

- Снижает стоимость решения и затраты на обслуживание системы с течением времени
- Упрощает процессы расширения архитектуры системы

Data Vault – это детальная, логически связанная, структура, хранящая полную историю изменения данных в источниках. Закрывает одну или несколько бизнес-областей предприятия. Гибридный подход, в котором применяется ЗНФ и схема «звезда».

Преимущество при переходе с ХД традиционной структуры (звезда или снежинка) на ХД, организованное по схеме Data Vault:

- **Гибкость.** Быстрое внесение массовых изменений в структуру ХД без переделки предыдущей структуры.
- **Масштабируемость.** Нет ограничений по увеличению размеров и масштабов ХД
- **Производительность.** Позволяет снижать стоимость и сложность процессов загрузки данных



H HUB_COUNTRY	
P	* HUB_COUNTRY_SEQ NUMBER (12)
U	* COUNTRY_ID CHAR (2)
	* HUB_Load_DTS DATE
	* HUB_Rec_SRC VARCHAR2 (12)
	👉 HUB_COUNTRY_PK (HUB_COUNTRY_SEQ)
	🔹 HUB_COUNTRY_UK1 (COUNTRY_ID)

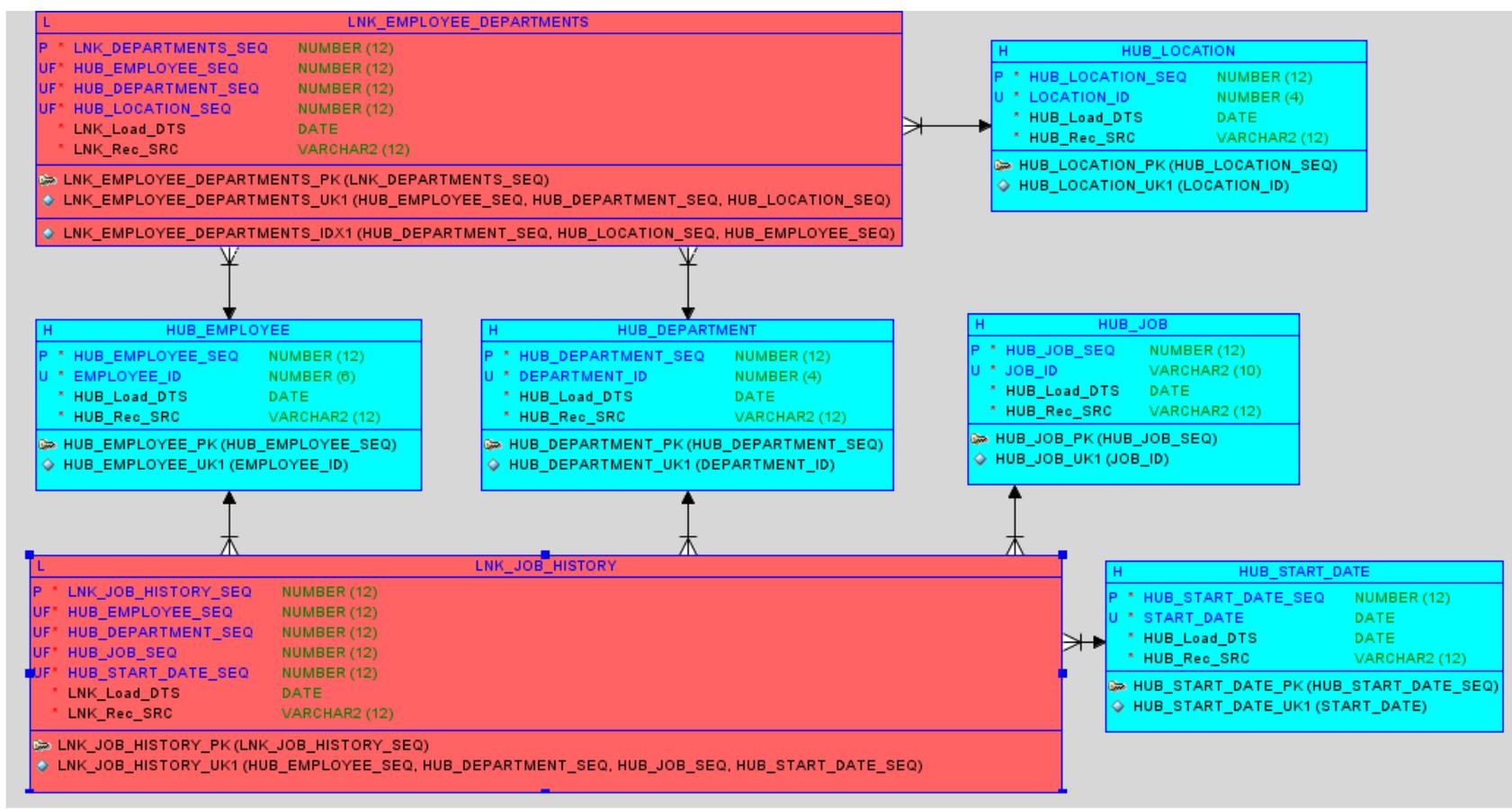
H HUB_DEPARTMENT	
P	* HUB_DEPARTMENT_SEQ NUMBER (12)
U	* DEPARTMENT_ID NUMBER (4)
	* HUB_Load_DTS DATE
	* HUB_Rec_SRC VARCHAR2 (12)
	👉 HUB_DEPARTMENT_PK (HUB_DEPARTMENT_SEQ)
	🔹 HUB_DEPARTMENT_UK1 (DEPARTMENT_ID)

