

Разработка **Real Time Decision Making** решения на базе **Hortonworks Data Platform & Data Flow**

Алексей Кузнецов

Data Engineer, Software Architect

Антон Ярмолюк

Data Engineer, Chapter Lead

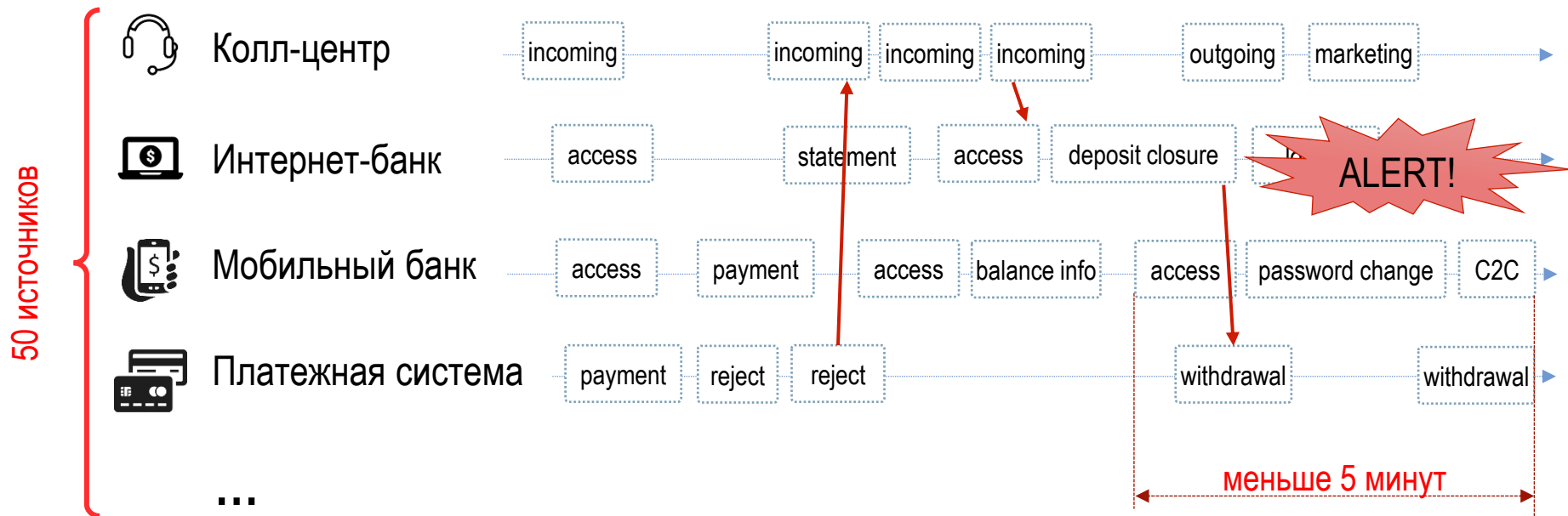


**Райффайзен
БАНК**



Задача

Источники событий



Аналитика Real-time

