



Подходы к организации работ с большими
данными в Сбербанке
Москва, 19 апреля 2018

Новая платформа Банка



Единая
Фронтальная
Система

ЕФС



Единое «Окно» контакта с клиентом

БИЗНЕС-
ХАБ



Профиль Клиента

«Интеллект» - логика
и система принятия
решений



ПРОДУКТОВЫЕ
ФАБРИКИ



«Производство» - продукты, транзакции, учет

ФАБРИКА
ДАННЫХ



Исторический
профиль Клиента

«Опыт» - накопление
клиентского опыта и
генерация новых знаний



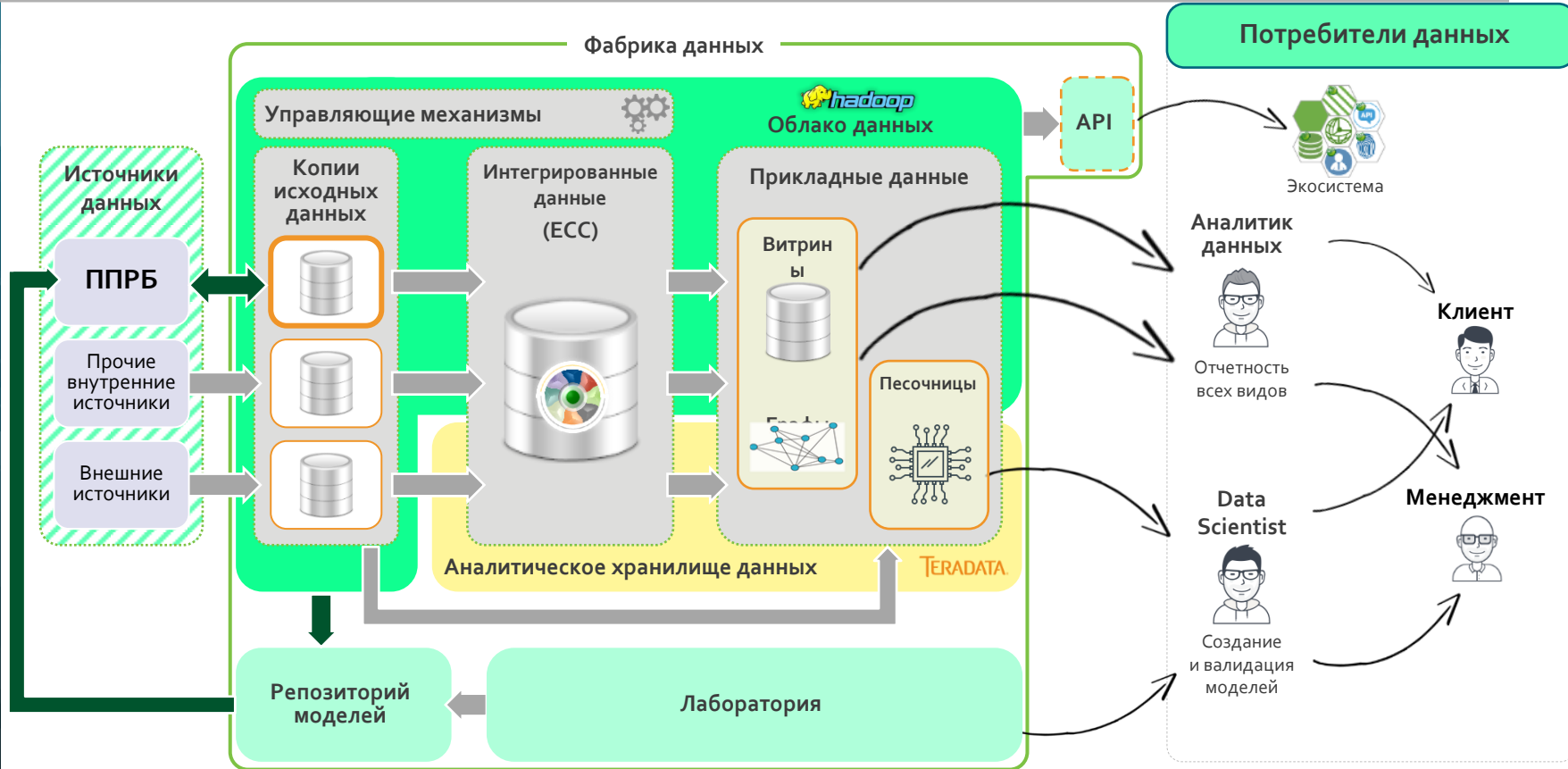
Трансформация в data driven модель

- ❑ За 2017 год создана новая инфраструктура, которая позволяет сегодня решать принципиально новые задачи, чтобы уже в ближайшем будущем встать в один ряд с технологическими компаниями