H HUB_EMPLOYEE	
P	* HUB_EMPLOYEE_SEQ NUMBER (12)
U	* EMPLOYEE_ID NUMBER (6)
	* HUB_Load_DTS DATE
	* HUB_Rec_SRC VARCHAR2 (12)
	👉 HUB_EMPLOYEE_PK (HUB_EMPLOYEE_SEQ)
	🔹 HUB_EMPLOYEE_UK1 (EMPLOYEE_ID)

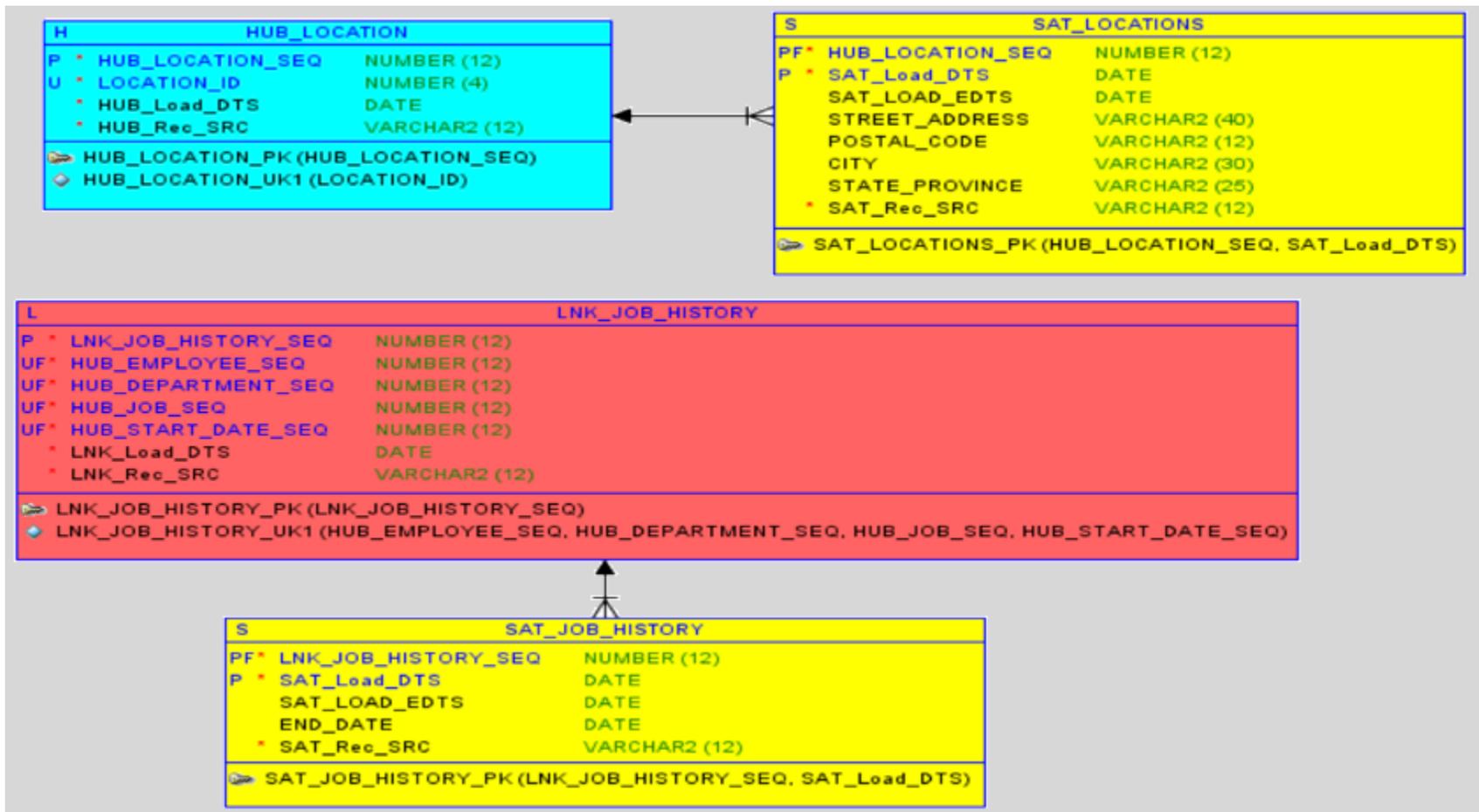
H HUB_LOCATION	
P	* HUB_LOCATION_SEQ NUMBER (12)
U	* LOCATION_ID NUMBER (4)
	* HUB_Load_DTS DATE
	* HUB_Rec_SRC VARCHAR2 (12)
	👉 HUB_LOCATION_PK (HUB_LOCATION_SEQ)
	🔹 HUB_LOCATION_UK1 (LOCATION_ID)

