

KORE DWH

PIMON 2024



+ | T I I I

09.24

Описание продукта



Набор методик и инструментов для автоматизации построения хранилища данных.



Архитектурный дизайн системы



Стандарты моделирования проектирования и разработки

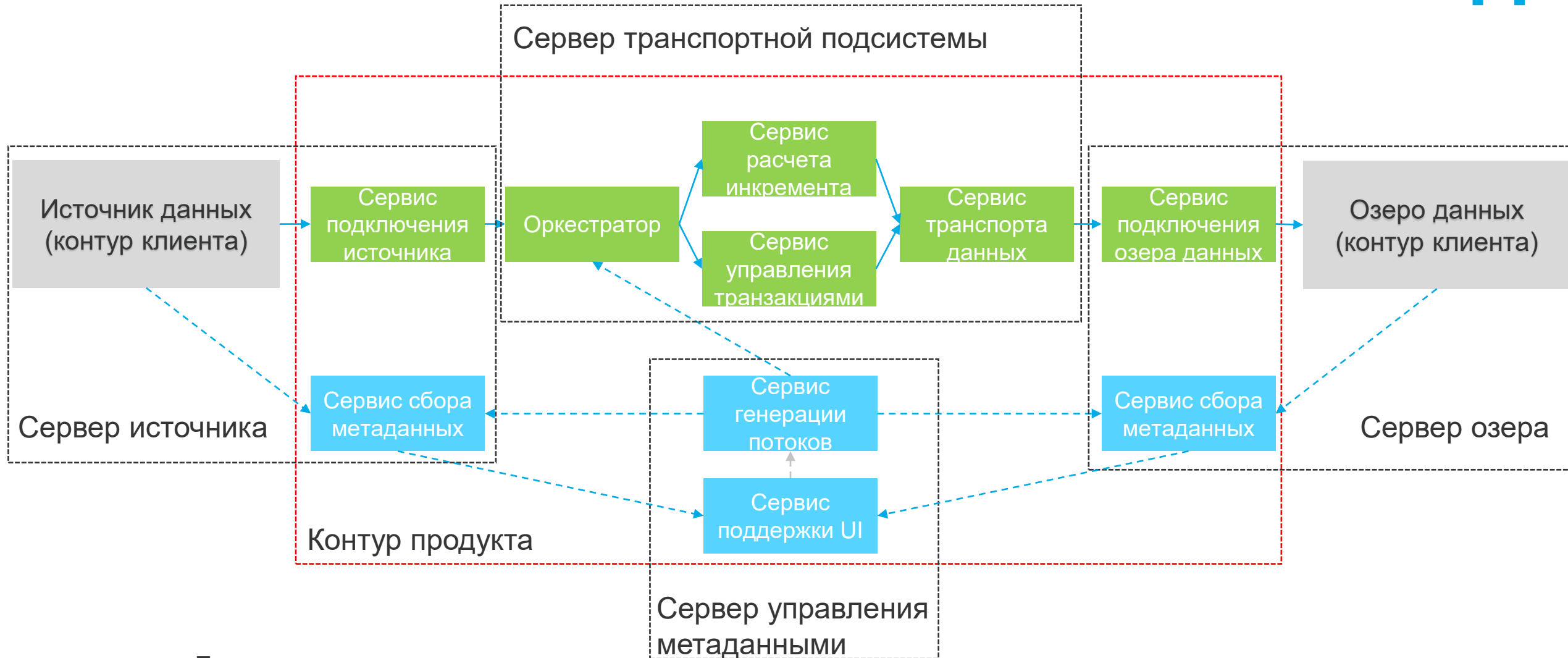


Готовые инструменты и технологии

Преимущества готового решения

- + Базис для создания децентрализованной платформы данных
- + Повышение гибкости и скорости для внесения изменений
- + Увеличение управляемости нагрузкой на источники и получатели данных
- + Устранение «бутылочных горлышек» в процессе обеспечения данными потребителей
- + Объединение усилий всех участников процесса предоставления данных

