



Spark Streaming

SELEZNEV ARTEM
HEAD OF CVM ANALYTICS @ MAGNIT

WELCOME
STREAMS



STREAMING

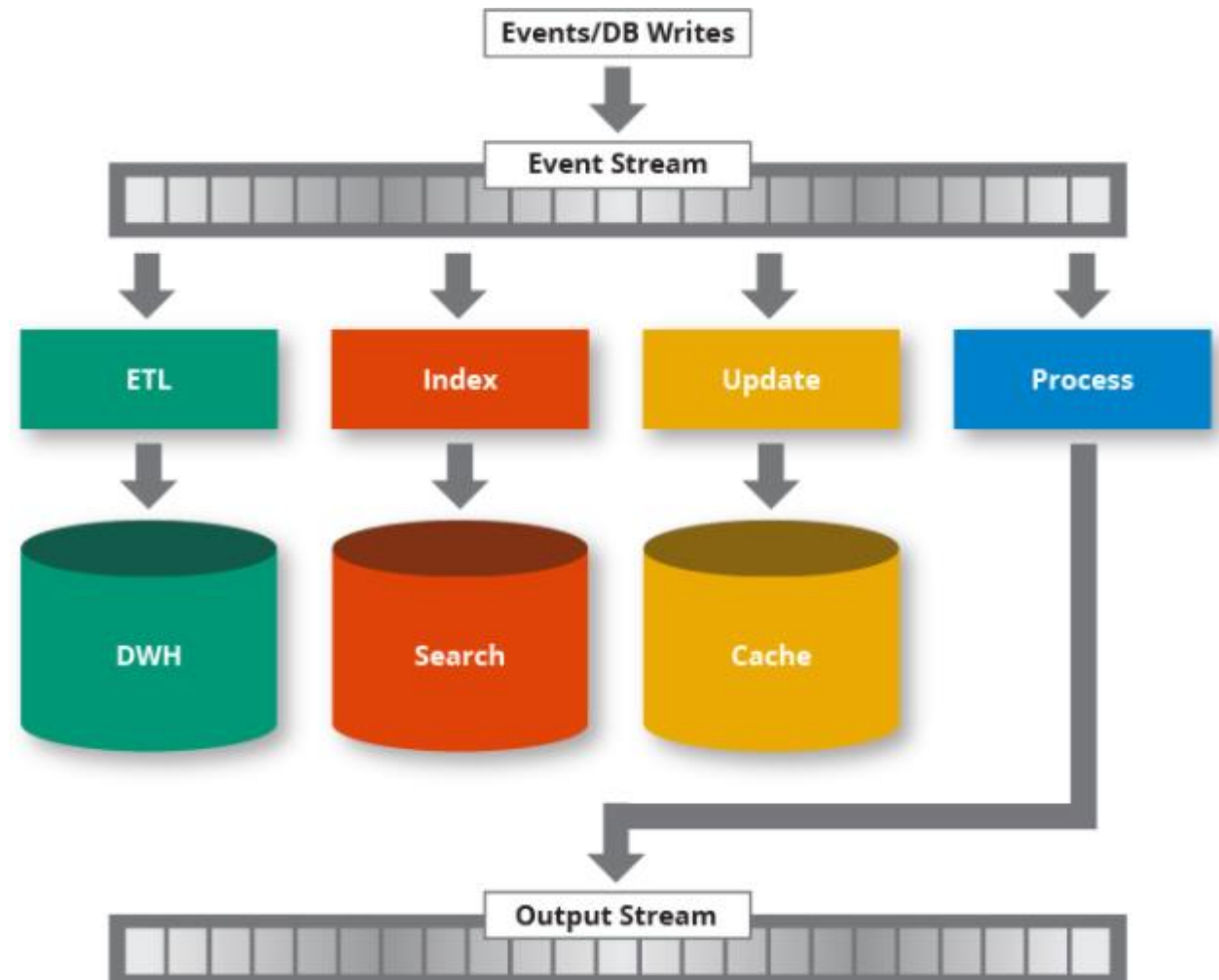


STREAMING



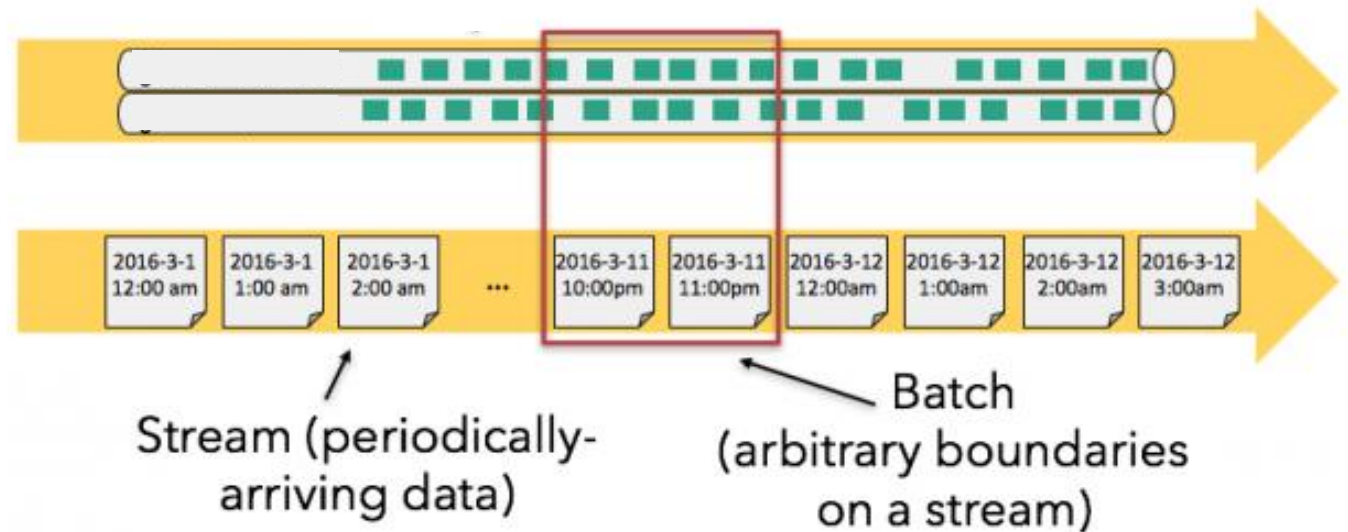
STREAMING TYPES | EVENT BASED

- Один event в промежуток времени
- Параллельные, но независимые
- Задержка ~10ms



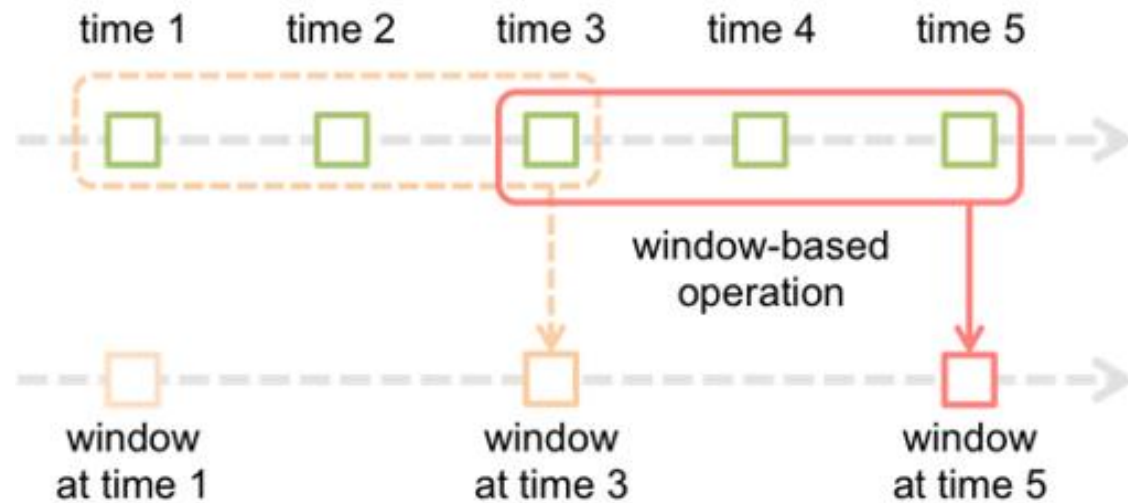
STREAMING TYPES | MICRO-BATCH

- Зависимость от размера батча
- Батчи выполняются последовательно
- Задержка $\gg 1s$



STREAMING TYPES | WINDOWS

- Окно – вид батча



STREAMING – ОСНОВНАЯ ИДЕЯ

```
from pyspark.sql import SparkSession
from pyspark.sql.functions import explode
from pyspark.sql.functions import split

spark = SparkSession.builder.appName("WordCount").getOrCreate()

lines = spark.readStream.format("socket").option("host",
"localhost").option("port", 9999).load()

words = lines.select(explode(split(lines.value, " ")) .alias("word"))

wordCounts = words.groupBy("word").count()
```

```
$ spark-submit wordcount.py localhost 9999
```

```
Batch: 0
```