H HUB_JOB	
P	* HUB_JOB_SEQ NUMBER (12)
U	* JOB_ID VARCHAR2 (10)
	* HUB_Load_DTS DATE
	* HUB_Rec_SRC VARCHAR2 (12)
	👉 HUB_JOB_PK (HUB_JOB_SEQ)
	🔹 HUB_JOB_UK1 (JOB_ID)

Hub хранит уникальные бизнес-ключи сущностей



Link определяет транзакционные, ассоциативные и ролевые связи сущностей



Satellite содержит:

- Данные.
- Историю изменения данных.
- Историю изменения связей между сущностями.

Содержание

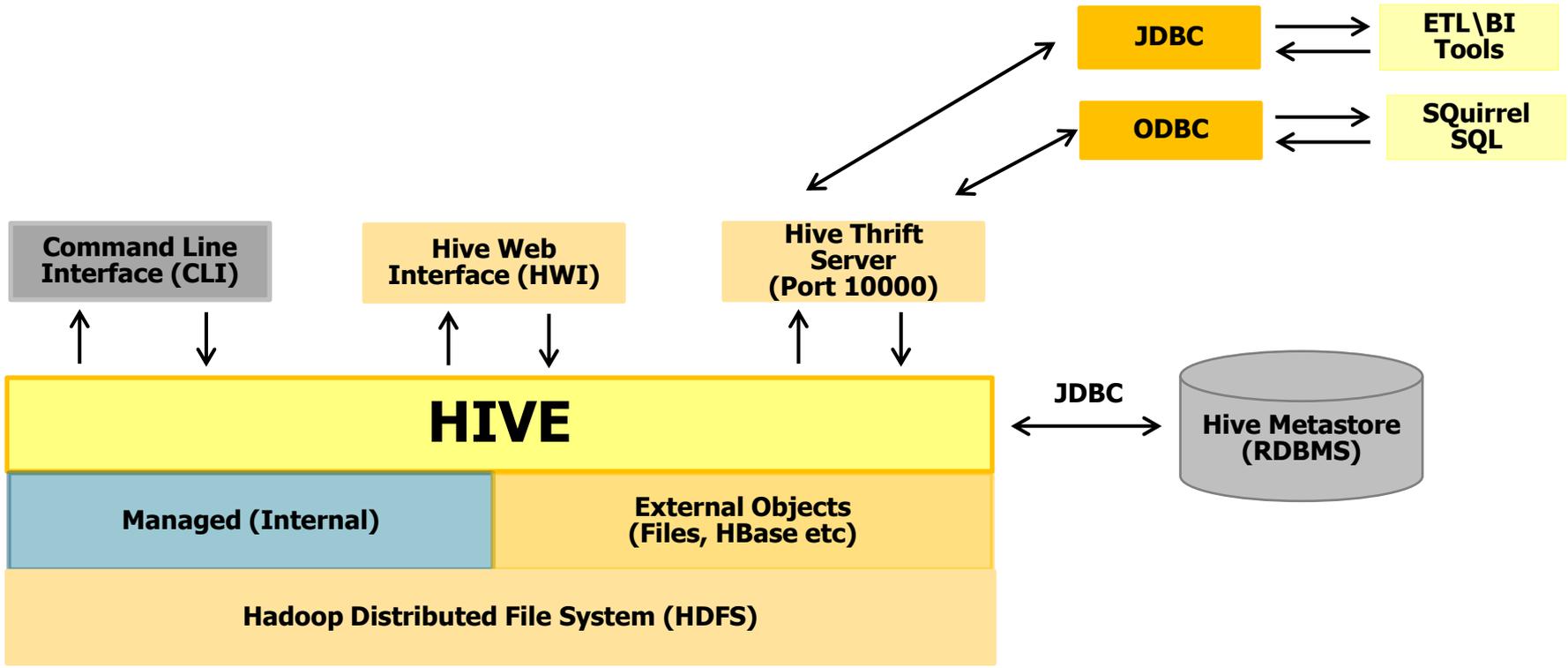
- Проблемы классических корпоративных хранилищ данных
- Корпоративные хранилища данных 2.0 (DWH 2.0)
- **Хранилища данных в HADOOP**
- Витрины данных на продуктах Oracle и SAP
- Витрины данных HIVE + HBase (альтернатива)
- BI Partner, предложение к сотрудничеству

Hadoop вносит изменения в традиционную парадигму хранилищ данных

- **Hadoop** не заменяет реляционные базы данных или платформы традиционных хранилищ данных, но его лучшее соотношение цена/производительность может помочь организациям снизить затраты, сохраняя при этом существующую инфраструктуру отчетности.
- **Hadoop** может содержать в себе все виды данных: структурированных, частично структурированных и не структурированных. Одно из важных свойств среды – низкая стоимость хранения и легкость масштабирования.
- Одним из эффективных способов использования технологий **Big Data** является применение этой технологии в организации корпоративных хранилищ данных. В этом случае организуется зона первичных данных, содержащая полную историю изменения данных в источниках (идеально подходит для ХД структуры **Data Vault**).