Архитектура. Поток данных



Поток данных

Поток метаданных

Организация и выполнение интеграционных потоков

Работа с ресурсами

Логическая сущность системы, с которой выполняет работу поток. Каждая система представляется набором ресурсов, для управления разрабатывается провайдер, реализующий набор API по спецификации

Работа с транзакциями

Утилита для создания и управления распределенными транзакциями загрузки данных. Менеджер транзакций взаимодействует со всеми провайдерами ресурсов, участвующих в транзакции, по выделенным API

Работа с инкрементами

Используется для регистрации изменений в областях данных, с которыми работают провайдеры ресурсов. Менеджер дельты регистрирует все операции над ресурсом и создает версии ресурса, обслуживаемых конкретным провайдером

Стандартные преобразования

Большое количество подготовленных функций и операторов, решающих задачи преобразования данных, работы с различными алгоритмами версионирования и историчности данных, возможность добавлять собственные функции

Интеграция на основе метаданных



Основной принцип построения интеграции с использованием платформы KORE – максимальное использование автоматизировано обрабатываемых метаданных для описания входов, выходов, настроек и правил трансформации

KORE DWH интегрирован с инструментом управления метаданными, который выступает в роли пользовательского интерфейса

Цель – снизить количество ручных операций по поиску и описанию данных и сократить расстояние от поиска данных до тестирования потока

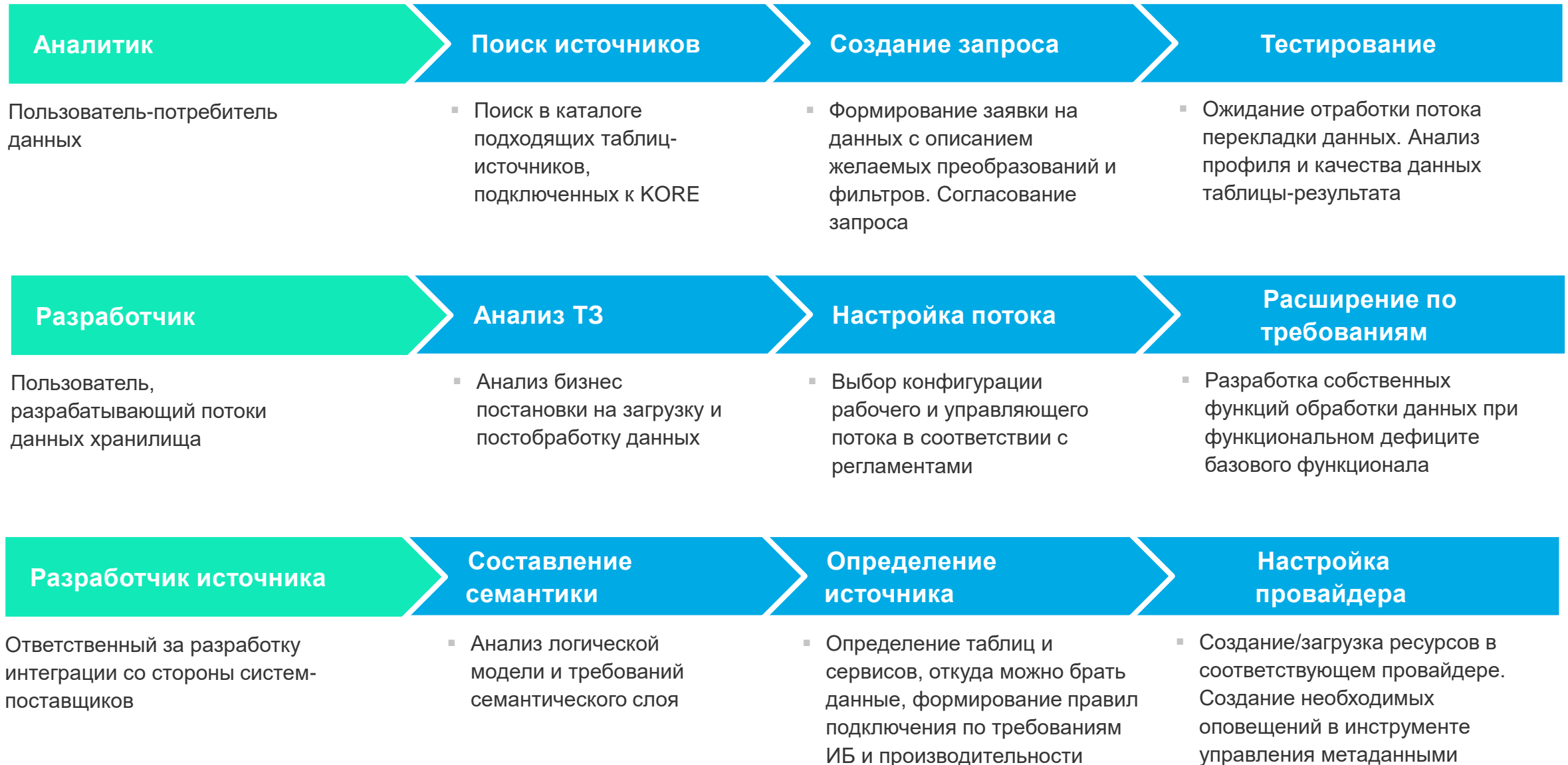
Технические метаданные используются для описания потенциальных объектов источников и получателей данных, а также для анализа графа прослеживаемости