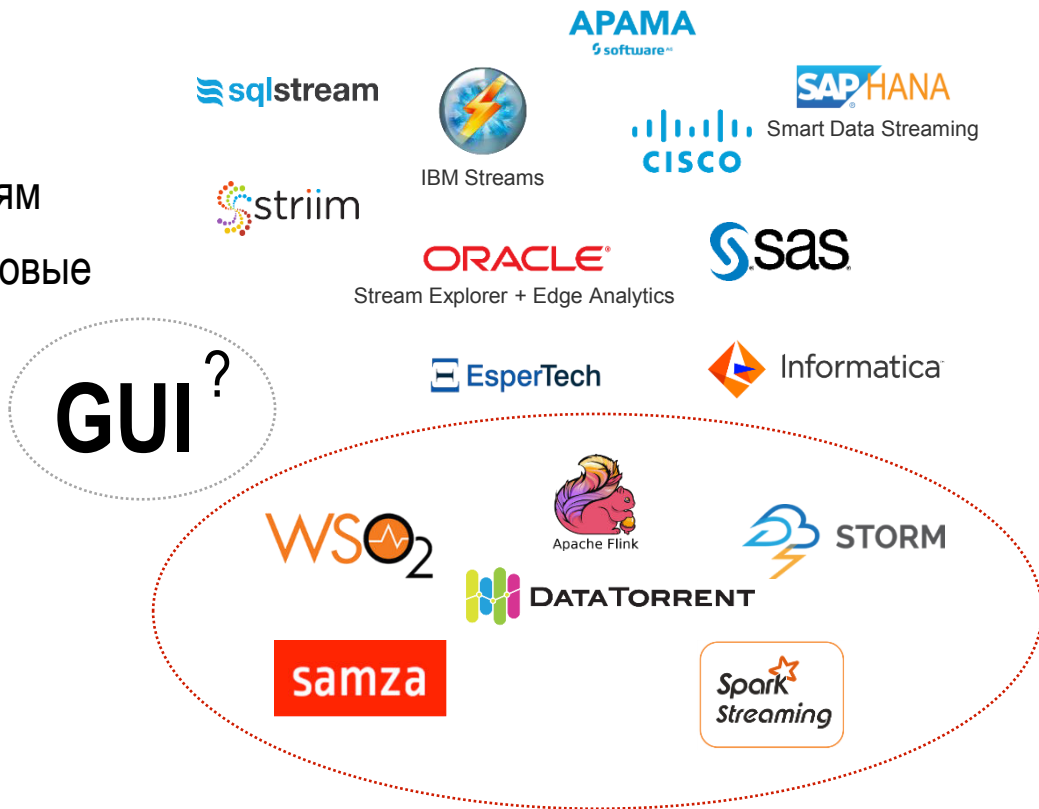
Real-time Decision Making

Complex Events Processing

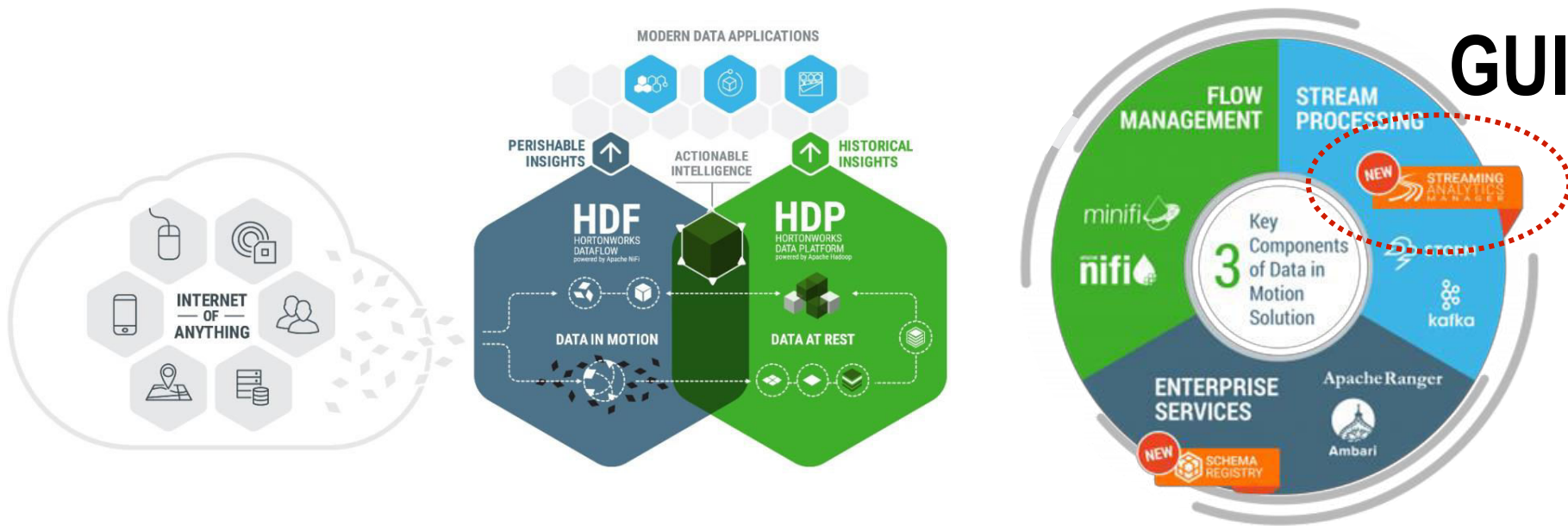


Наш кейс

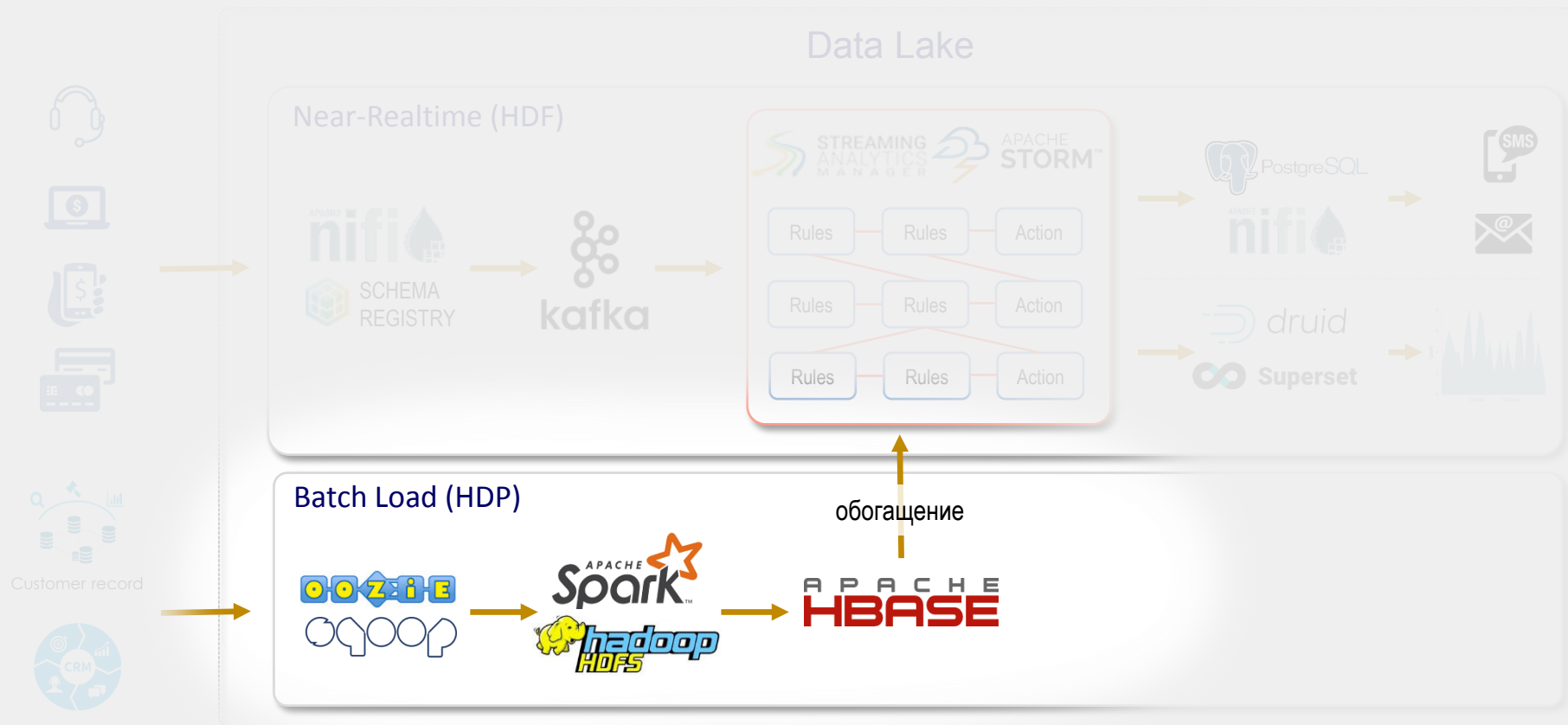
Дать возможность бизнес-пользователям
самостоятельно проводить маркетинговые
кампании на платформе RTDM



