

ОНЛАЙН-ОБРАЗОВАНИЕ

Не забыть включить запись!



Меня хорошо слышно && видно?



Напишите в чат, если есть проблемы!

Ставьте + если все хорошо

Правила вебинара



Активно участвуем



Задаем вопрос в чат или голосом



Off-topic обсуждаем в Slack #канал группы
или #general



Вопросы вижу в чате, могу ответить не сразу

Занятие 25. Case studies. Кейсы компаний.



Цели вебинара

После занятия вы сможете:

1 познакомиться с некоторыми кейсами аналитики больших данных

2 вынести для себя уроки из чужих ошибок

3 вдохновиться для новых проектов 😊

- КХД в крупном банке
- Аналитика телеком данных
- Планирование закупок в энергетической сфере
- Операционный datalake для банка и LC&I
- Streaming платформы в крупном банке
- Deep learning на производстве
- Datalake в фармацевтической компании

КХД в крупном банке

- Крупный банк в Восточной Европе
- Программа КХД с замахом на все источники данных (только крупных источников – 300 типов систем)
- Первые два проекта с суммарной длительностью около года интегрируют 20+ типов источников
- Сложные бизнес-правила реконсиляции данных
- Инструментарий:
 - СУБД: Teradata
 - ETL: Informatica
 - BI/Reporting: BusinessObjects

Сложность («трудность») – это размер усилий, которые необходимо затратить на решение задачи

Как вы думаете, от чего зависит сложность типичные задачи построения аналитической системы?

- Физически консолидировать данные в одном месте («скачать данные в одно место»)
 - Интегрировать данные («сделать из двух таблиц клиентов одну»)
 - Историзовать данные («понять какая запись – это новый клиент, а какая – изменение определенного атрибута»)
 - ...
-
- Объем данных: гигабайт в начальной загрузке, ежедневно, в секунду (в пике)
 - Номенклатура: количество источников, сущностей, связей, таблиц, колонок

В классической компании все данные в источниках можно поделить на два класса:

1. Human generated – данные, которые создаются людьми, и входят в бизнес-процесс: анкеты абонентов в телекоме, описание проблемы в хелп-деске, операции с банковским счетом и проч
2. Machine generated – данные, которые сгенерированы машиной как фиксация события внутри системы: CDR в телекоме, запись в логе веб-сайта или мобильного приложения, показатель с сенсора



- Отдельно оценивать сложность консолидации и интеграции
- Разделение задач консолидации и интеграции по разным задачам (или командам) – путь к успеху
- Знайте (и умейте работать) с ограничениями вашего ETL инструмента
- Автоматизируйте стандартные операции: получение данных, интеграции
- Для масштабного внедрения нужны Continuous Delivery и метаданные с самого начала