Бизнес-метаданные используются для логических объектов платформы KORE и связанных описаний

Связь с инструментами проверки качества данных для сквозной информации по источникам данных и результате загрузки

Профилирование данных для контроля работы потока: количества обрабатываемых записей

Группы пользователей



Преимущества



Классические ETL решения

Инструмент стандартизует и автоматизирует интеграцию на уровне шага процесса передачи данных

Инструмент для разработчика, частично для системного аналитика

Область применения – только команда хранилища или озера данных

Метаданные – справочная информация, использование ограничено

Инструменты классического стека работы с данными, хорошо применимо для классических архитектур ХД



Платформа KORE DWH

Уровень стандартизации – тип потока перекладки данных. Платформа определяет и инкапсулирует все стандарты работы с потоком

Платформа содержит функционал для совместной работы всех групп пользователей.

Разделяет функции процесса работы с данными между командами ХД и источника

Метаданные – ядро системы, важный рабочий инструмент

Инструменты modern data stack, лучше применимо для распределенных архитектурных шаблонов, таких как data fabric, data mesh и т.д.

Технологический стек



Провайдер ресурсов

Обеспечивает взаимодействие с системой. Реализуется на базе компонентов **KORE.DWH (Python)**. Реализует проверки доступности и готовности источника и предоставляет необходимую информацию для подключения



Шаблоны потоков и трансформаций

Обеспечивает хранение структурированных настроек и их трансформацию в конфигурации объектов транспортного слоя. Реализуется на базе **Airflow**, **DBT**, а также компонентами системы. Версионирование осуществляется в **Gitlab**



Предоставление метрик

Обеспечивает на стороне источника и платформы расчет агрегированных значений, маскирование данных, семплирование перед отправкой данных в платформу. Реализуется на базе **Python** или **DBT**



Управление метаданными

Обеспечивает ведение пользовательских настроек и просмотр метаданных подключённых систем. Единая точка пользовательского интерфейса. Реализуется на базе **React JS** + **OpenMetadada (Arenadata Catalog)**



Транспортный слой

Выполняет подготовленные потоки данных в необходимом режиме взаимодействия. Для реализации могут использоваться компоненты решений Заказчика (например **PXF**) или комбинация **NiFi** или **Airflow** + **Spark**



Хранение наборов данных

Зона хранения подготовленных наборов данных и предоставления их пользователям-заказчикам на базе **S3**, **HDFS**, **GreenPlum**, возможно использование реляционной базы заказчика, совместимой с **DBT**