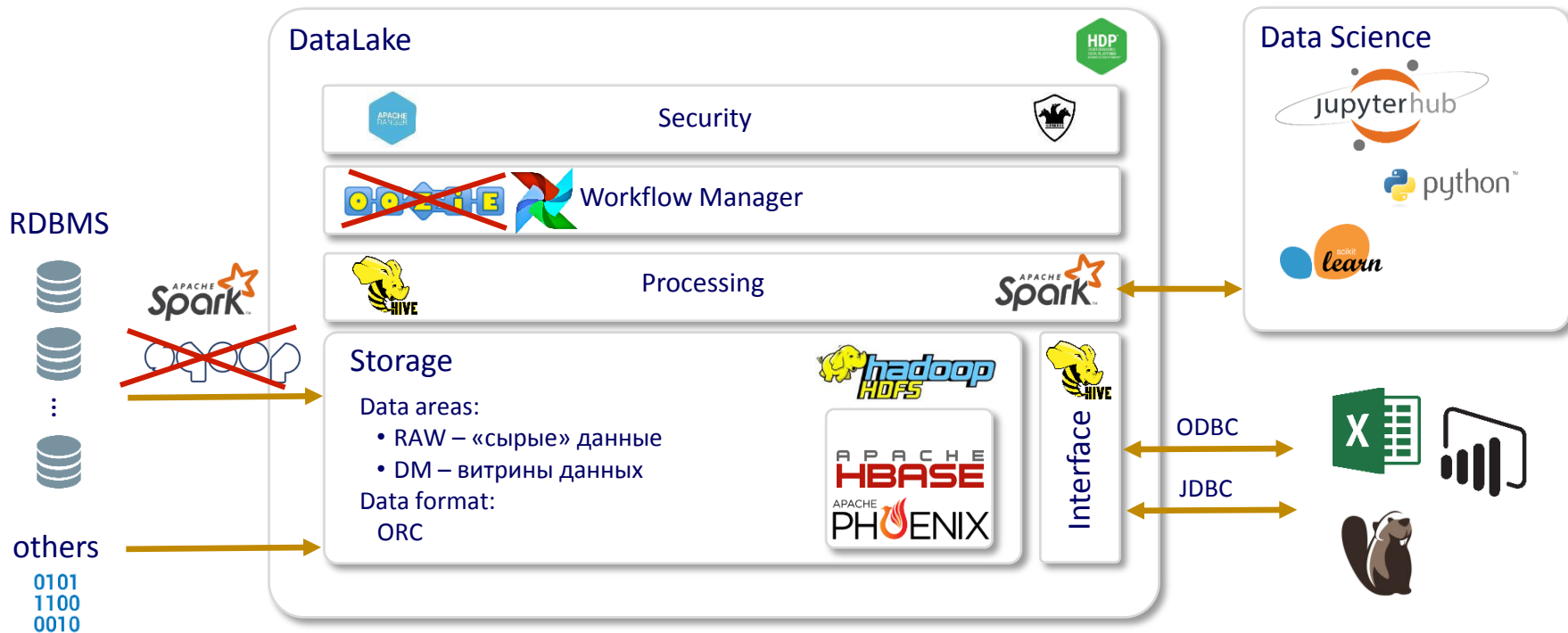
Hortonworks Data Flow 3.0



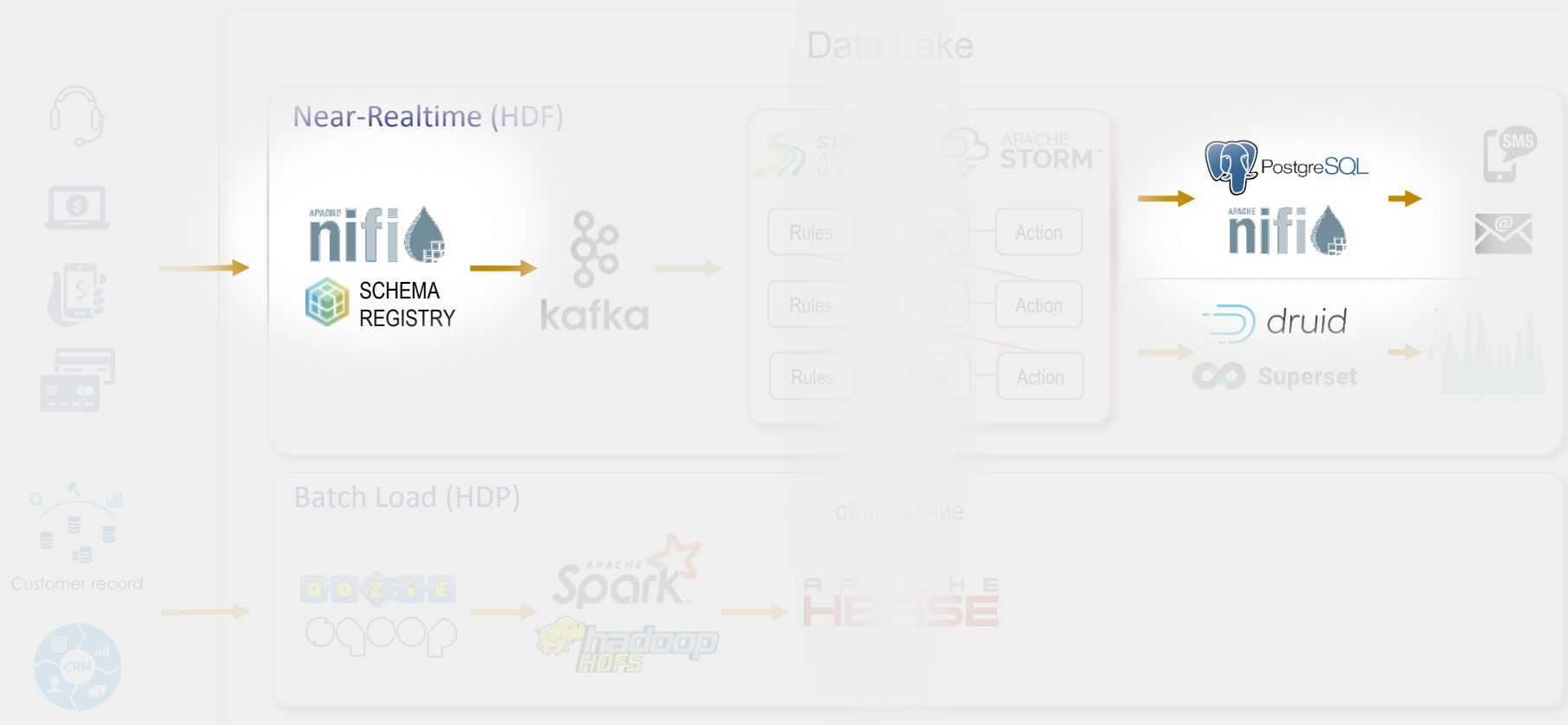
Архитектура RTDM-решения



HDP - Batch Load

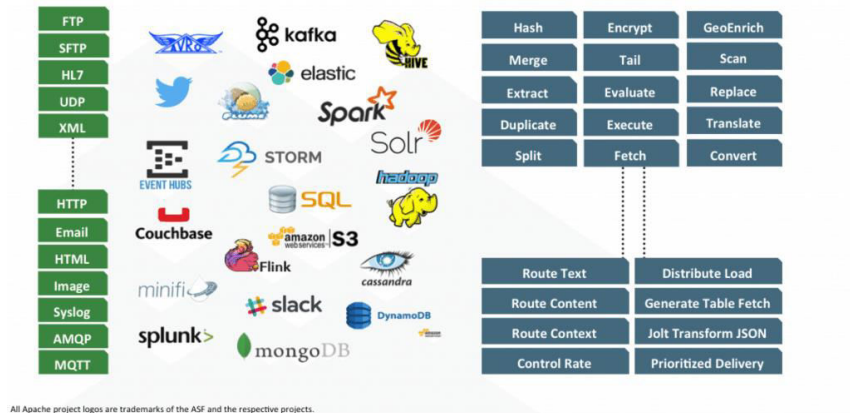


Архитектура RTDM-решения



HDF – Apache NiFi

Connecting Data Between Ecosystems Without Coding: 170+ Processors



All Apache project logos are trademarks of the ASF and the respective projects.

- Разработка АНБ, передан в Open Source в 2014 году
- Реализует концепцию Flow Based Programming (FBP)
- Визуальный интерфейс разработки

- UUID
- Name
- Size
- Entry Time

HEADER

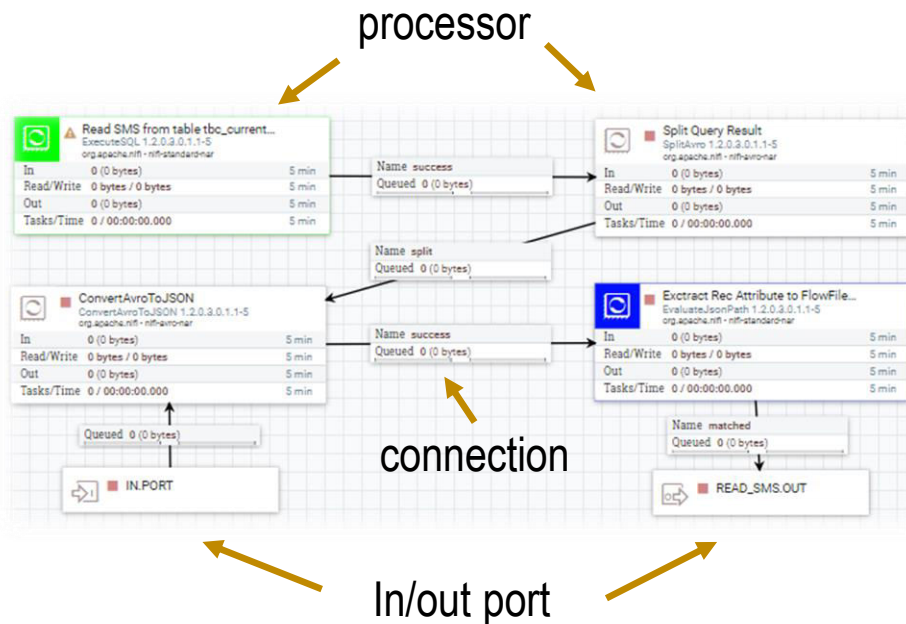
Attribute Map
[Key,Value]

CONTENT

- Types
 - Events
 - Objects
 - Files
 - Messages
 - Media
- Formats
 - JSON
 - Avro
 - Text
 - Mp4
 - Proprietary
- Sizes
 - Bytes to GBs

