# Distributed Time Travel for Feature Generation

Prasanna Padmanabhan

DB Tsai

Mohammad H. Taghavi

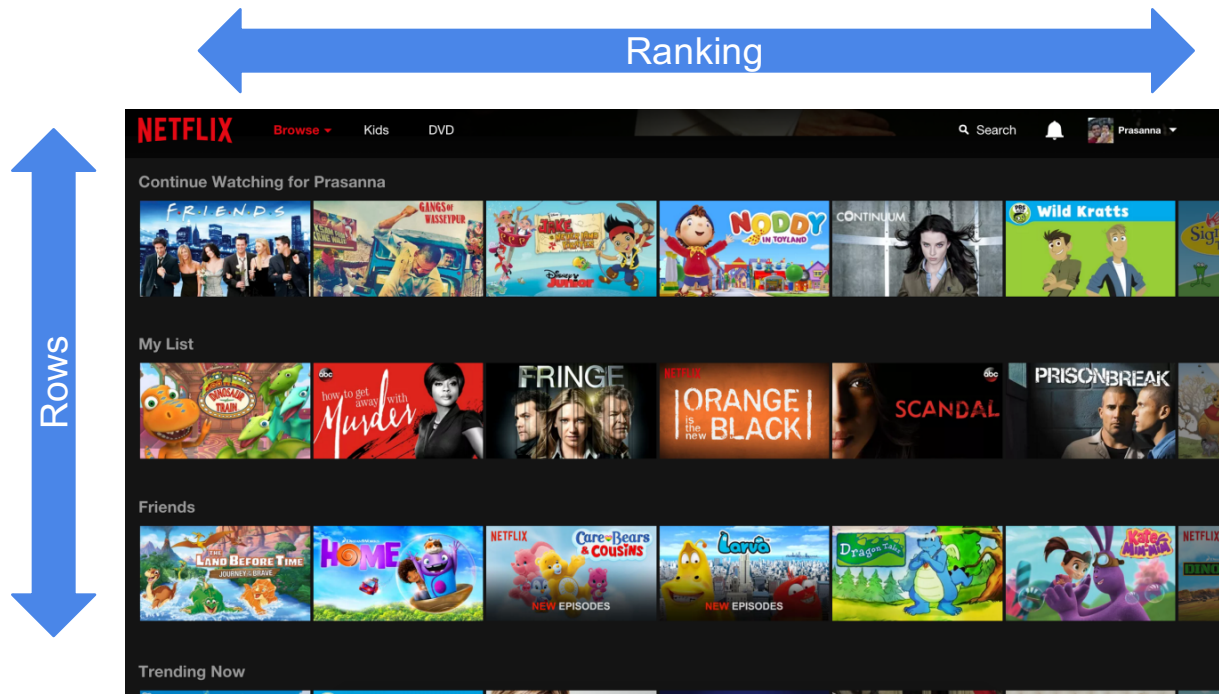**NETFLIX**

Turn on Netflix, and the **absolute best content for you** would **automatically** start **playing**

NETFLIX

# Everything is a Recommendation



**Over 80%** of what members watch comes from our recommendations

Recommendations are driven by **Machine Learning Algorithms**

# Data Driven

- Try an **idea offline using historical data** to see if it would have made better recommendations

- If it did, deploy a live **A/B test** to see if it performs well in Production

NETFLIX

# **Why** build a Time Machine?

NETFLIX

**Quickly** try ideas on **historical data** and **transition** to online A/B test

NETFLIX

# The Past

- Generate features based on event data logged in Hive
  - Need to reimplement features for online A/B test
  - Data discrepancies between offline and online sources

- Log features online where the model will be used
  - Need to deploy each idea into production

- Feature generation calls online services and filters data past a certain time
  - Works only when a service records a log of historical events
  - Additional load on online services

NETFLIX

DeLorean image by JMortonPhoto.com & OtoGodfrey.com

NETFLIX

# Time Travel using Snapshots

- Snapshot online services and use the snapshot data offline to generate features