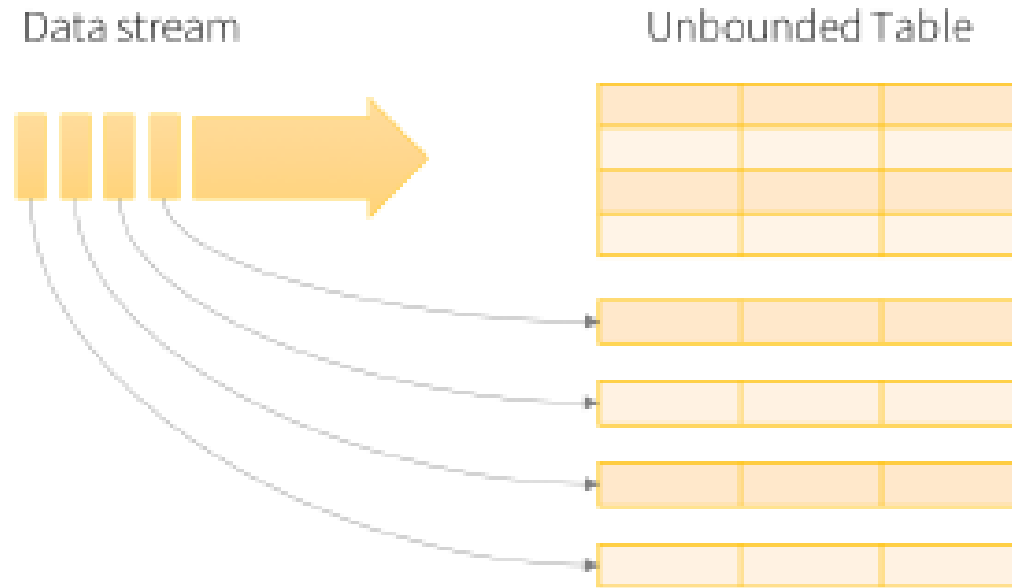
value	count
apache	1
spark	1

```
Batch: 1
```

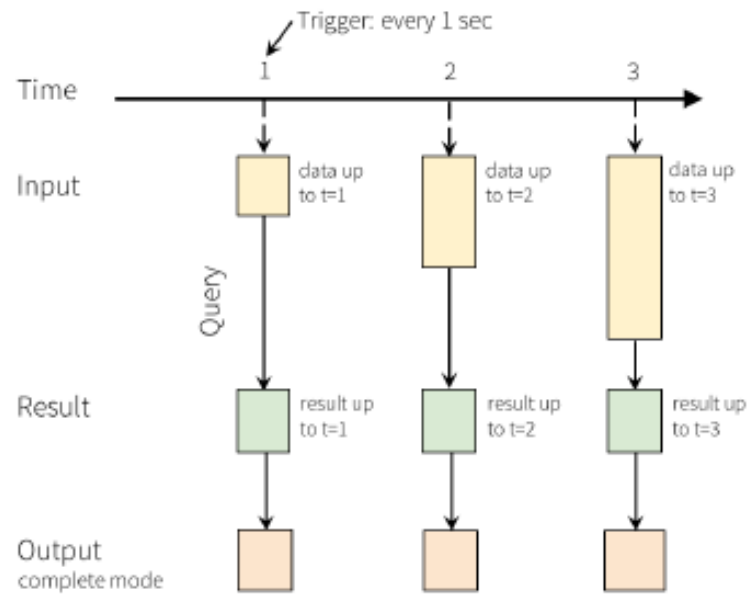
value	count
apache	2
spark	1
hadoop	1

```
---
```