Основная единица передаваемых данных - FlowFile

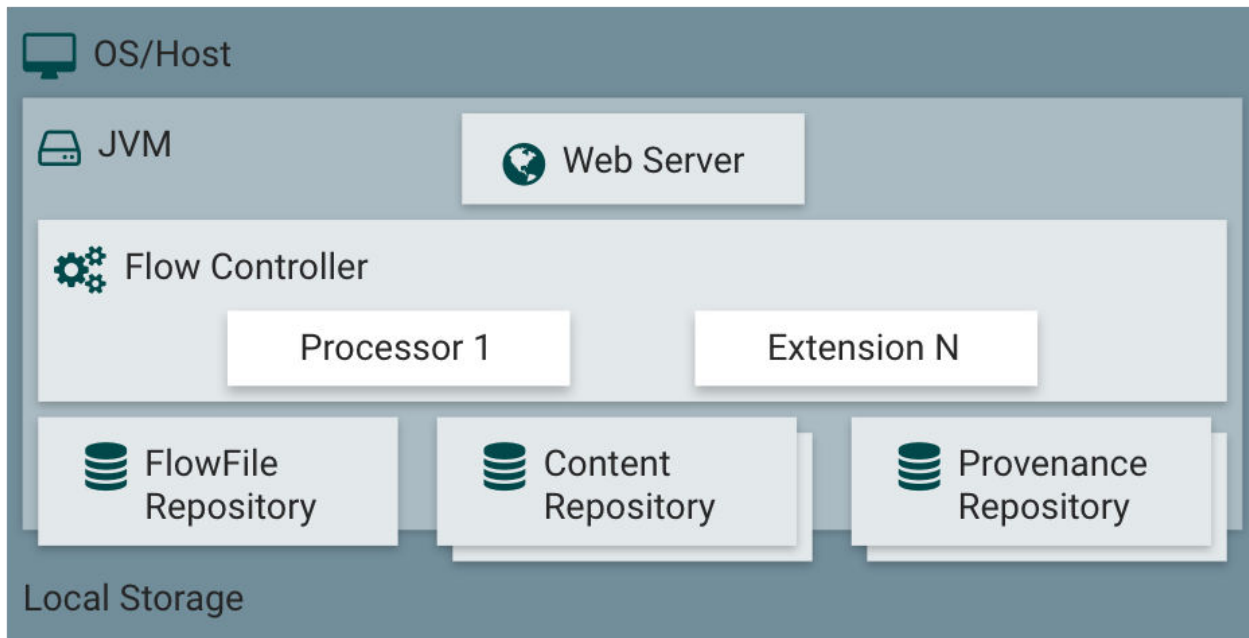
HDF – Apache NiFi



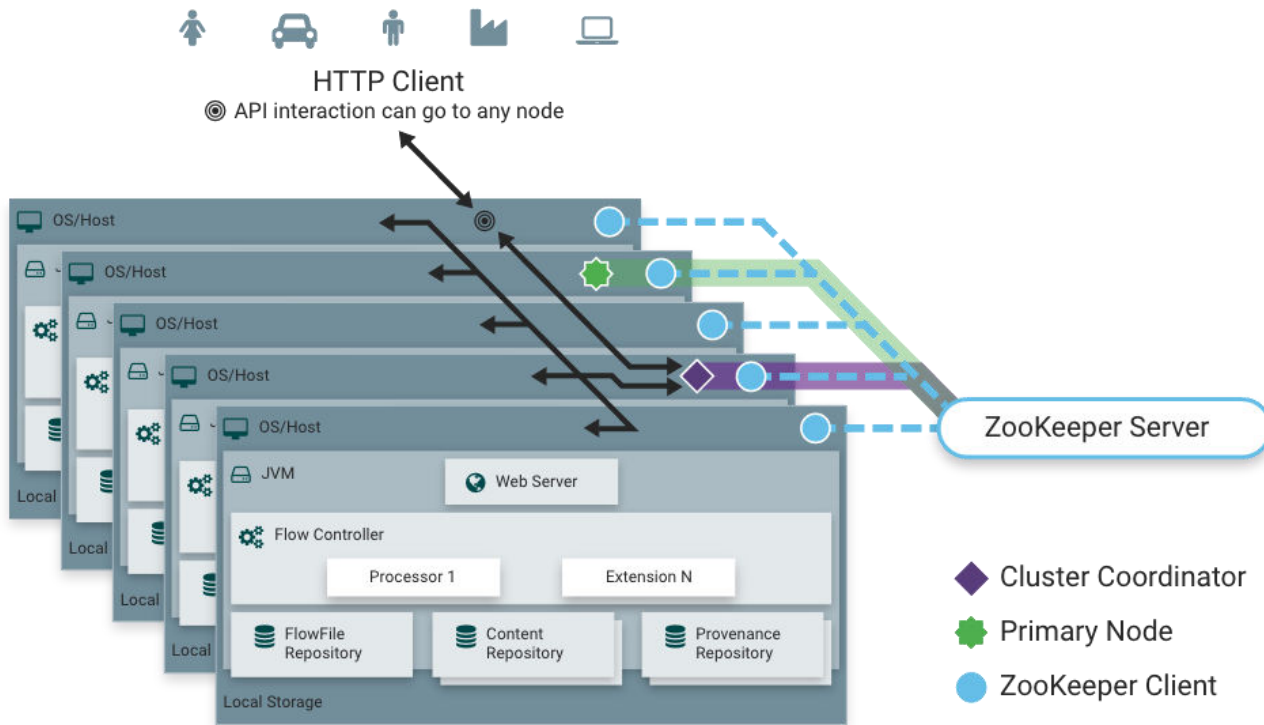
Processor Group

Read SMS from Table			
Queued	0 (0 bytes)		
In	0 (0 bytes) → 1	5 min	
Read/Write	0 bytes / 0 bytes	5 min	
Out	1 → 0 (0 bytes)	5 min	
No comments specified			