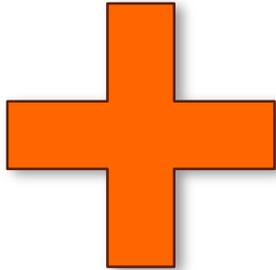


Apache Hive – хранилище данных для Hadoop



- **HiveQL** – SQL подобный язык запросов к данным.
- Разработан для работы с данными большого объема (PB).
- Предназначен для анализа данных в т.ч. с использованием ad-hoc запросов.
- Умеет работать с DDL операторами: `database \ table \ partition \ bucket`.
- Поддерживает стандартные и комплексные типы данных.
- Легко расширяется.
- Не предназначен для небольших баз данных и OLTP систем.





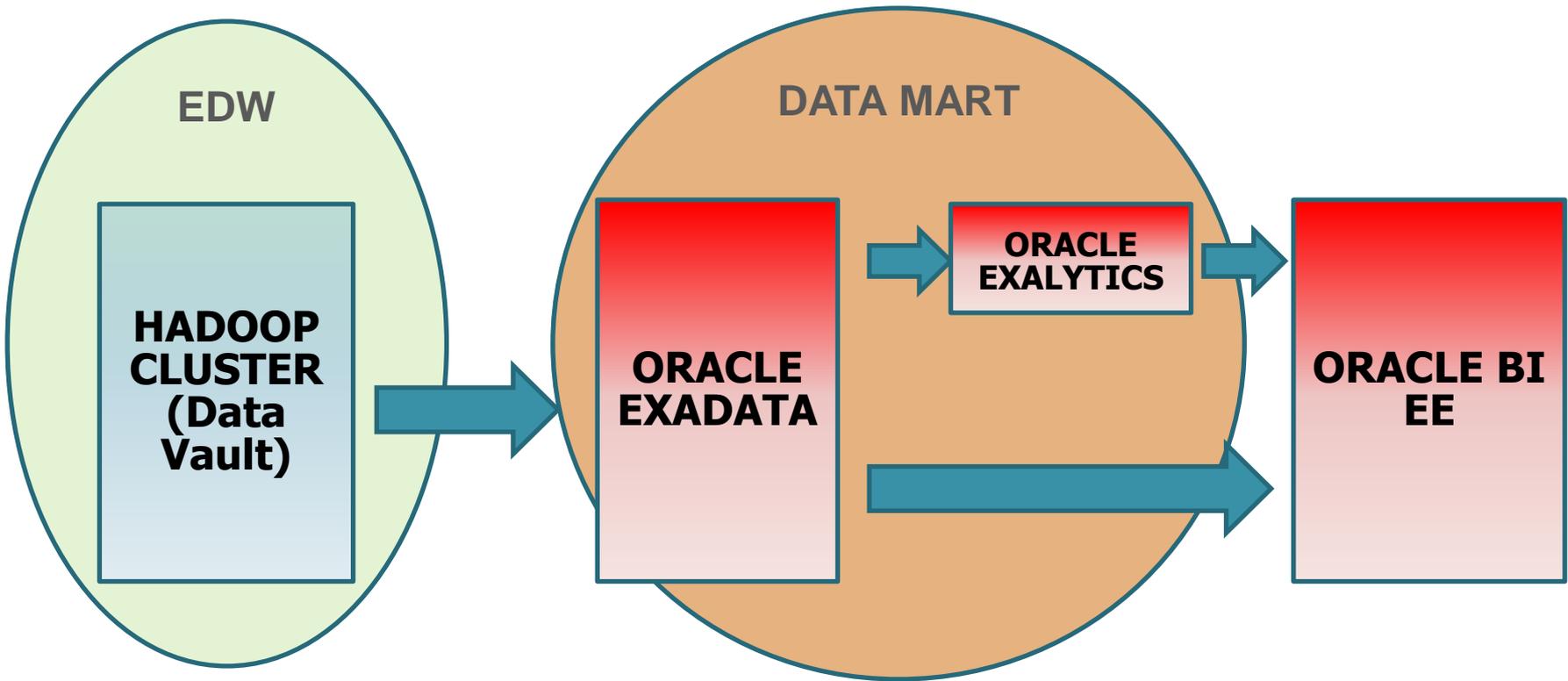
- Централизация данных
- MPP и MAP/REDUCE
- Низкая стоимость хранения данных
- Можно использовать сегментирование данных (partitioning)
- Можно использовать сжатие данных
- Можно выбирать формат хранения данных
- Гибкость