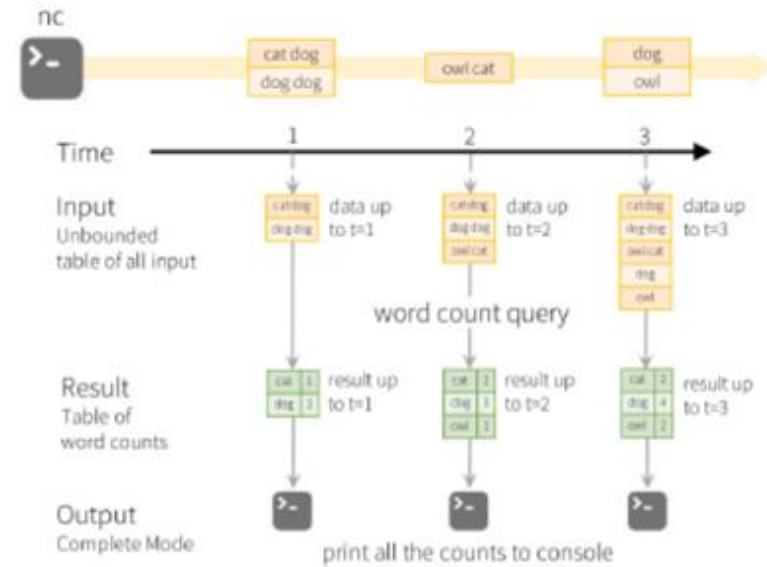
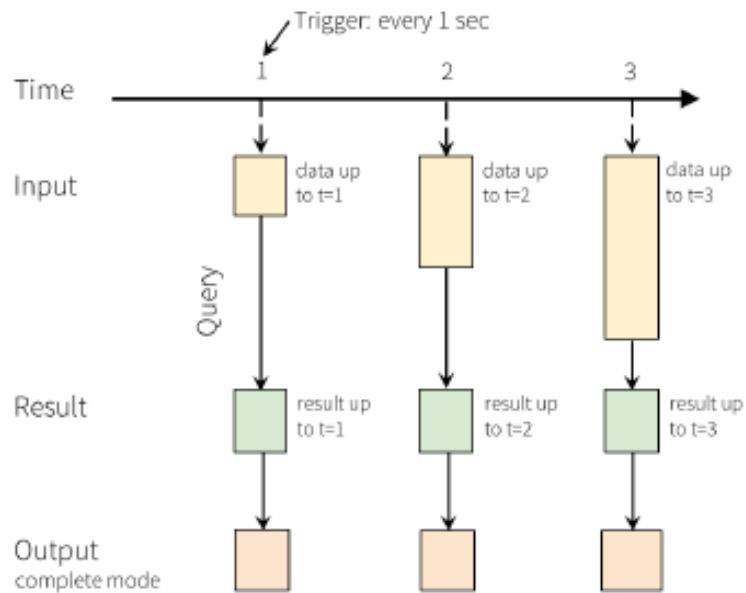

STREAMING – ОСНОВНАЯ ИДЕЯ



STREAMING – ОСНОВНАЯ ИДЕЯ



STREAMING – ОСНОВНАЯ ИДЕЯ

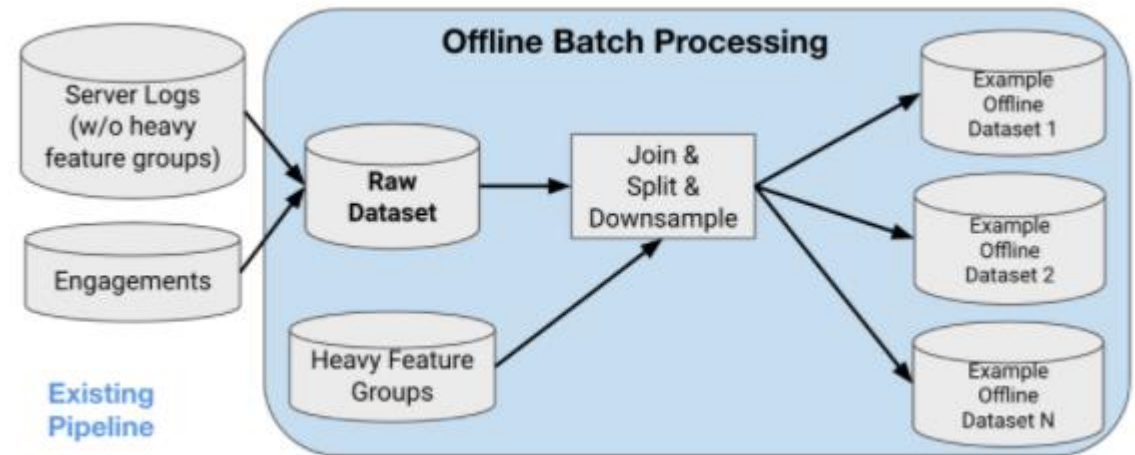


STREAMING CASE



CASE

- Мы можем получить TOP твитов в вашем аккаунте
- Эти твиты выводятся вам с помощью рекомендательной системы (онлайн)
- Сейчас модель переобучается 1 раз в 7 дней, задача переобучать 1 раз в день



CASE

- Почему надо переобучить?

Score падает с каждым днем!

	T-4	T-3	T-2	T-1
<i>RCE (Relative Cross Entropy)</i>	38.81	38.82	38.84	39.12

RCE?