**Объем
данных**

Human generated

Machine
generated

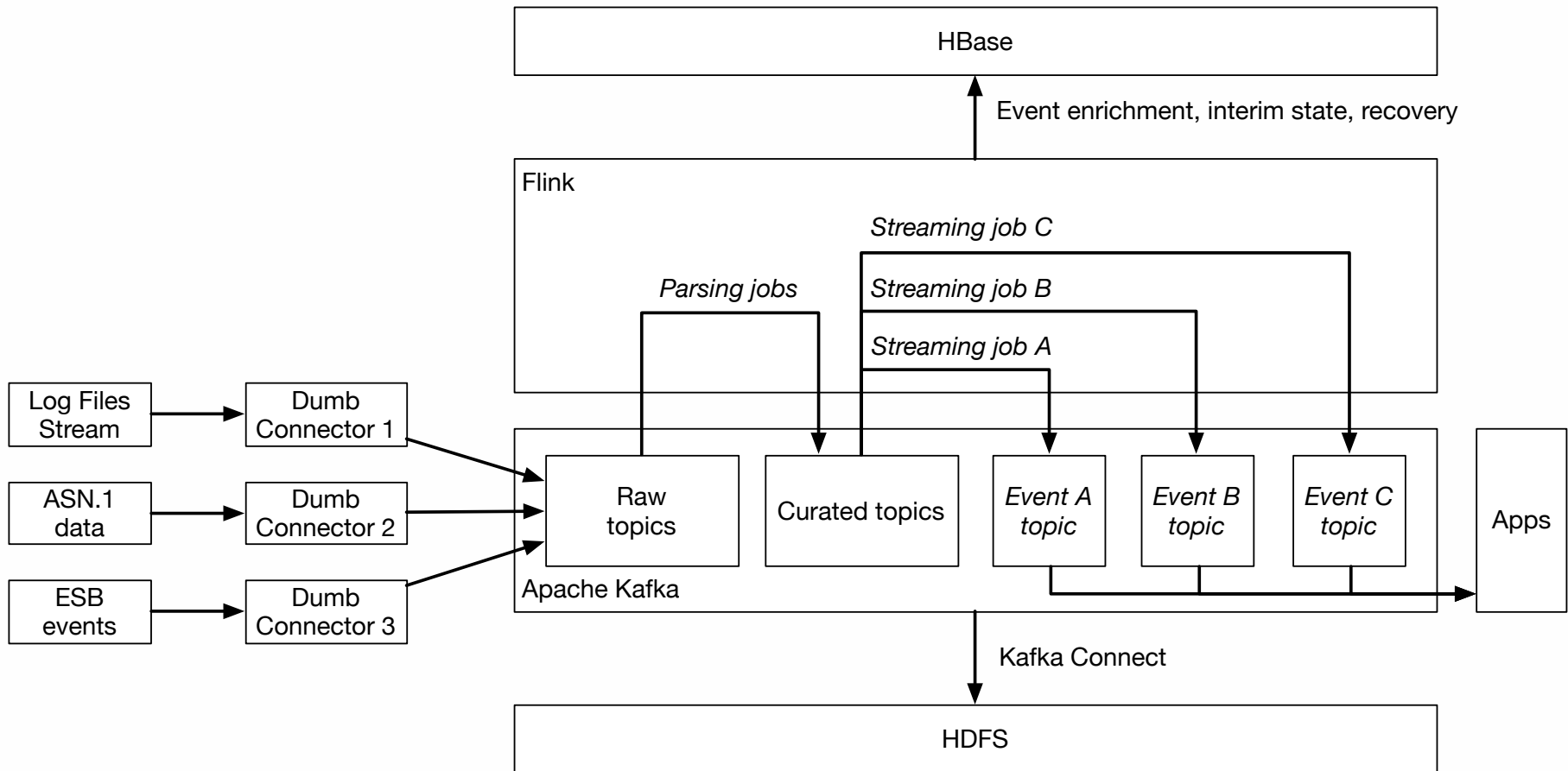
Номенклатура

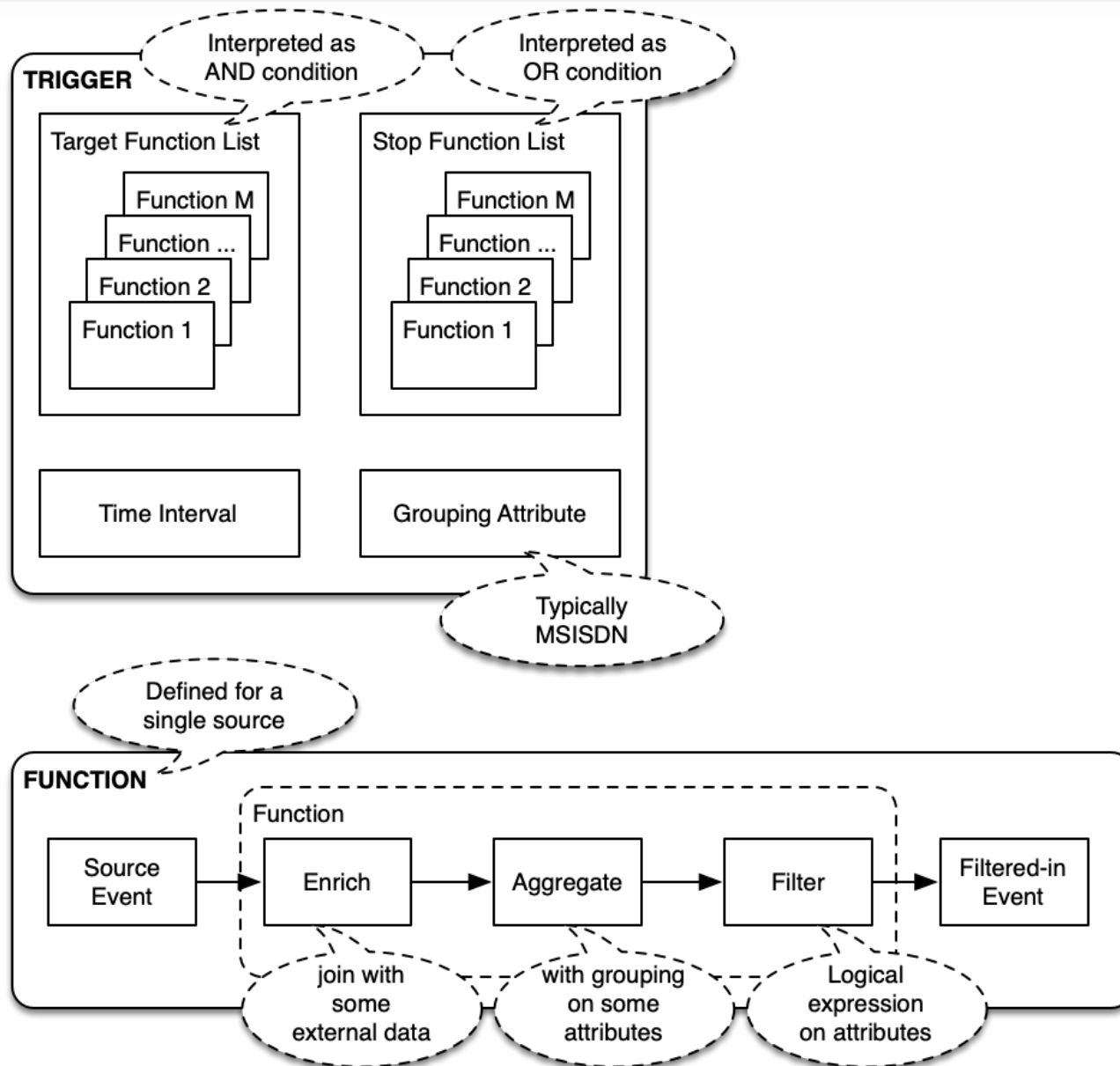
Human generated

Machine
generated

Аналитика телеком данных

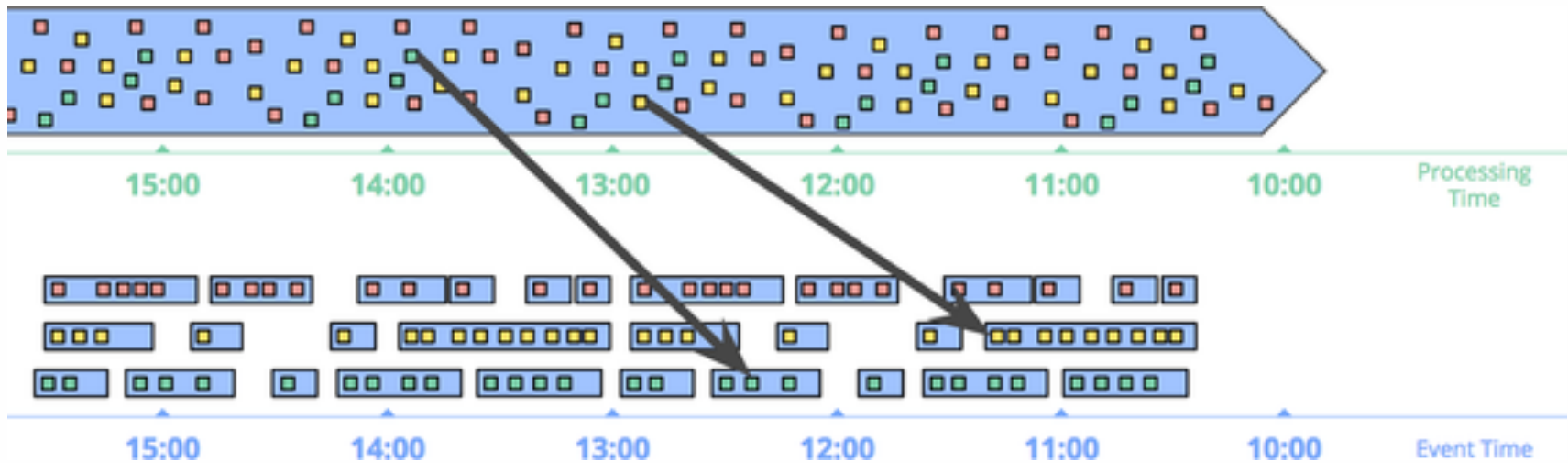
- Крупный мобильный оператор
- Программа Big Data сервисов для внутренних и внешних кейсов
- Платформа операционной аналитики, Machine Learning as a Service, с поддержкой streaming аналитики и CEP (Complex Event Processing)
- Инструментарий:
 - Hadoop, HBase
 - Kestrel -> Kafka, RabbitMQ
 - Storm -> Flink
 - MapReduce -> Spark, Impala