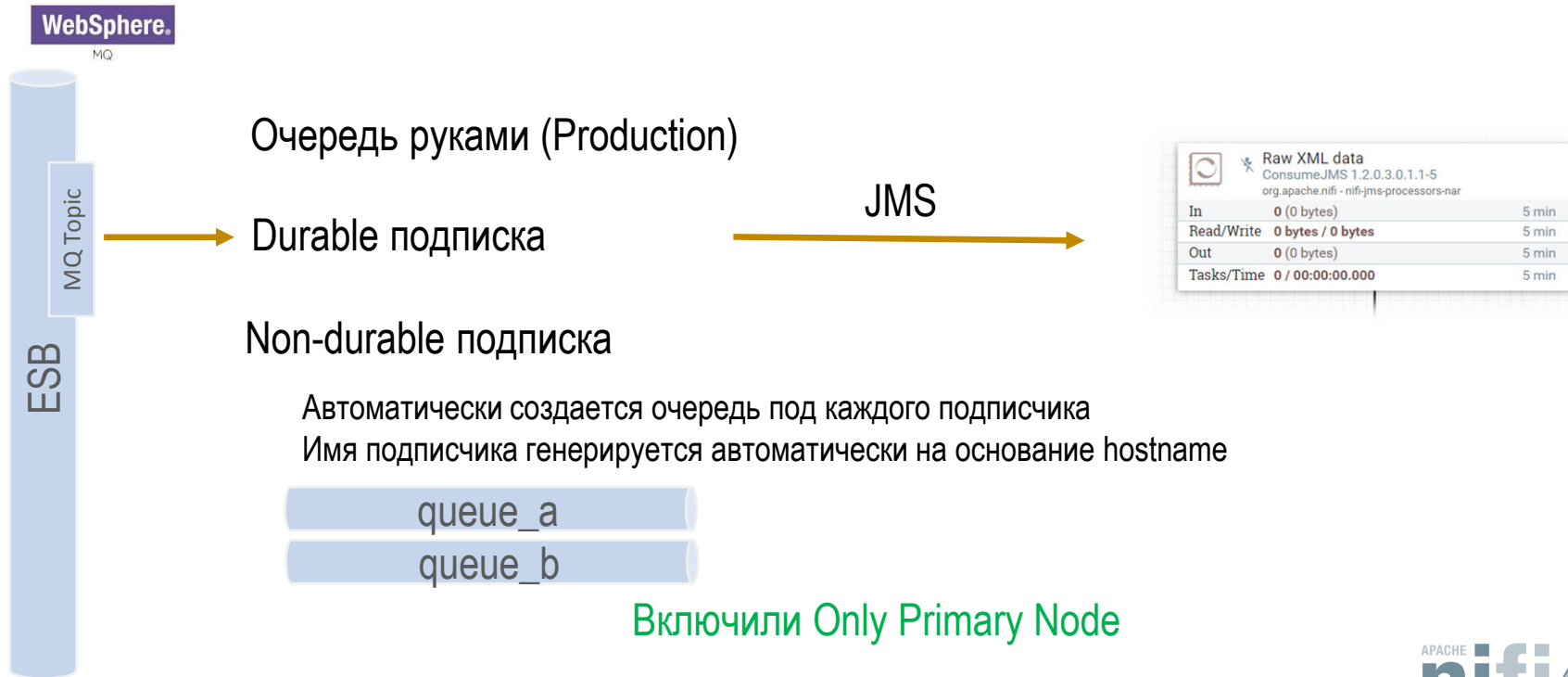
HDF – Apache NiFi



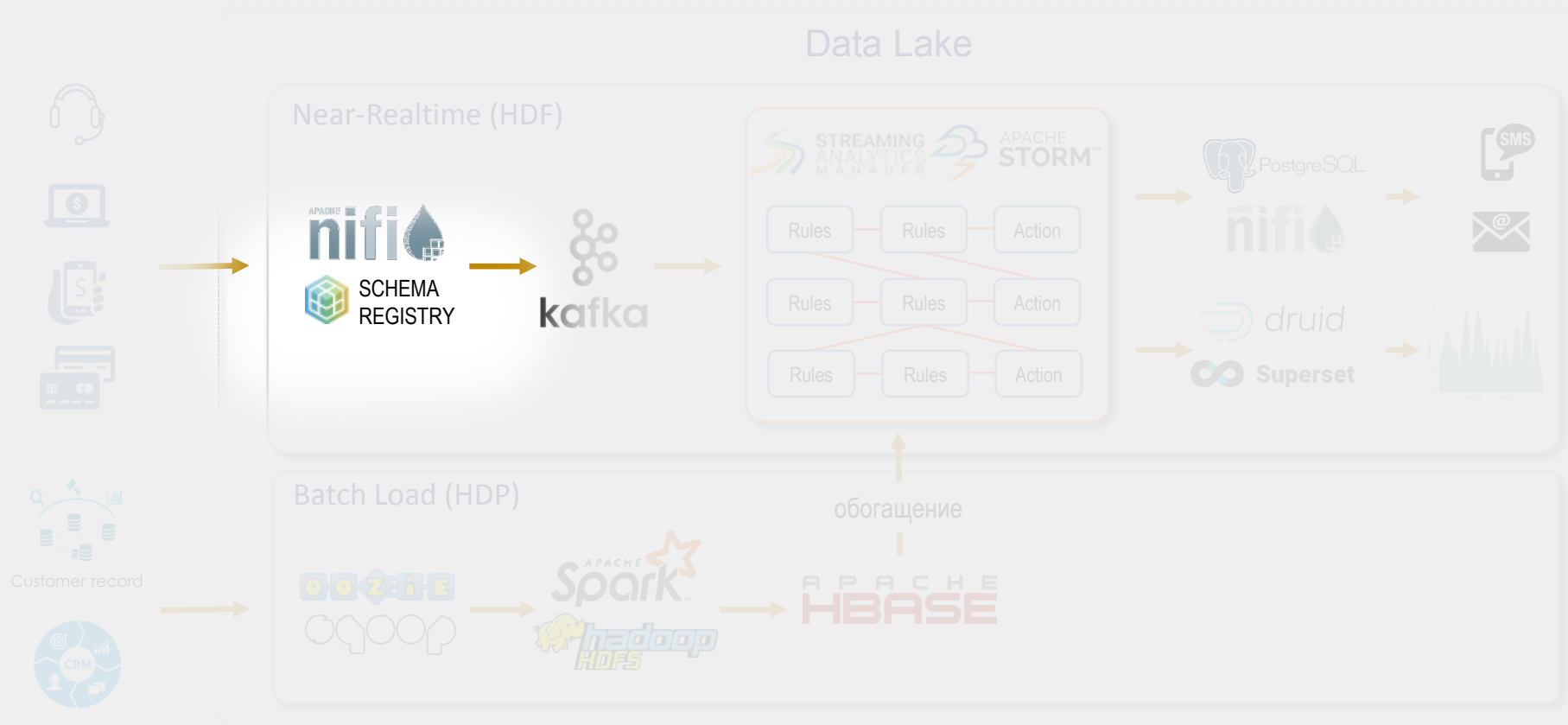
HDF – Apache NiFi



HDF – Apache NiFi



Архитектура RTDM-решения



Schema Registry



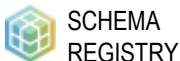
Topic

message #1

011110001100

message #2

011110001100



Schema

```
{
  "type": "record",
  "namespace": "ru.raiffeisen.bigdata.sam",
  "name": "SOURCE.RUB_TRANSFER",
  "fields": [
    {
      "name": "clientID",
      "type": "string"
    },
    {
      "name": "amount",
      "type": "double"
    }
  ]
}
```



Stream

clientID

amount

client #1

500 p

clientID

amount

client #2

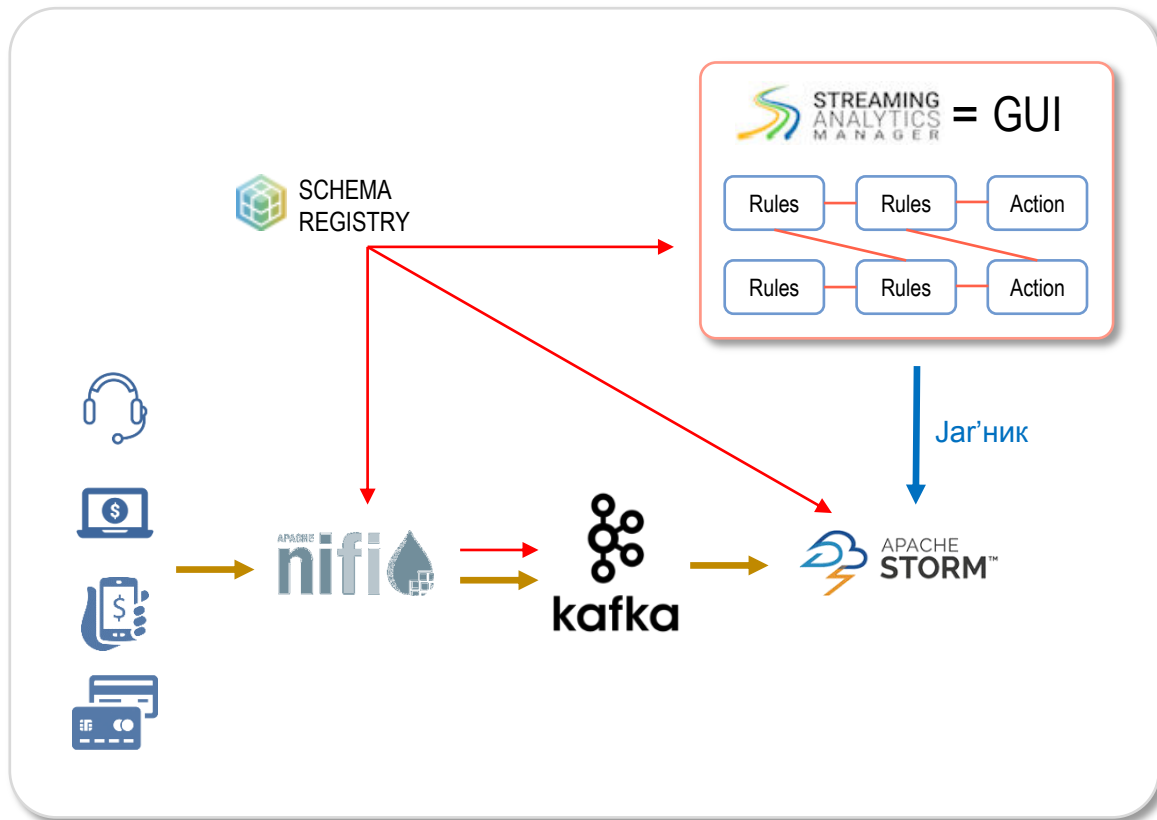
1000 p



SCHEMA
REGISTRY GUI

	TYPE	GROUP	VERSION	SERIALIZER & DESERIALIZER
RONLINE.C2C_FROM_THIRD_PARTY_BANK	avro	Kafka	1	0
DESCRIPTION: Перевод Card to Card со стороннего банка				
CHANGE LOG: v1 23d 18h 5m 51s ago CREATED				
<pre>1 { 2 "type": "record", 3 "namespace": "ru.raiffeisen.bigdata.sam", 4 "name": "RONLINE.C2C_FROM_THIRD_PARTY_BANK", 5 "fields": [6 { 7 "name": "requestId", 8 "type": "long" 9 }, 10 { 11 "name": "clientId", 12 "type": "string" 13 }, 14 { 15 "name": "amount", 16 "type": "double" 17 } 18] 19 }</pre>				
RONLINE.INTERNAL_RUB_TRANSFER	avro	Kafka	1	0
RONLINE.DEPOSIT_OPEN	avro	Kafka	1	0
RONLINE.DEPOSIT_REPLENISHMENT	avro	Kafka	1	0

Schema Registry



- **Schema Registry**

- создаем и описываем схему "стрима"