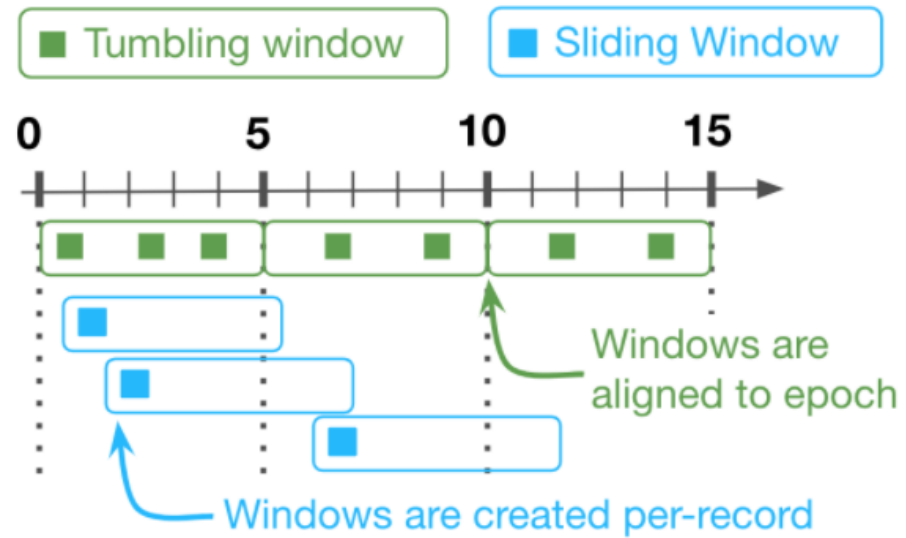
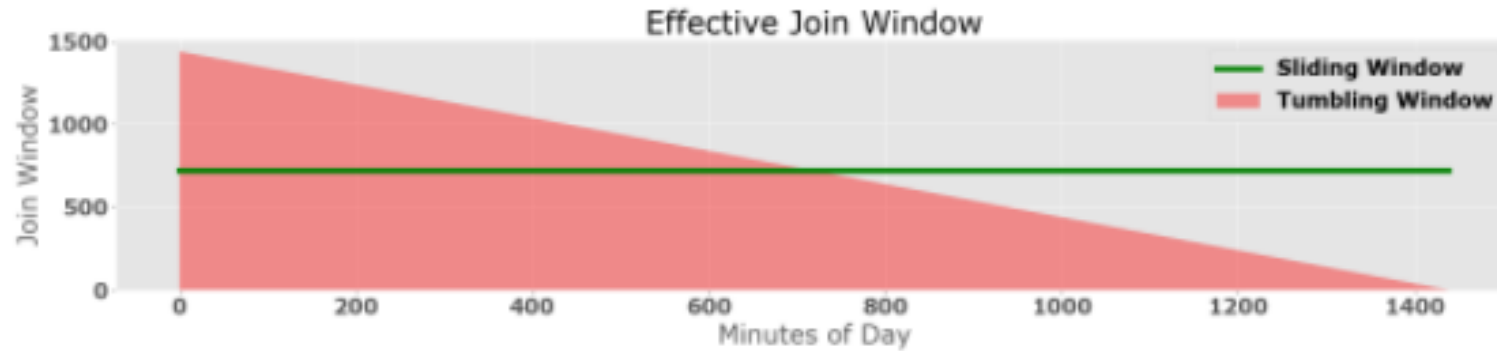
CASE

- Почему надо переобучить?
Score падает с каждым днем!

	T-4	T-3	T-2	T-1
<i>RCE (Relative Cross Entropy)</i>	38.81	38.82	38.84	39.12

rce — на сколько предсказание лучше, чем базовая модель «без фич по пользователю»
(straw man of naïve prediction)

CASE



ПОПРОБУЕМ
САМОСТОЯТЕЛЬНО

ТЫ НЕ ДЕЛАЕШЬ ЭТО НЕПРАВИЛЬНО



**ЕСЛИ НИКТО НЕ ЗНАЕТ, ЧТО
КОНКРЕТНО ТЫ ДЕЛАЕШЬ**

SPARK ПРОЕКТ



SPARK ПРОЕКТ

- Задачи реализуемые на Spark ML:
 - Создайте единый набор данных: акции + новости с Reddit
 - Преобразуйте текст в «фичи» на которых сможет обучиться алгоритм
 - Сделайте выбор окна для обучения: 3 дня, 7 дней, 14 дней, 21 день
 - Определите, что вы будете определять: фактическую цену / % изменения