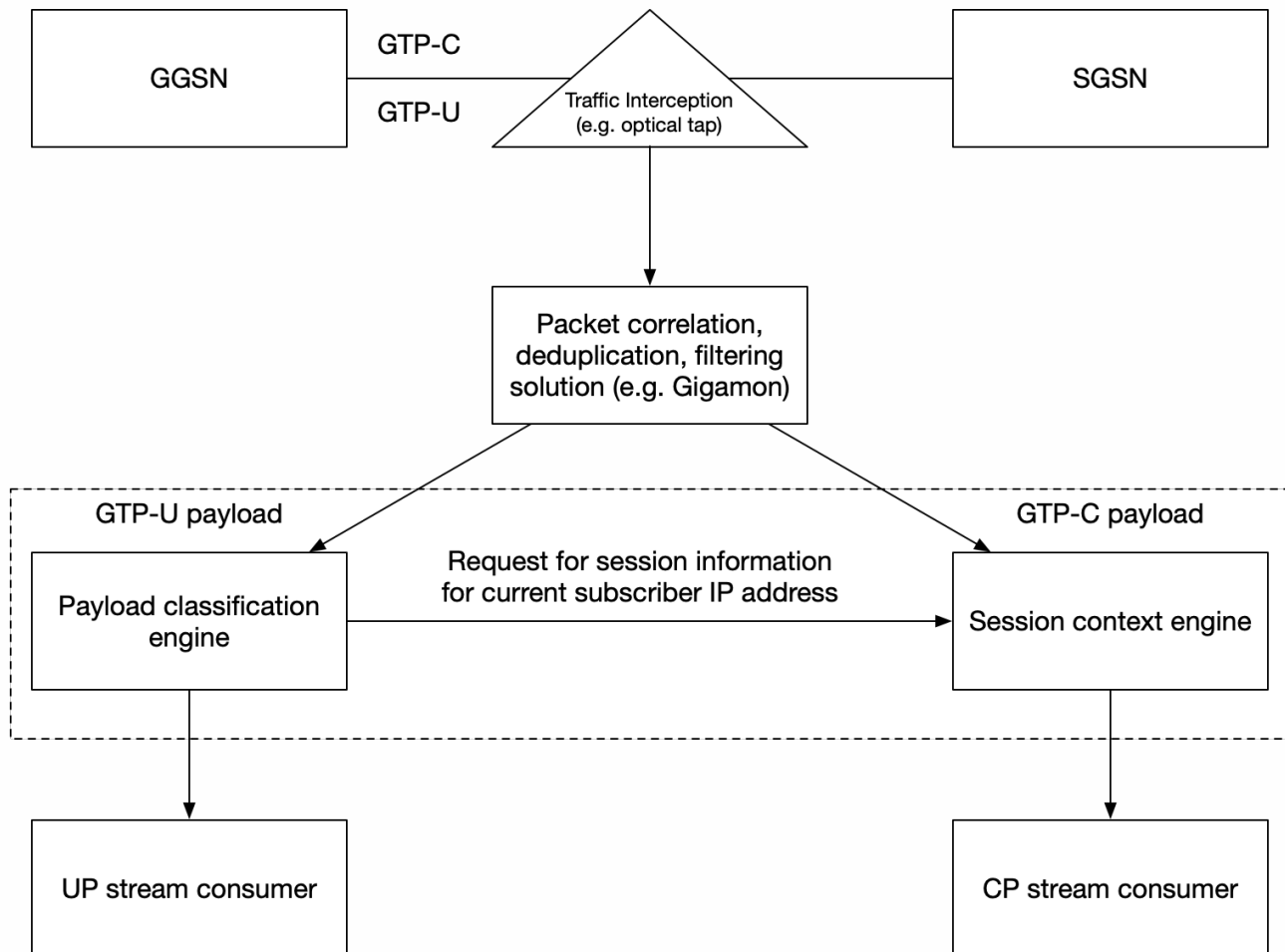


- Внутренние сервисы: финансы, сеть, маркетинг
- Внешние сервисы: финансы, smart city, реклама
- На пике активности оператор отчитался что 1% годовой выручки приносили Big Data сервис, речь о новом бизнесе, без учета улучшений в маркетинге и обслуживании клиента

- Группировка событий по определенным бизнес-правилам
- Сущность «сессия» крайне важна для целого ряда метрик, оценивающих вовлеченность пользователя, качество сервиса и других аспектов взаимодействия
- Сессонизацию сложно делать на больших данных в батче, так как она требует джоина по условию больше или меньше



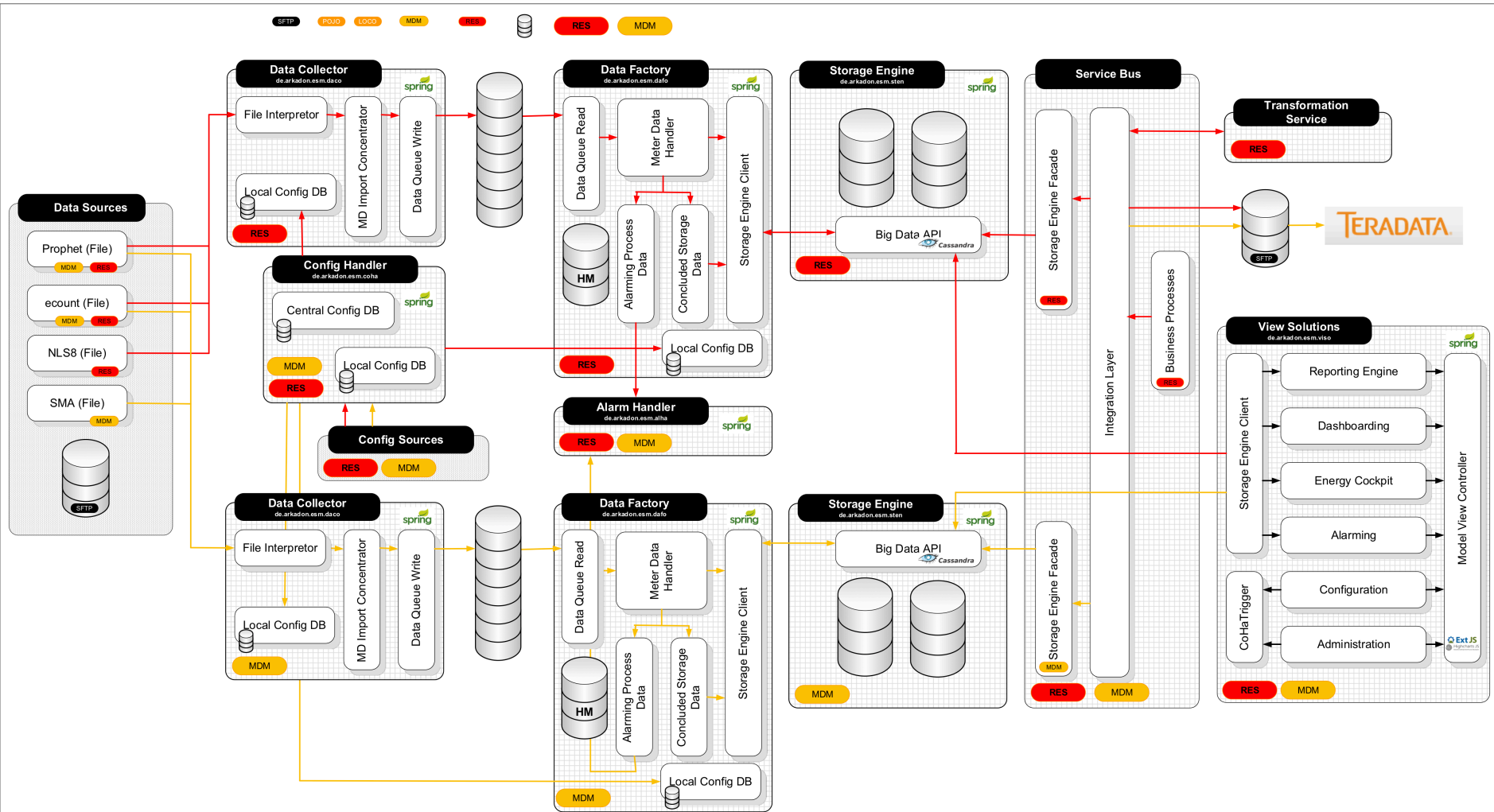
- Join «очень большой таблицы» с другой «очень большой таблицей»
- Bucketed Tables могут помочь, но все-равно медленно
- На помощь Streaming Join