- **Nifi**

- читаем схему "стрима" (метаданные)
 - сериализуем и отправляем в Kafka

- **Kafka**

- не знает о том, что в ней лежит

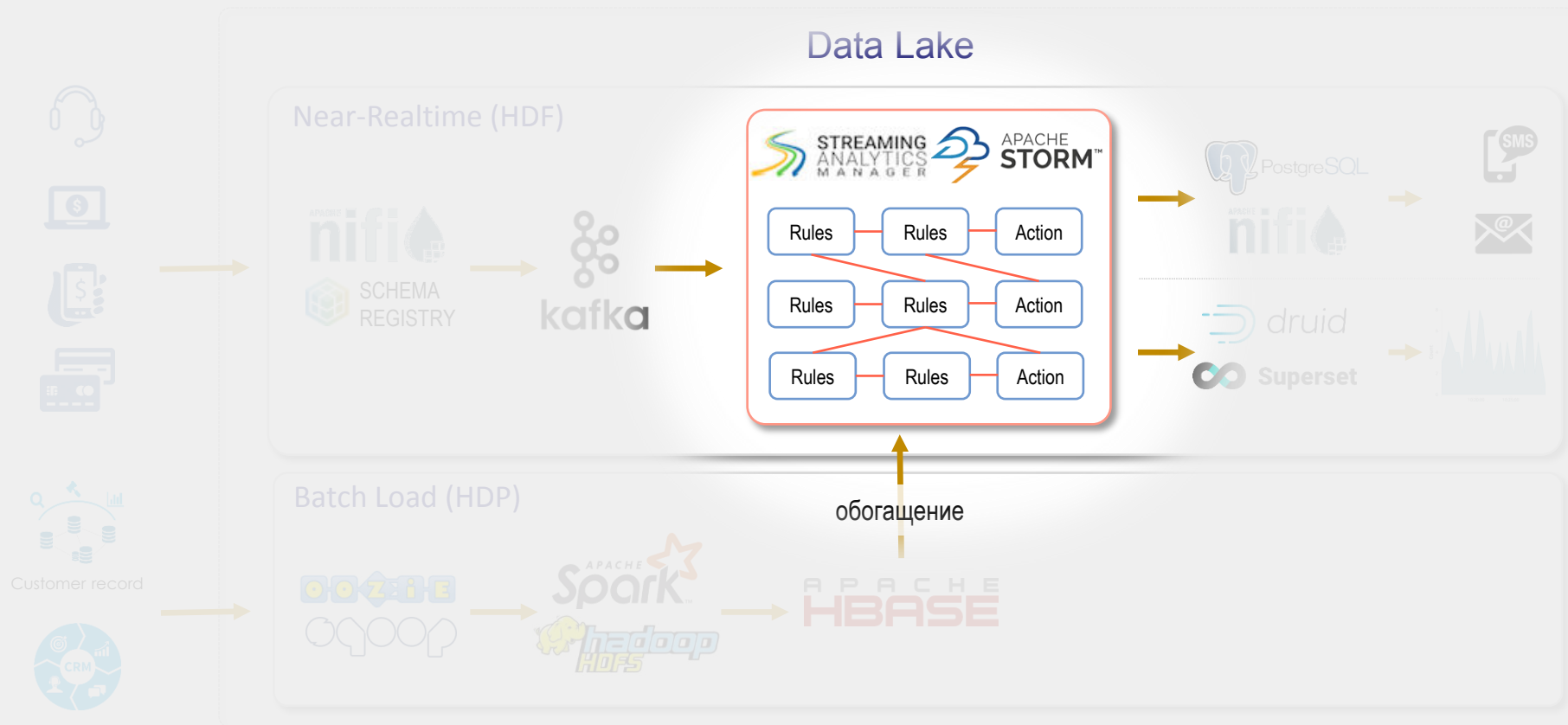
- **SAM**

- читаем схему "стрима" (метаданные)

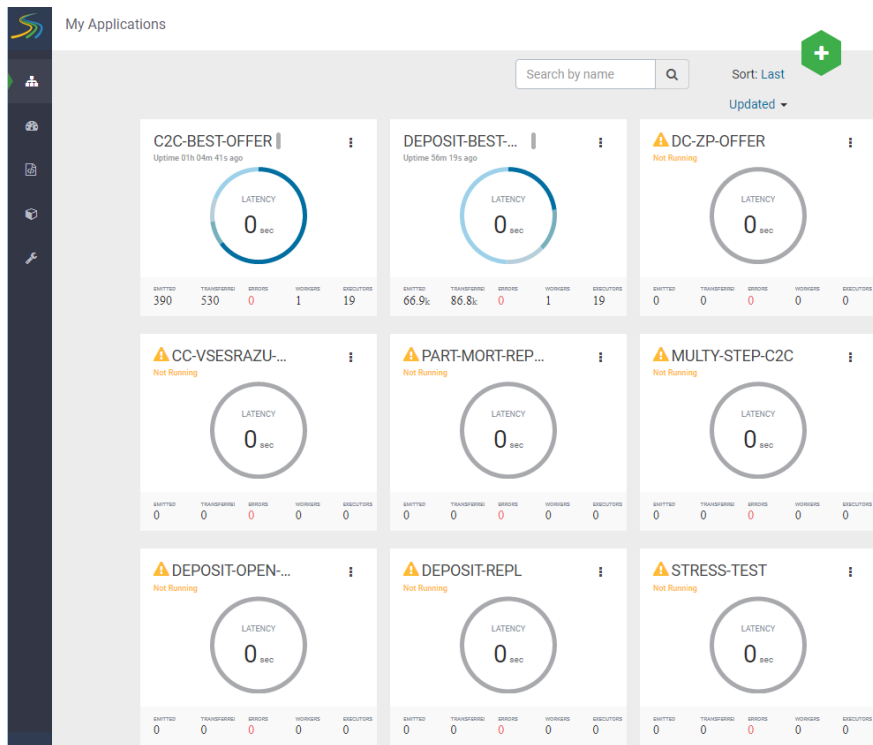
- **Storm**

- де-сериализуем "стрим" в "данные"

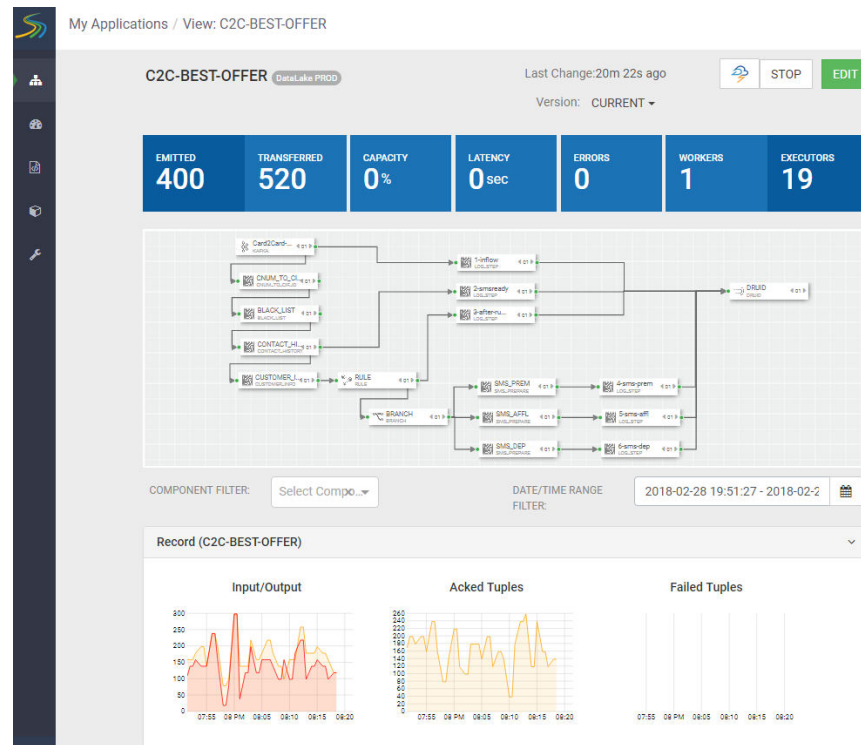
Архитектура RTDM-решения



Streaming Analytics Manager (SAM)