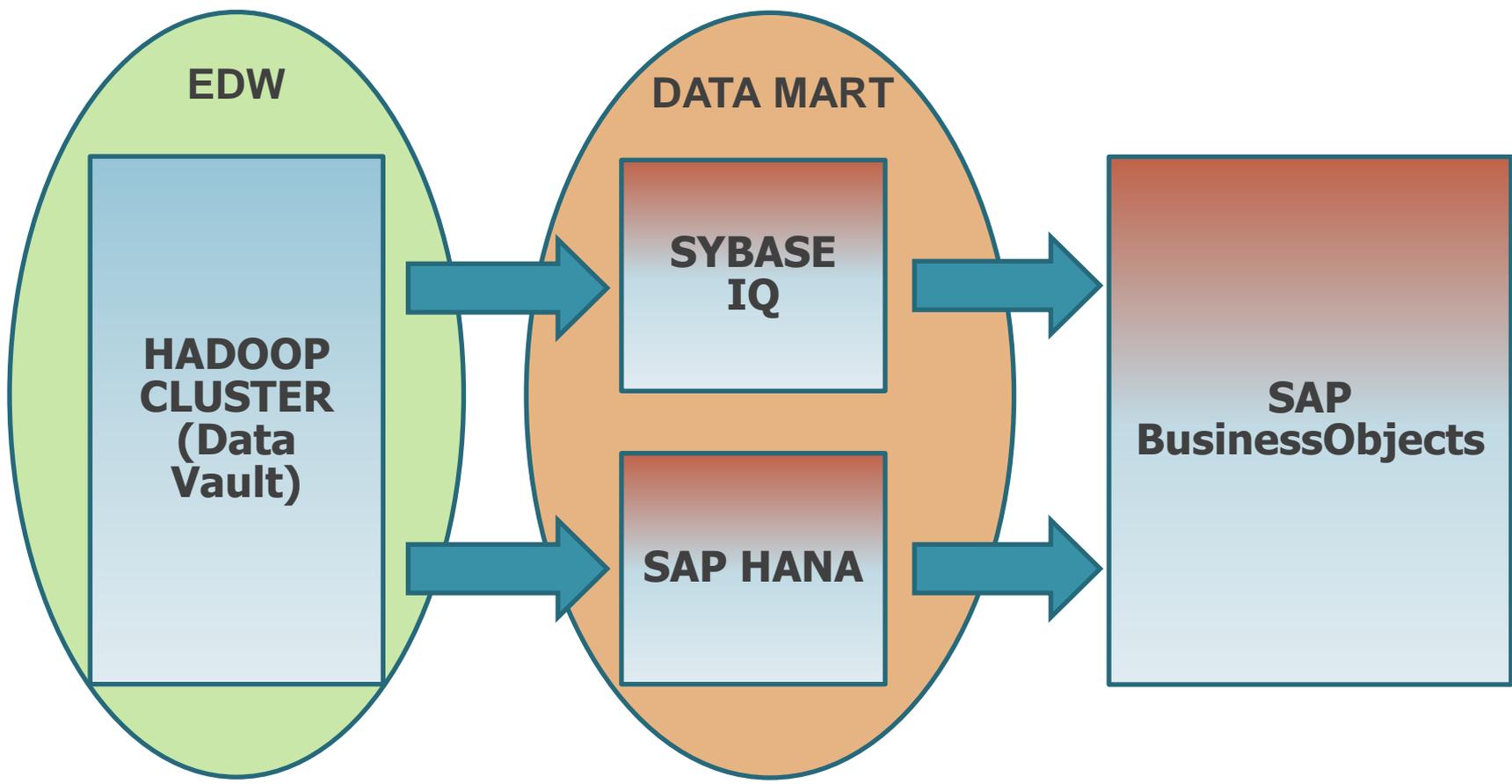


- Файловое управление данными
- Данные нельзя обновлять (update) и удалять (delete)
- Для разработки приложений необходимы навыки высококлассного программиста
- Трудно оптимизировать

Содержание

- Проблемы классических корпоративных хранилищ данных
- Корпоративные хранилища данных 2.0 (DWH 2.0)
- Хранилища данных в HADOOP
- **Витрины данных на продуктах Oracle и SAP**
- Витрины данных HIVE + HBase (альтернатива)
- BI Partner, предложение к сотрудничеству