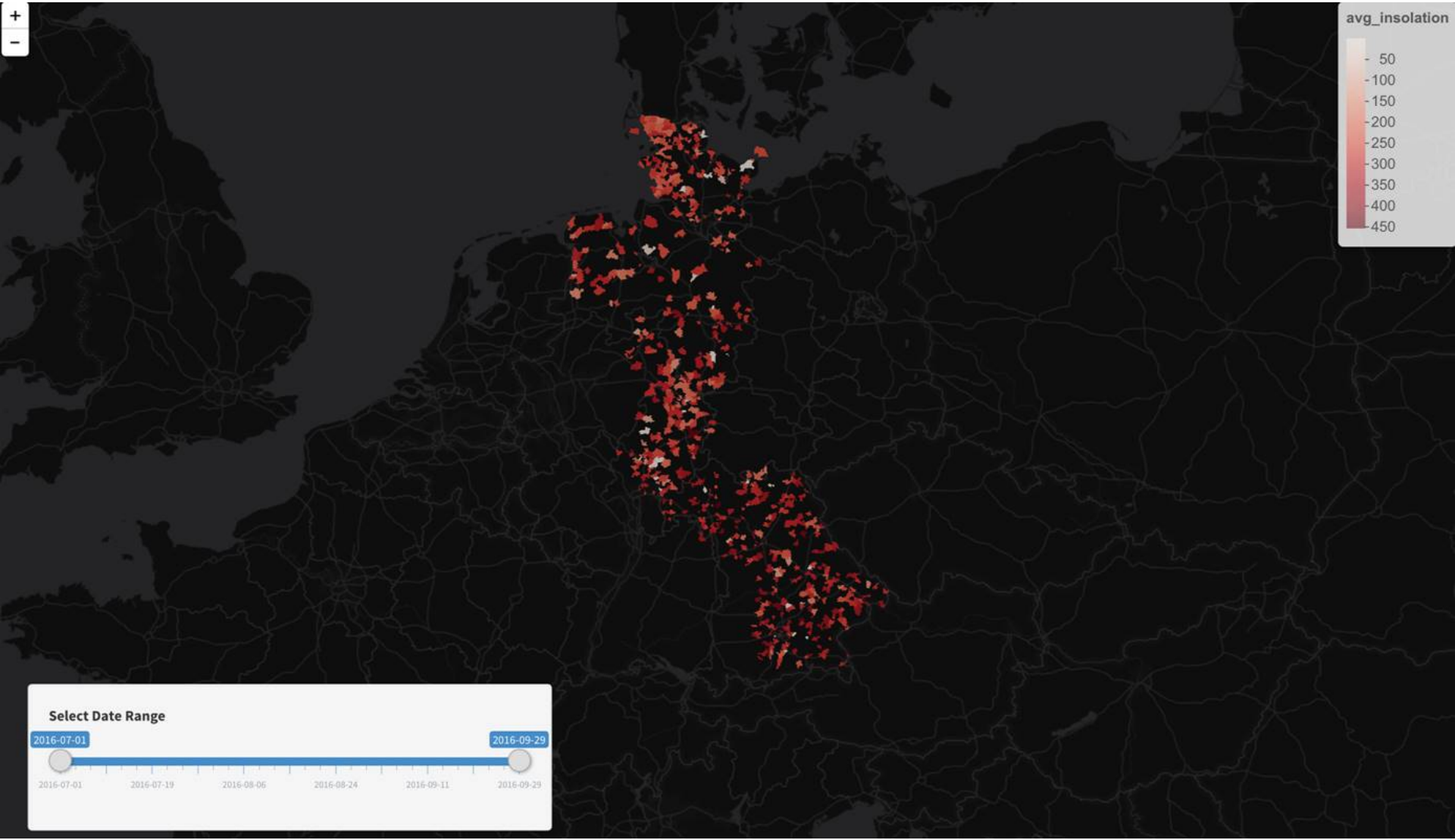


- Делая решение – думайте и о платформе
- Streaming может быть не усложнением, а упрощением конкретного решения
- Качество должно быть таким же ориентиром как производительность или скорость выкатки

Планирование закупок в энергетической сфере

- Распределительная компания в Европе обязана покупать всю энергию от воспроизводимых источников (ветер, солнце, гидро)
- При этом энергия ветра и солнца очень нестабильны, и для компенсации компания вынуждена закупать энергию на спотовом рынке
- Чем точнее предсказаны спрос (и предложение на рынке RES), тем точнее может быть сделана ставка, сохраняя деньги компании
- Инструменты:
 - Cloudera Hadoop
 - Cassandra
 - Spark
 - R, Shiny





- Forecast
- Performance
- Benchmark
- SMA

Zeithorizont:

2015-01- bis 2016-11-

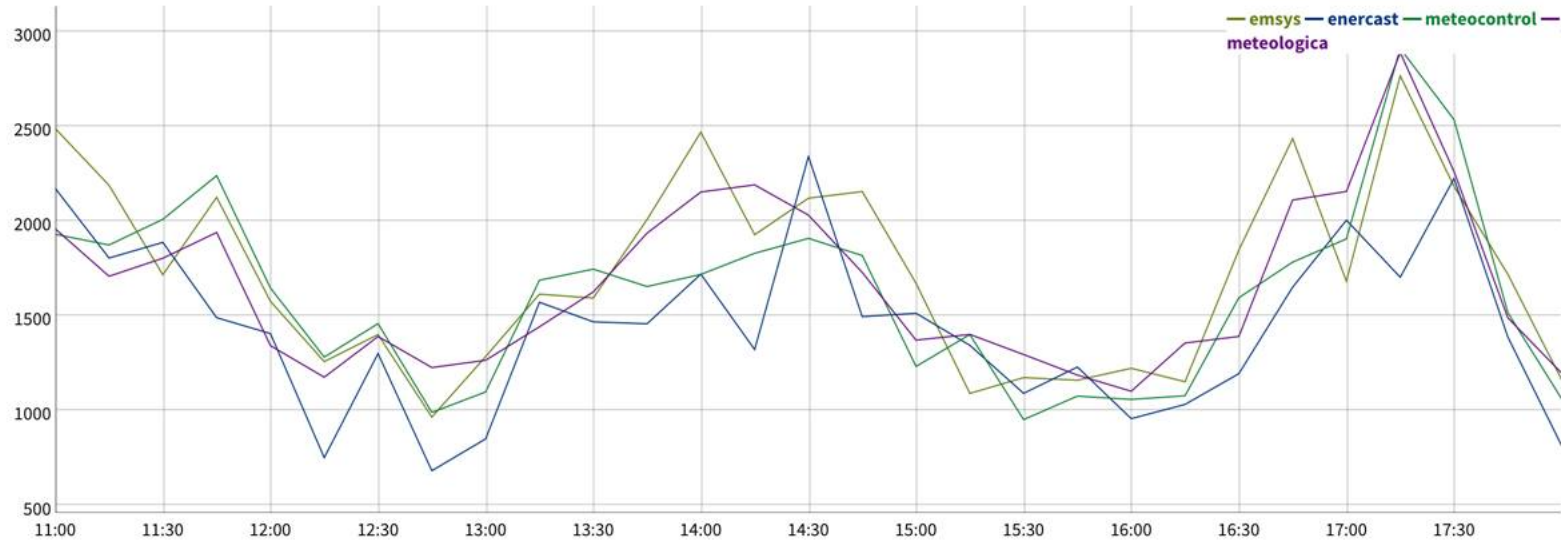
Ist-Datum:

2016-11-01

Provider

- Emsys
- Enercast
- Meteologica
- Meteocontrol

PV - Forecast



 EMSYS
21

 ENERCAST
20

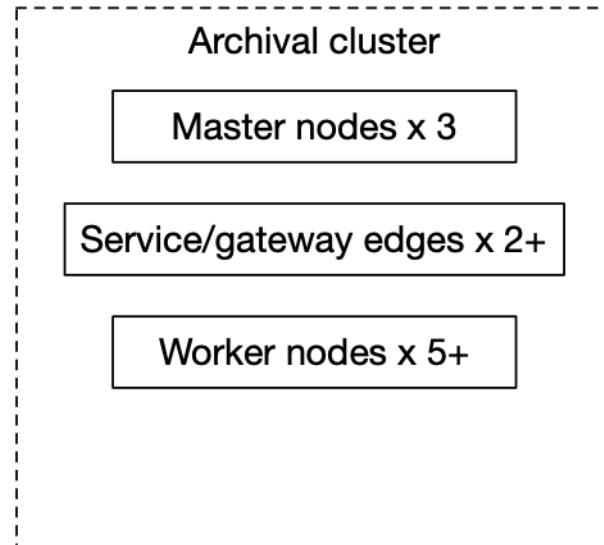
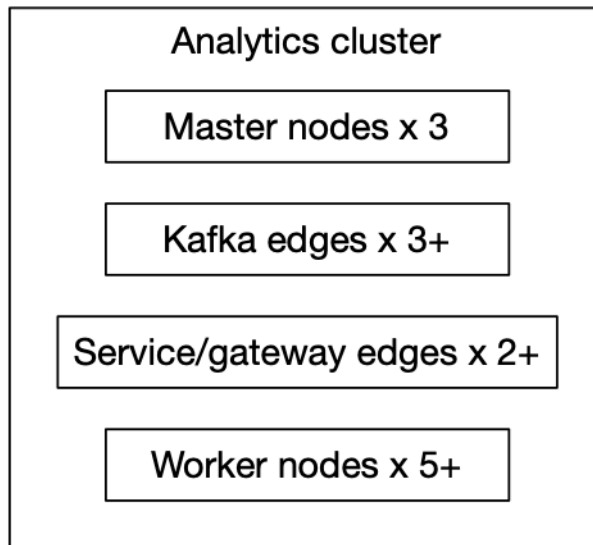
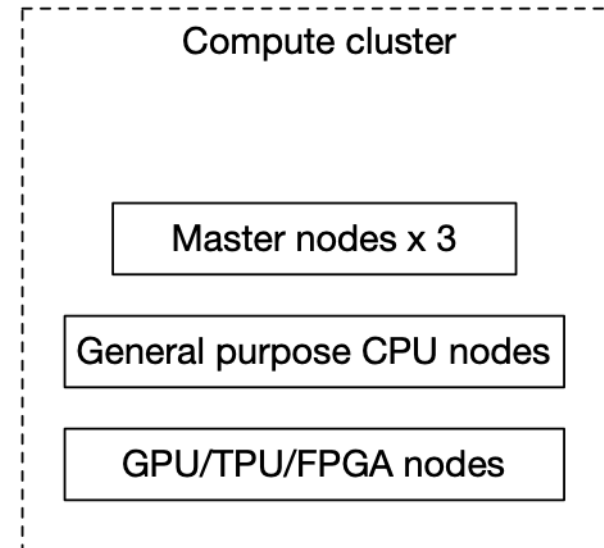
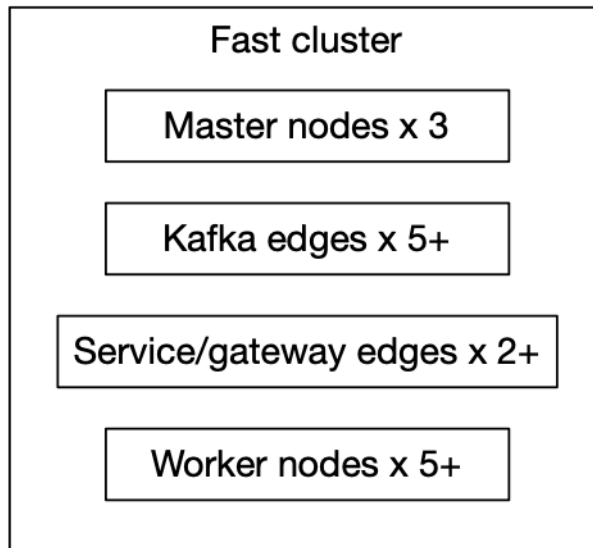
 METEOLOGICA
25