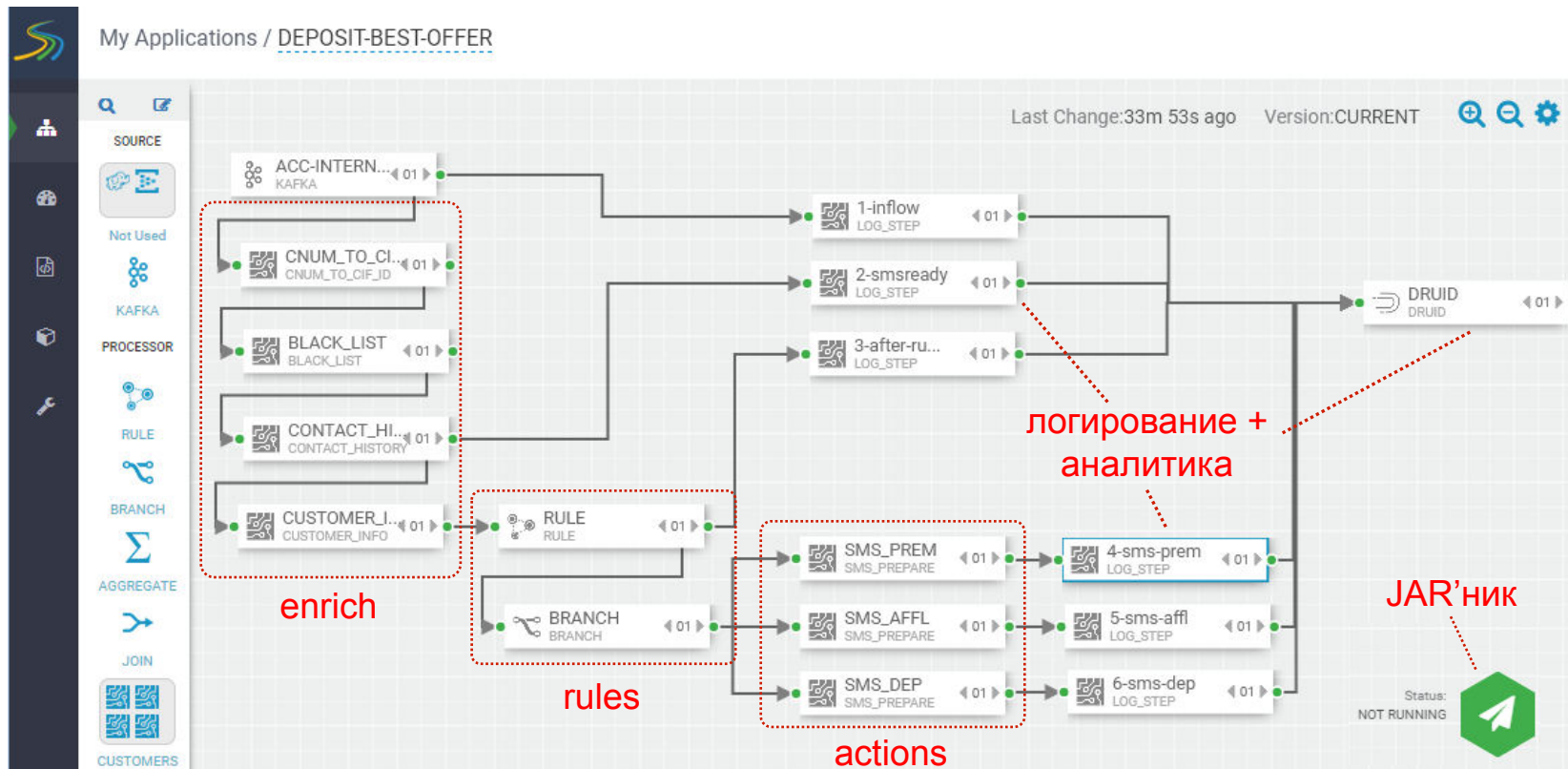


Список кампаний



Сводная информация по кампании

Streaming Analytics Manager (SAM)



Storm (via SAM)

Are you sure want to continue with this configuration?

GENERAL ADVANCED

NUMBER OF WORKERS
2

NUMBER OF ACKERS
1

TOPOLOGY MESSAGE TIMEOUT (SECONDS)
20

TOPOLOGY WORKER JVM OPTIONS

Cancel Ok

Annotations:
=0 → At most once
→ At least once
~~Exactly once~~

Are you sure want to continue with this configuration?

GENERAL ADVANCED

FIELD NAME*	FIELD VALUE*	
topology.max.task.parallelism	10	+
topology.max.spout.pending	5	+ -
backpressure.disruptor.high.wate	0.8	+ -
backpressure.disruptor.low.water	0.5	+ -

Cancel Ok



• кол-во Executors

- Готовность к дублям
- Идемпотентность операций

- Крайне мало предустановленных процессоров
- Разработка кастомных процессоров

SAM + DRUID

My Applications / DEPOSIT-BEST-OFFER

1-C2C_IN

CONFIGURATION NOTES

Input

- requestId*
LONG
- clientId*
STRING
- requestSubtypeId*
INTEGER
- amount*
DOUBLE
- currency*
STRING
- cif_id*
STRING
- sch*
STRING

ИМЯ ШАГА *

1-ALL_INCOME_C2C

INPUT SCHEMA MAPPING

OUTPUT FIELDS*

- x requestId x clientId x amount x currency x cif_id
- x log_action

Output

- requestId*
LONG
- clientId*
STRING
- amount*
DOUBLE
- currency*
STRING
- cif_id*
STRING
- log_action*
STRING

Cancel Ok

Last Change: 33m 53s ago Version: CURRENT

логирование + аналитика

DRUID DRUID

4-sms-prem LOG_STEP

5-sms-affl LOG_STEP

6-sms-dep LOG_STEP

Druid + Superset: Аналитика и Дашборды

Superset

Sources

Slices

Dashboards

Refresh Save User

Deposit-Best-Offer ★

bo_metrics_pre

1.94k

+24.0%

bo_metrics_pre5

682

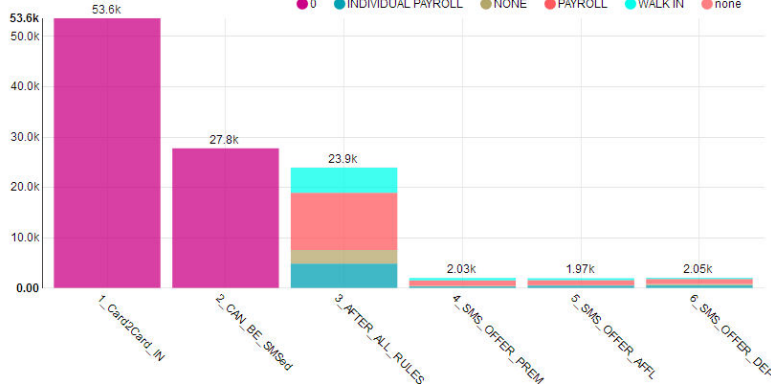
-75.7%

bo_metrics_sal

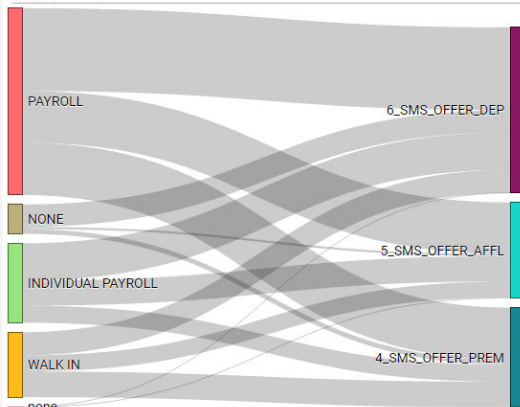
472

-45.4%

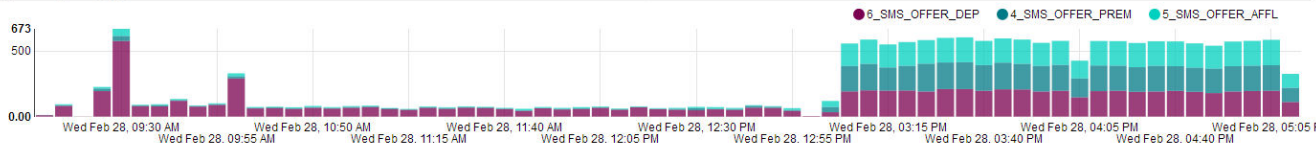
bo_step_cnts



bo_payroll_to_sms



bo_sended_sms_by_time

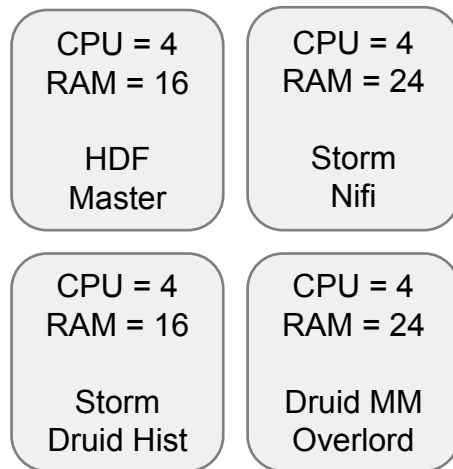


Druid (by Hortonworks?)