- ❑ Data Driven организация:
 - ✓ Целостное сочетание традиционной аналитики и больших данных
 - ✓ Аналитика как неотъемлемый компонент ведения бизнеса
 - ✓ Быстрое и гибкое обеспечение решения
 - ✓ Аналитические инструменты доступны в точке принятия решений
 - ✓ Аналитика интегрирована в операционные процессы

- ❑ Данные – актив, которым нужно управлять: Каталог, Качество, Доступность
- ❑ Новые роли и компетенции – Data Scientist, Data Engineer
- ❑ Новые подходы к решению задач

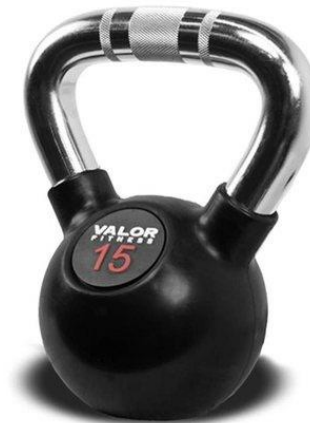
Архитектура Фабрики данных



Облако данных. Факты

Железо:

- ❑ 250 узлов в составе кластера
- ❑ 170 узлов Hadoop data node (HDFS + YARN):
 - ✓ 6800 CPU Cores
 - ✓ 109 Тб оперативной памяти
 - ✓ 10 Пб HDD



Данные:

- ❑ 2.5 Пб - объем данных, план – рост до 10 ПБ в течение 2018г.
- ❑ 15.5 Тб в сутки – объем получаемых изменений
- ❑ 2000 – 5000 транзакций в секунду
- ❑ 400 параллельных заданий загрузки данных
- ❑ 200 ТБ - объем ежесуточно обновляемой информации в репликах
- ❑ 6 часов - время на ежесуточное обновление данных

Использованные технологии

SQL аналитика
на данных
Hadoop



Spark SQL



cloudera
IMPALA

Вычисления
в памяти



NoSQL
СУБД



ETL



Индексный
поиск



Интеграция и
поточная
обработка



Администр.
Управление
Координаци
я



cloudera manager



Управление ресурсам кластера

YARN

Распределенная файловая система кластера



ЕСС и Презентационные слои

Единый семантический слой (ЕСС):

- Единая версия правды
- В основе - Клиентоцентрическая модель сущностей
- Является основным источником для потребителей данных
- Хранит данные в Parquet формате
- Может обновляться batch и streaming процессами



Специализированные витрины данных:

- Рассчитываются на основе ЕСС и прочих источников
- Расчет только batch процессами
- Всегда конкретный заказчик
- Могут являться входом и выходом какого-либо процесса
- Данные хранятся в формате Parquet

Работа с моделями

ФП Библиотека Моделей:

- ❑ Хранение кода в git репозитории
- ❑ Версионирование
- ❑ Передача кода на исполнение

ФП Пакетное Исполнение Моделей:

- ❑ Поддержка моделей на SAS, Python, Scala, Java
- ❑ Извлечение моделей из Библиотеки
- ❑ Исполнение моделей под Spark
- ❑ Сохранение результатов работы модели



ML, AI. Применение в банке

Дескриптивная аналитика

Что происходит?

Описание

Выделение ключевых характеристик, группировка данных

Принцип



Примеры

Сегментация клиентов
Классификация типов событий

Предиктивная аналитика

Что дальше?

Прогнозирование вероятности наступления будущих событий

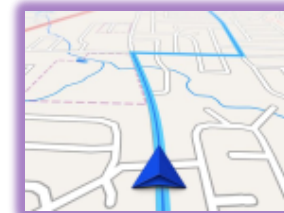


Прогноз показателей
Прогноз банкротства
Предсказание надежности клиента

Предписывающая аналитика

Как мы можем повлиять?

Рекомендация управляющих действий



Персонализация предложений
Блокировка мошеннических транзакций

Спасибо за внимание



Борис Рабинович

Сбербанк-Технологии

Директор центра компетенции
развития VI технологий

BRabinovich.sbt@sberbank.ru