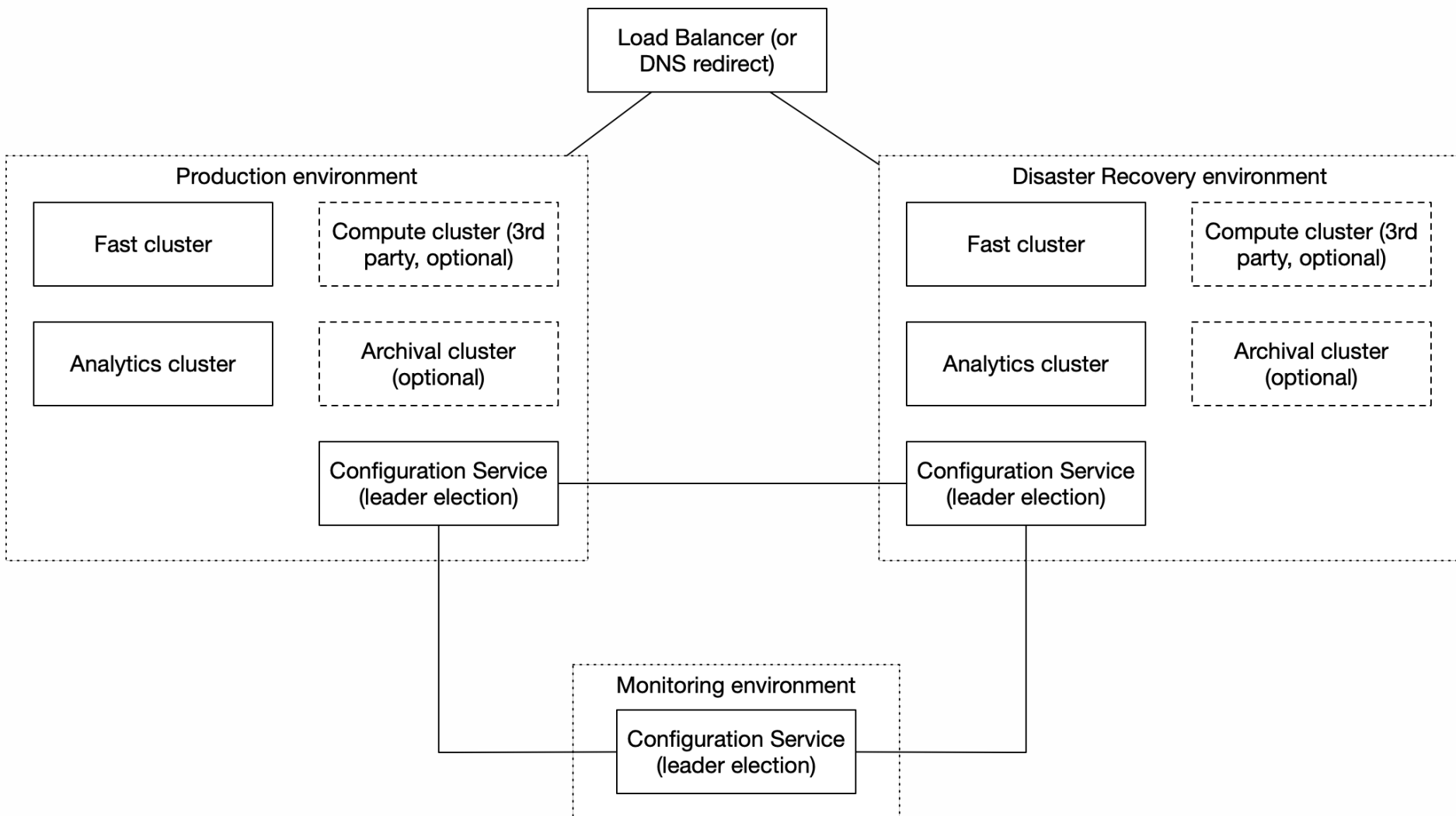
 METEOCONTROL
19

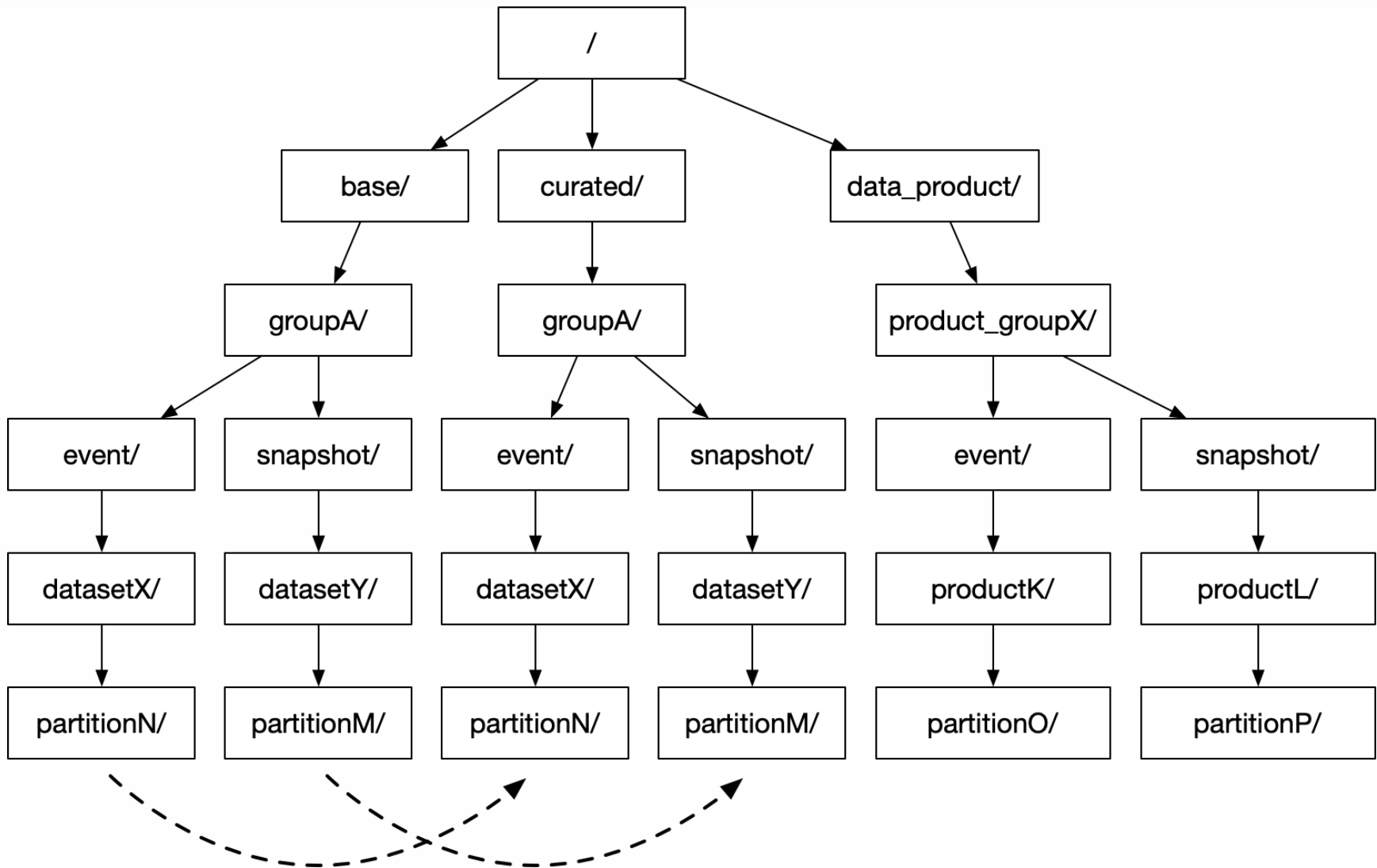
- Работа с временными рядами требует отдельного навыка от аналитиков
- Сенсорные данные нужно агрегировать
- Для кейсов сопряженных с физическим миром важны внешние данные, они стоят денег