- Druid Historical
- Druid Broker
- Druid Coordinator
- Druid Router
- Druid Overlord
- Druid MiddleManage



HDFS



- ✓ Иногда может зависнуть без видимых причин, загружая поток со скоростью 1-2 события в секунду
- ✓ Есть подозрение, что это связано с тем что не хватает ресурсов на индексацию
- ✓ Открывается слишком много “временных окон” для индексации? – ~~не~~ факт

Развертывание в Enterprise

SAM -> Maven

Собирает JAR-файл для Storm'a

при каждом запуске SAM-кампании



доступ в интернет
с Prod-серверов

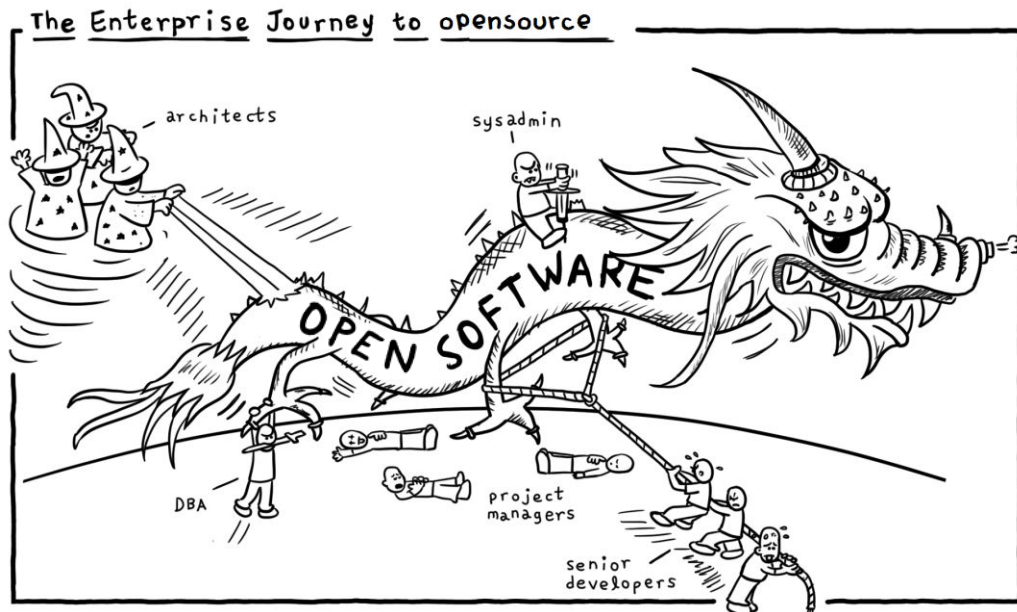
решение:

1. скачать библиотеки и сложить в локальную папку
2. патч SAM (переписать класс)

Hortonworks (HDP+HDF)

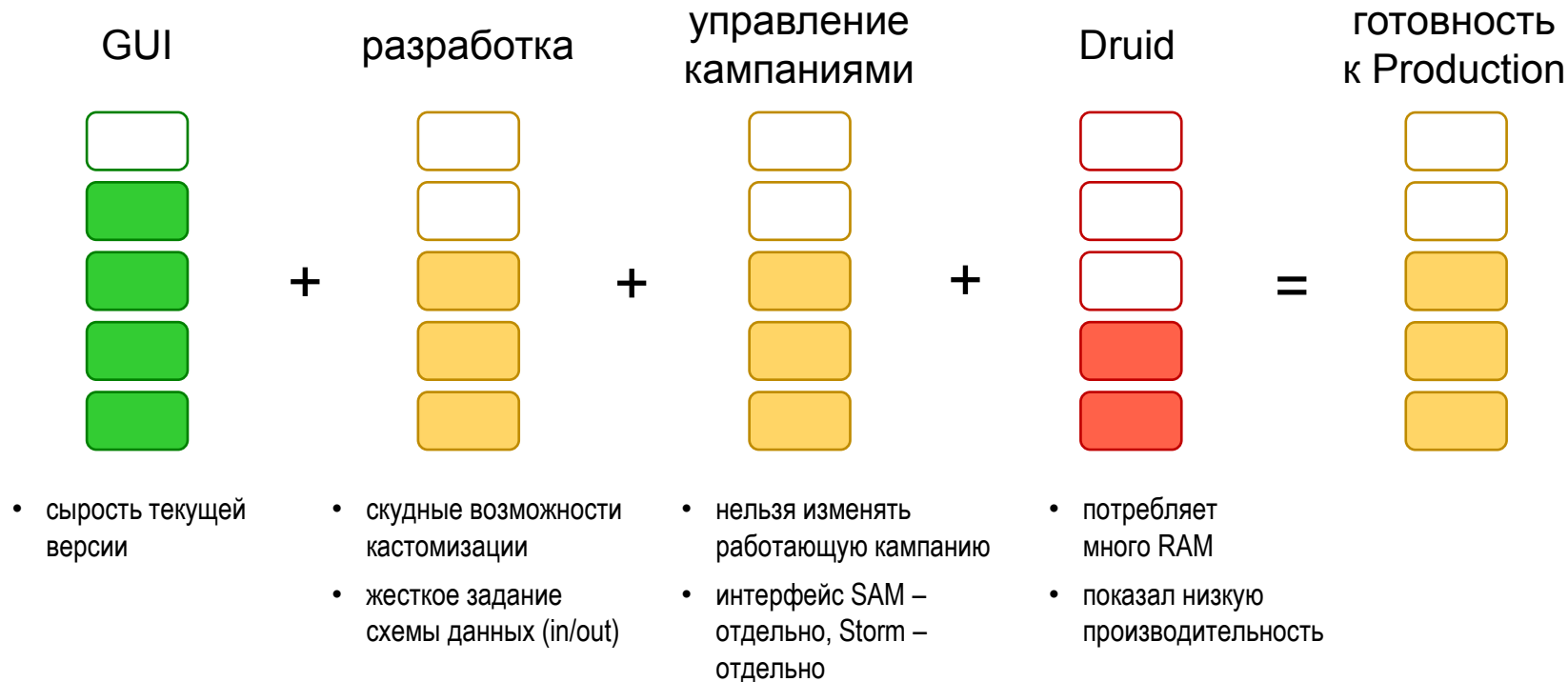
Каждое обновление платформы "как в первый раз":

- Необходимых для работы конфигов/JAR-файлов может не быть в поставке
- Регулярные конфликты старых и новых версий библиотек (вычищаем руками)



Pic source: <https://builttoadapt.io/enterprise-journey-to-the-cloud-aa9e9024a976>

Опыт HDF после 2-х недель пилота (и 3-х месяцев развертывания)



Выводы

ОНО реально РАБОТАЕТ !

НО, Вам придется:

- запастись терпением
- успокоительными
- ввести запрет на мат на рабочем месте

ЧТО, можно сделать лучше:

- DRUID → ClickHouse
- Storm → Spark Structured Streaming
- Автоматическое логирование



Спасибо!

aleksey.kuznetsov@raiffeisen.ru

anton.yarmolyuk@raiffeisen.ru