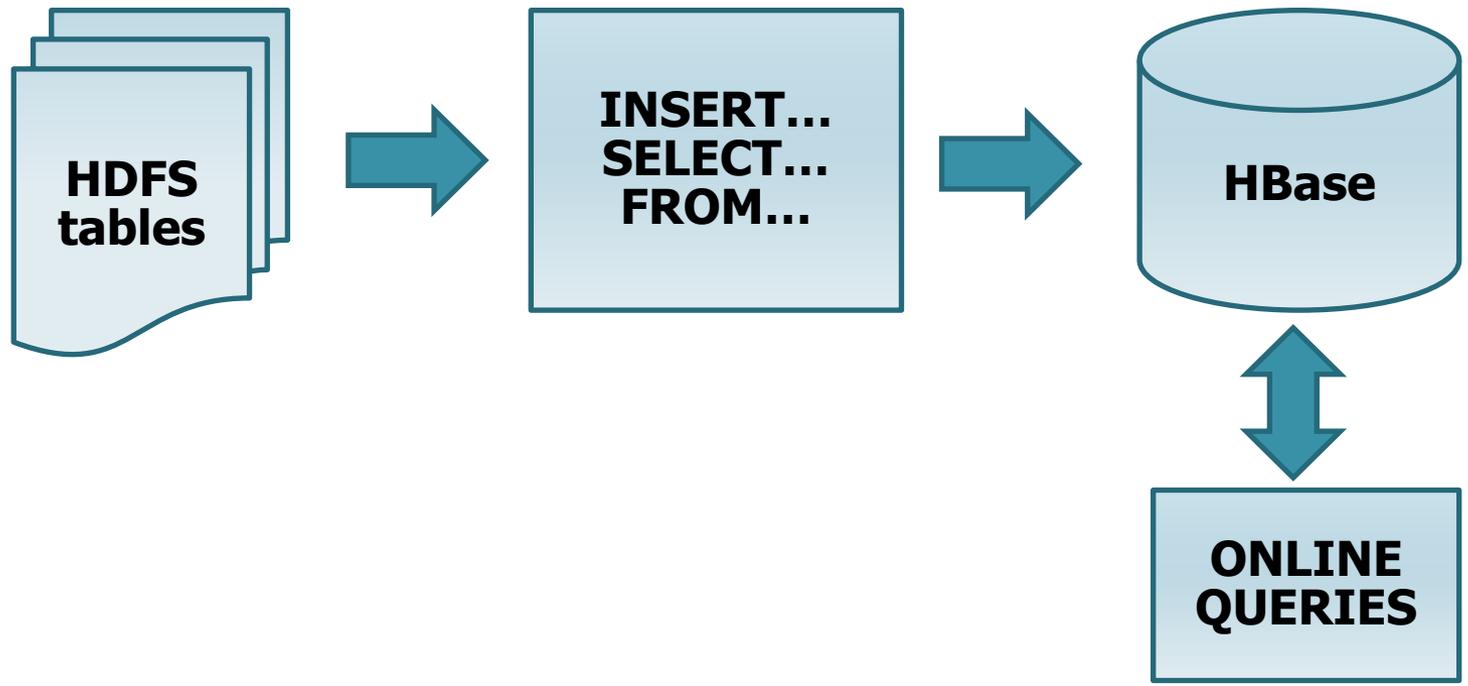


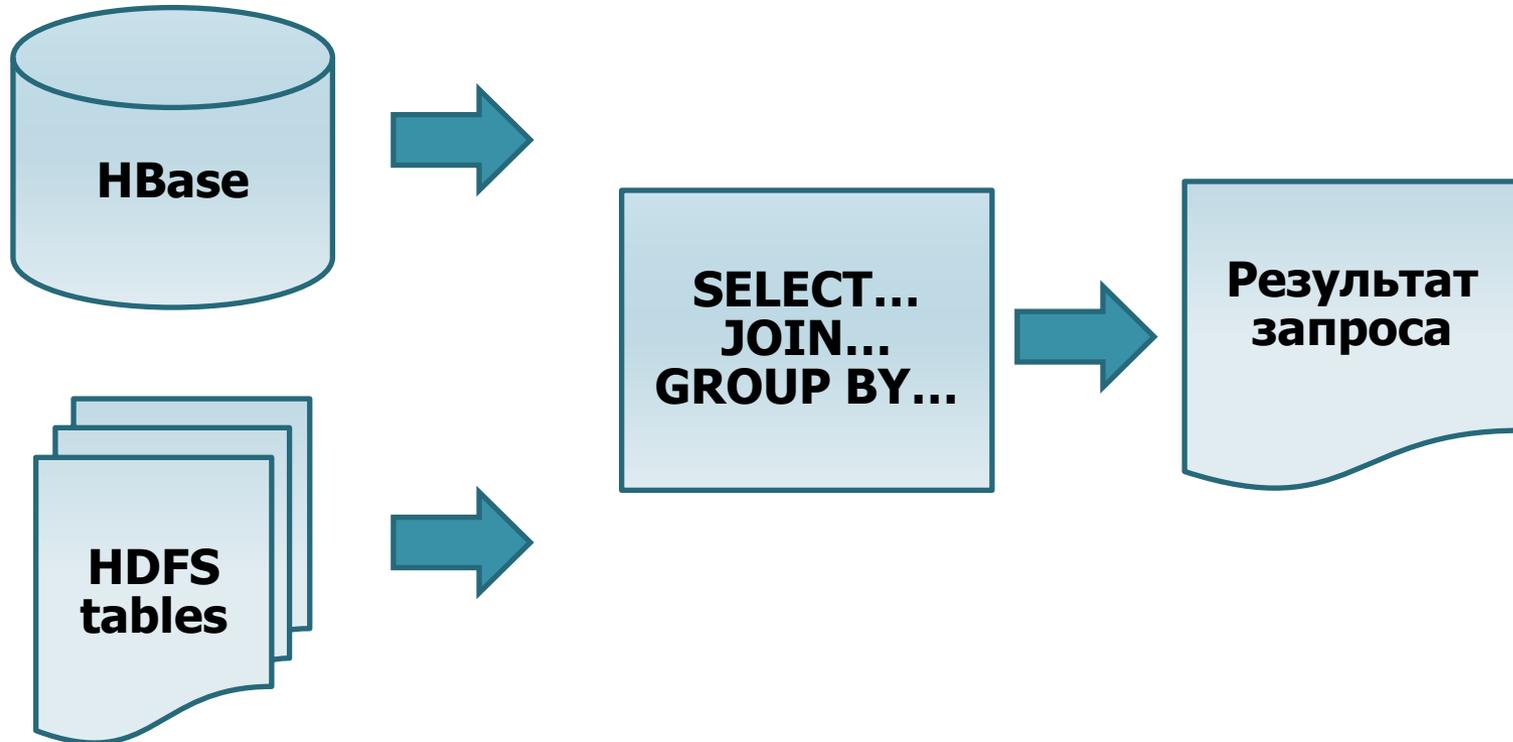
Содержание

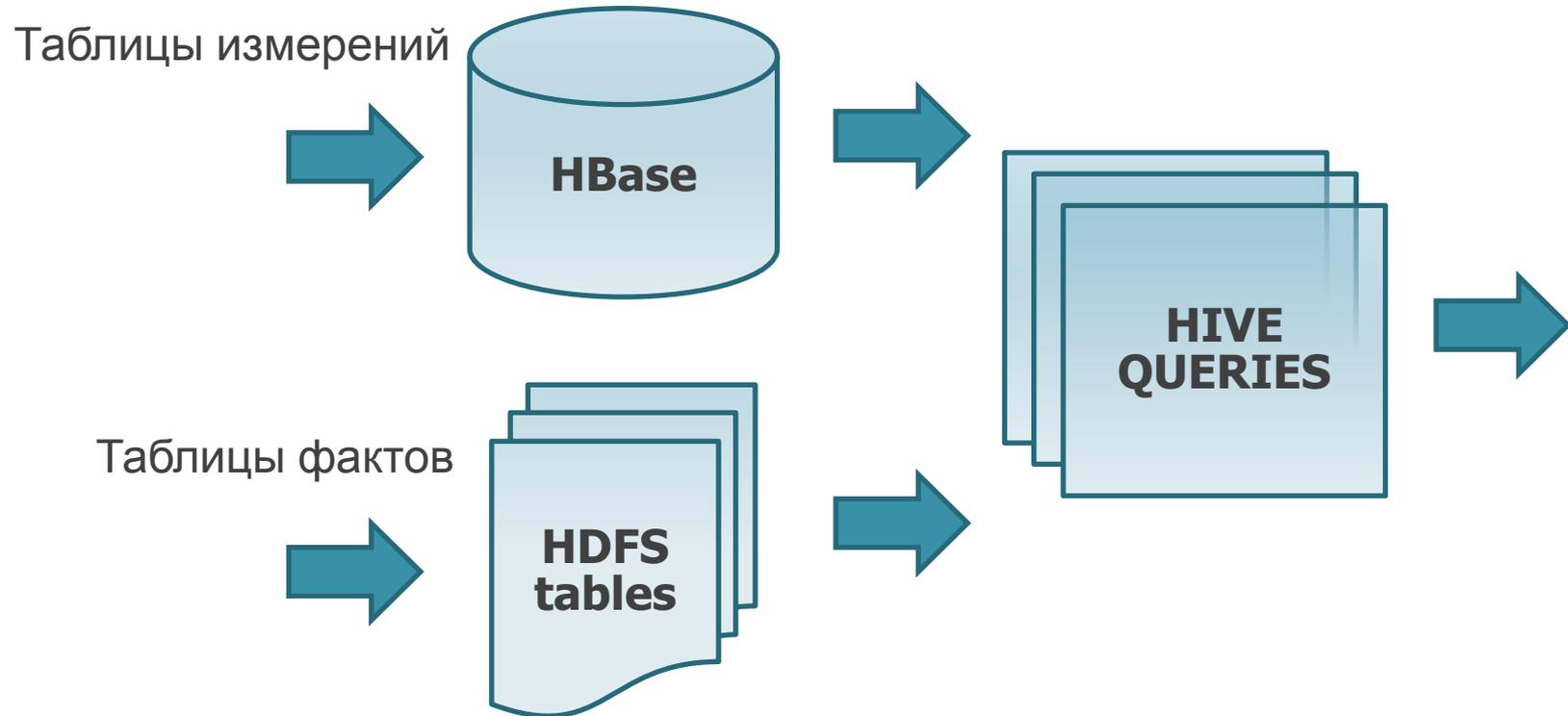
- Проблемы классических корпоративных хранилищ данных
- Корпоративные хранилища данных 2.0 (DWH 2.0)
- Хранилища данных в HADOOP
- Витрины данных на продуктах Oracle и SAP
- **Витрины данных HIVE + HBase (альтернатива)**
- BI Partner, предложение к сотрудничеству