Операционный datalake для банка и LC&I

- Крупный универсальный банк в Северной Европе
- Сильная Data Science практика
- Один из бизнесов – портфельное управление Large Corporations & Institutions
- Datalake не только решает задачи аналитики, но и работает бэкендом для многих операционных сервисов
- Инструментарий:
 - Hortonworks Hadoop
 - Confluent Kafka
 - Open Shift, Nvidia GPUs







1. Обновление данных: для качественной аналитики нужно иметь онлайн обновление рыночных данных
2. Complex Event Processing: возможность детектировать события по определенной бизнес-логике
3. Microbatch/batch on-demand: возможность запускать быстрые пересчеты по триггерам

Для CEP важным подспорьем является использование хэш-структур данных:

- HyperLogLog реализует count distinct
- Count-min sketch реализует groupBy,count
- И несколько других

Эти структуры можно сохранять, комбинировать, использовать также в батчевой аналитике

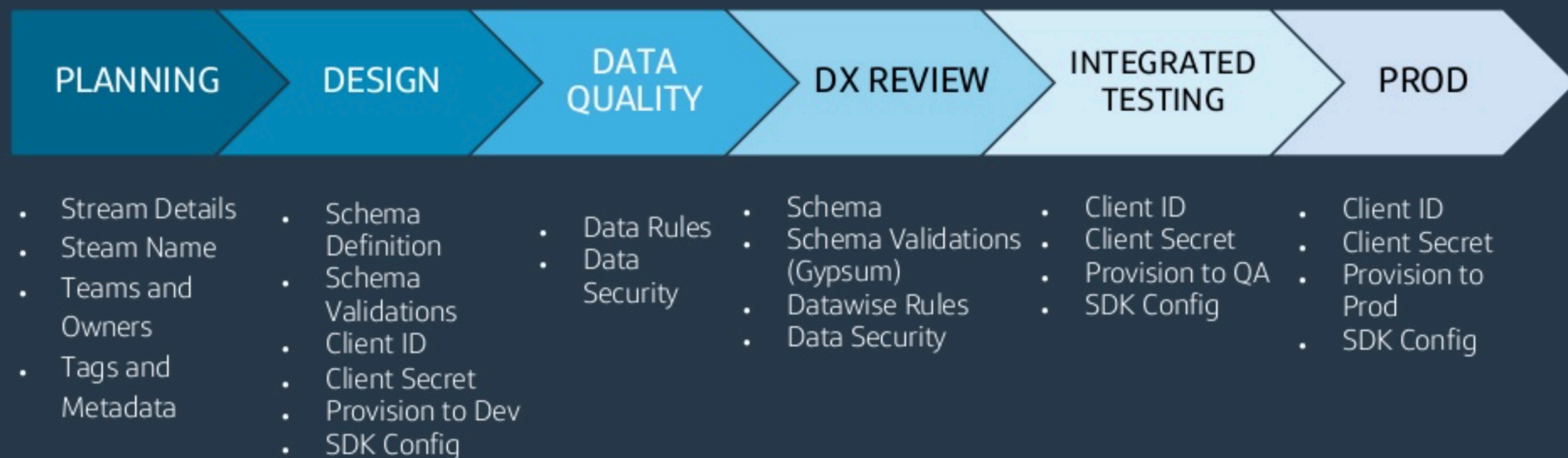
- Необходимость использовать свое железо не означает отсутствие эластичности и отказоустойчивости
- ... и свое железо очень быстрое!
- Data locality не нужна (для пакетной аналитики)
- Определенные задачи можно решать гораздо быстрее если пожертвовать точностью
- Если на раннем этапе создать хорошую модель безопасности, то внедрение может идти гораздо быстрее и спокойнее

Streaming платформы в крупном банке

- Крупный банк в США
- Отдельная организация (со своим ВП) отвечает за стриминг, как основу бизнеса
- <https://www.confluent.io/kafka-summit-SF18/building-an-enterprise-streaming-platform-capital-one>
- Инструментарий:
 - Confluent Kafka
 - Spark
 - ...