- **Apache™ Hbase** - не реляционная (NoSQL) база данных, которая работает поверх Distributed File System (HDFS) в Hadoop.
- Колоночная база данных, которая обеспечивает отказоустойчивость хранения и быстрый доступ к большим объемам разреженных данных.
- Добавляет транзакционные возможности для Hadoop, что позволяет пользователям проводить обновление, вставку и удаление данных.
- Поддерживает автоматическую балансировку нагрузки.
- Поддерживает кэширование в памяти.
- Поддерживает репликацию через центр обработки данных.





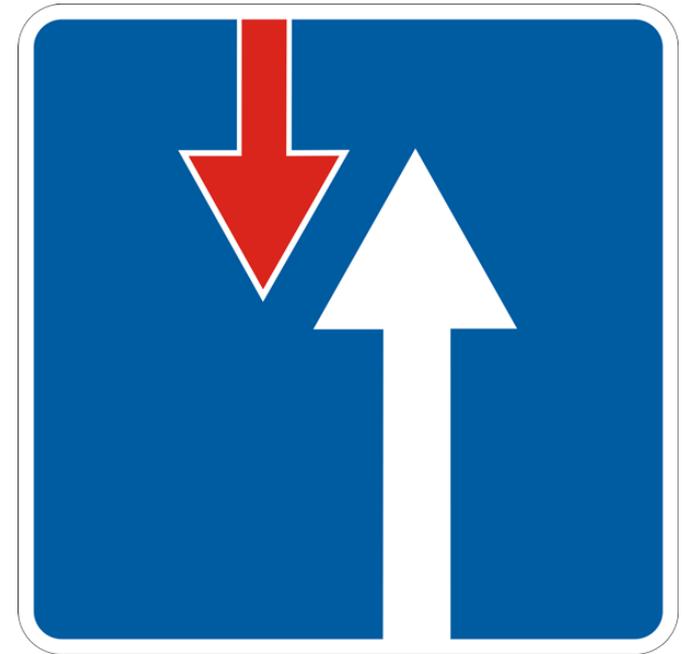




- Для input\output format, getSplits() и т.д. используются базовые классы HBase
- Функции выбора столбца и некоторые фильтры могут быть перенесены вниз к HDFS
- Таблицы HBase могут быть использованы с другими Hadoop-объектами в SQL конструкциях
- DDL операторы Hive преобразуются в DDL операторы HBase



- **Непомерно высокая стоимость масштабирования традиционных ХД** приводит к экономически неоправданным затратам, и даже если это масштабирование достигнуто, производительность традиционных систем не позволяет производить обработку большого объема данных.
- Применение технологий Big Data **многократно снижает стоимость аппаратных средств** для организации ХД, по сравнению с традиционными СУБД на аналогичных объемах данных.



Содержание

- Проблемы классических корпоративных хранилищ данных
- Корпоративные хранилища данных 2.0 (DWH 2.0)
- Хранилища данных в HADOOP
- Витрины данных на продуктах Oracle и SAP
- Витрины данных HIVE + HBase (альтернатива)
- **BI Partner, предложение к сотрудничеству**

- Год основания: **2002**
- **Первая** в России консалтинговая компания, специализирующаяся на бизнес-аналитике
- Более **100 проектов** в области ХД и ВІ
- Партнер **SAP** по системам управления эффективностью бизнеса
 - SAP Reseller
 - SAP PCoE (Partner Centre of Expertise)
 - SAP Preferred Training Partner
 - SAP Subscription-based Hosting Program Provider
- Партнер **Oracle** по базовым технологиям и хранилищам данных
 - Oracle Platinum Partner
 - Специализация **Datawarehousing**
 - Специализация **Oracle Database Performance Tuning**
 - Специализация **Oracle Database**
- Партнер **IBM** по бизнес-аналитике и бюджетному управлению
 - IBM Premium Business Partner





- Реализовать пилотный проект по построению ХД на технологиях **Big Data**
- Помочь вам развернуть кластер среднего уровня – до 10 узлов на платформе **Hortonworks** или **Cloudera**



www.bipartner.ru

consulting@bipartner.ru

Сергей Сухарев, руководитель BI-практики