Creating a Public Stream via DevExchange

How to create a stream during Design Time



Provisioning occurs during Design, QA and Prod

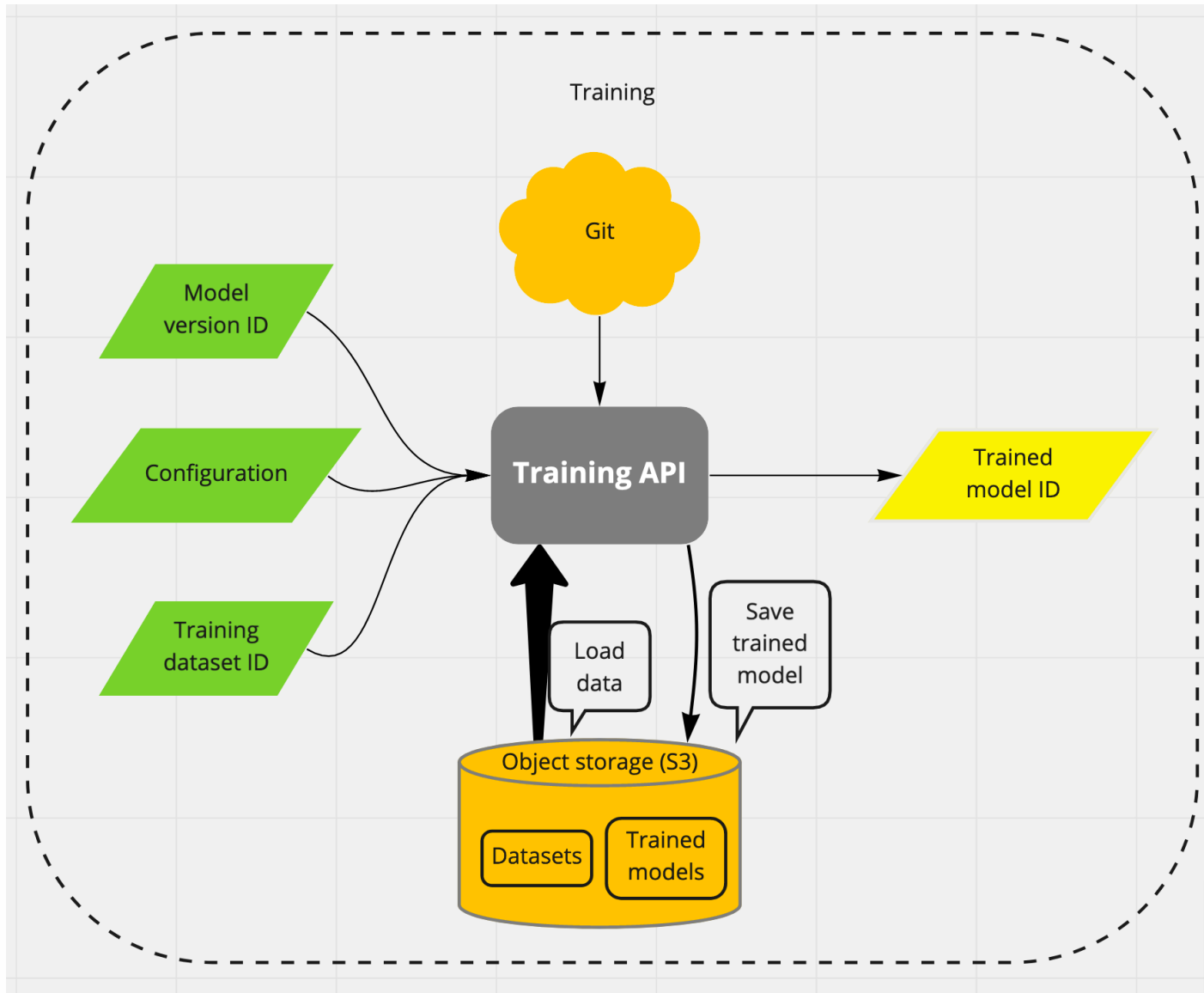
SDK Features

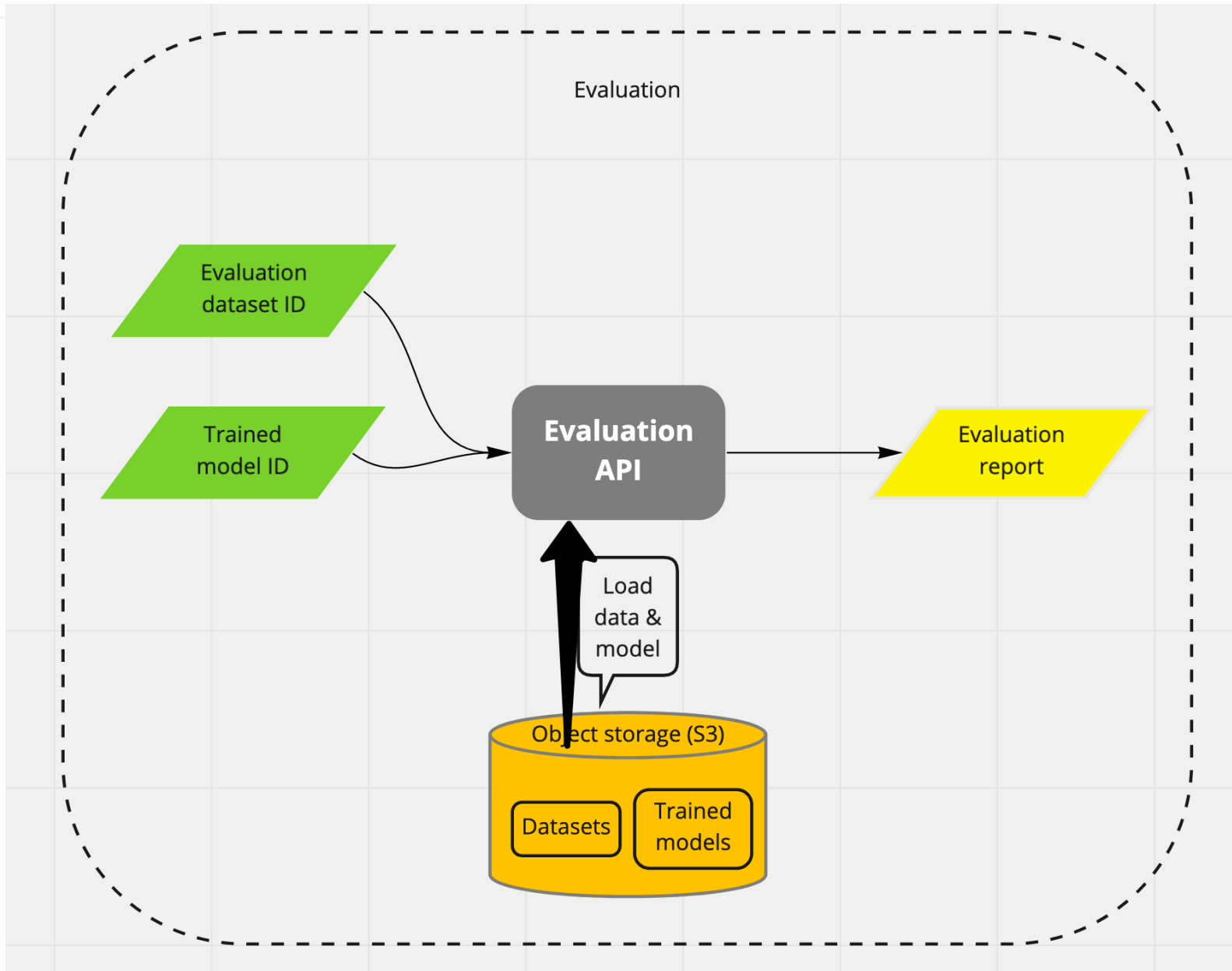
Feature Description	SCALA	PYTHON	SPARK
Creation of PEM files or Java Keystore files from DX credentials	✓	✓	✓
Data Validation against registered Avro schema	✓	✓	✓
Avro conversions	✓	✓	✓
Execution of data rules and watermarking	✓	✓	✓
Auto fill Enterprise Envelope fields	✓	✓	✓
Secondary cluster support for Producers	✓	✓	-
Tracking Consumer positions between clusters	✓	✓	-

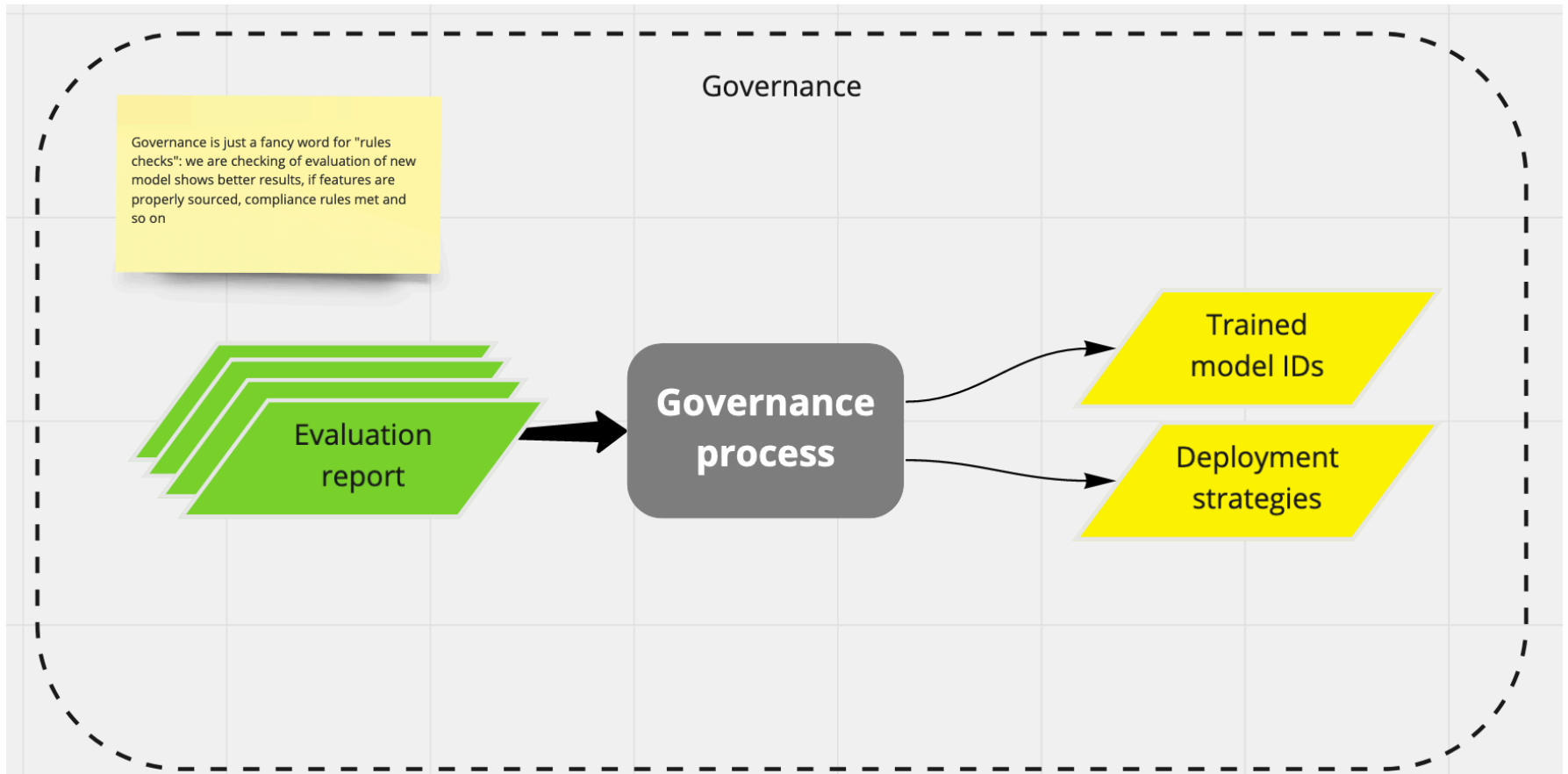
- Платформа как внутренний продукт со своим SDK
- Можно реализовать требования «на клиенте» – не стоит этим пренебрегать
- Помогает быть отдельной структурой внутри компании

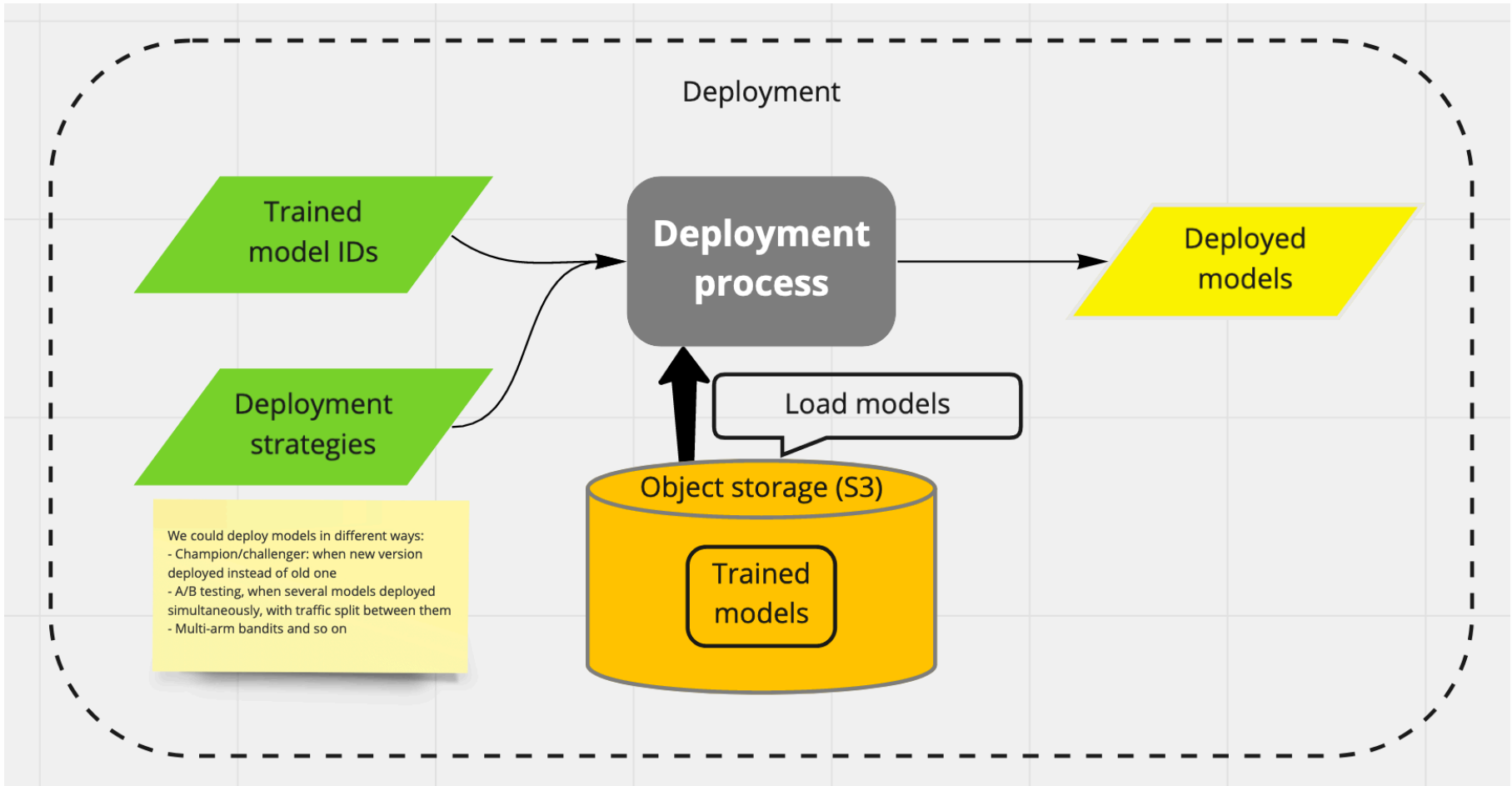
Deep learning на производстве

- Автопроизводитель
- Точечная сварка дает сбои, нужна аналитика которая обнаруживает плохие швы по фотографии
- Deep learning модель обнаруживает плохие швы на основе данных собранных ультразвуковым исследованием
- Инструменты:
 - AWS стэк
 - Seldon
 - Kubeflow









- Machine learning это процесс параллельный дата/ПО инжинирингу, и должен быть формализован
- Метаданные это ключ к успеху

Datalake в фармацевтической компании

- Компания которая разрабатывает и производит лекарства и медицинское оборудование
- Американская, с присутствием в Европе: HIPAA + GDPR + ...
- Datalake(s) на AWS
- Инструментарий:
 - EMR
 - S3

- Данные хранятся в S3, структурированные по владельцу
- Владелец управляет набором IAM ролей
- Специальное ПО Tenant Manager отслеживает какой группе пользователь относится, и создает динамическую Security Configuration
- Пользователю предоставляется API которое позволяет запустить джоб или кластер только с подключенной Security Configuration

- Замена ключей суррогатными / хэшами с фиксированной солью или сокрытие их через REVOKE не является анонимизацией, это маскирование или псевдоанонимизация
- Анонимизация это отрыв персоны от события, и требует удаления ключа как такового
- Существует два семейства моделей анонимизации:
 - k-анонимность
 - дифференциальная приватность

- Определение: массив данных называется k-анонимным, если для каждого подмножества атрибутов существует, для любой строки существует k-1 строка с таким же значением атрибутов
- Позволяет рассуждать о защите персональных данных в отторгаемых массивах данных
- Реализуется группировкой записей и перекодированием (генерализацией) значений атрибутов
- Развитие темы:
 - ρ -чувствительная k-анонимность
 - l-разнообразии
 - t-близость
- Все модели k-анонимности подвержены атаке связывания

- Определение (упрощенное): рандомизированный запрос M к набору данных D считается дифференциально приватным, если для любой пары D и D_i (получен удалением одной строки из D) невозможно идентифицировать удаленную строку
- Реализуется при помощи следующих механизмов:
- Сэмплирование исходного набора данных
- Добавление шума со сходным распределением для маскирования исходных значений

- Аналитика на анонимных данных возможна
- Развитая модель безопасности и ролей, связанная с инфраструктурной автоматизацией позволяет избежать строительства большого количества разрозненных песочниц

Следующий вебинар

Тема: Бонус. Дальнейшее развитие Hard skills + Soft skills.



Понедельник 2019.09.09 в 20.00



Ссылка на вебинар будет в ЛК за 15 минут

**Заполните, пожалуйста,
опрос о занятии**



**Спасибо
за внимание!**





ОНЛАЙН-ОБРАЗОВАНИЕ