- Share facts and features between experiments without calling live systems

NETFLIX

# **How** to build a Time Machine

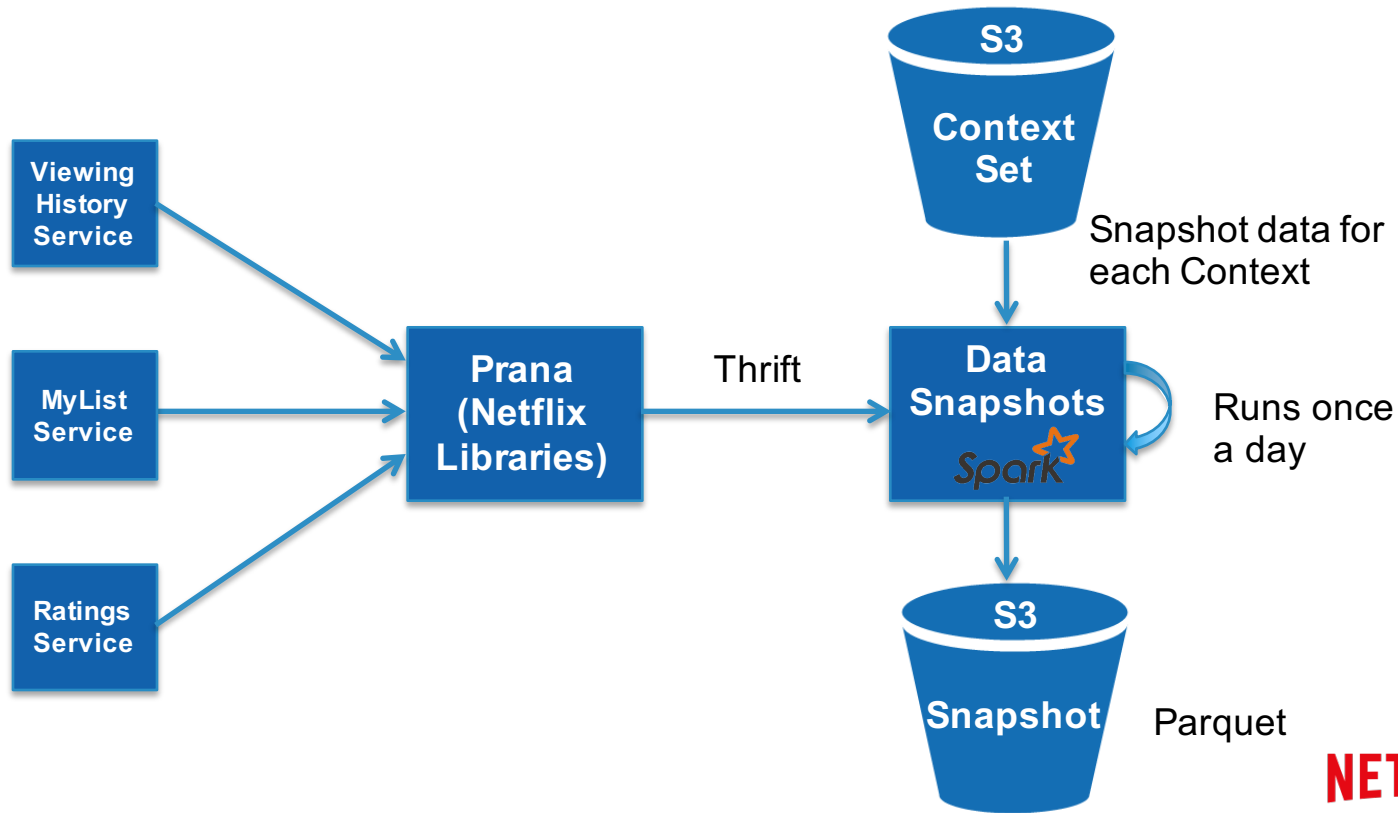NETFLIX

# Context Selection

# Data Snapshots

# APIs for Time Travel

# Context Selection

# Data Snapshots

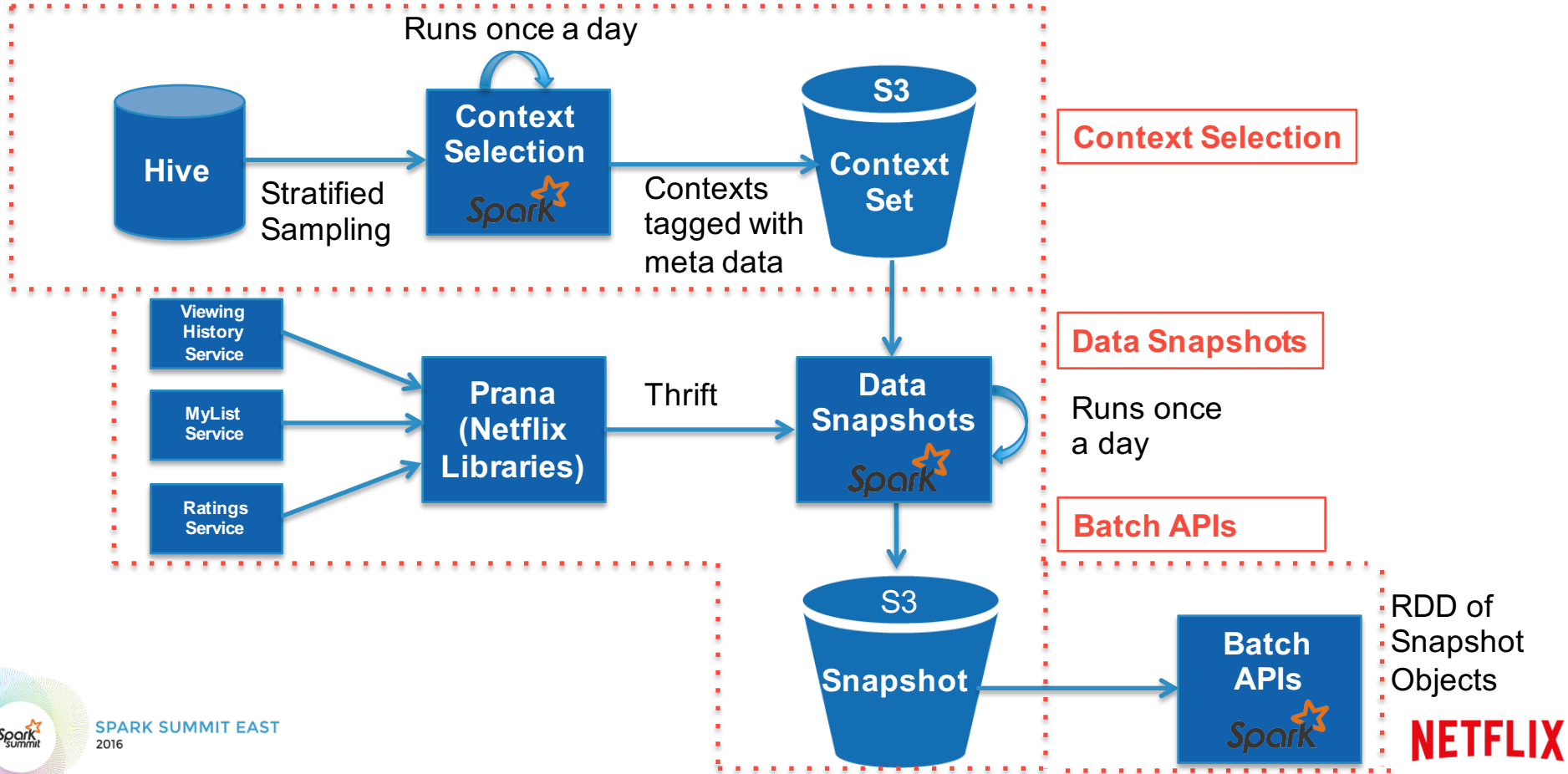# APIs for Time Travel

```
scala>  val snapshot = new SnapshotDataManager(sqlContext)
                       .withTimestamp(1445470140000L)
                       .withContextId(OUTATIME)
                       .getViewingHistory
snapshot: org.apache.spark.rdd.RDD[(Long, com.netflix.viewinghistory.ViewingHistory)]
```

NETFLIX

# Data Architecture

# Generating **Features** via **Time Travel**

NETFLIX

# Great Scott! There's the DeLorean!

- DeLorean: A time-traveling vehicle

  – uses data snapshots to travel in time

  – scales with Apache Spark

  – prototypes new ideas with Zeppelin

  – requires minimal code changes from
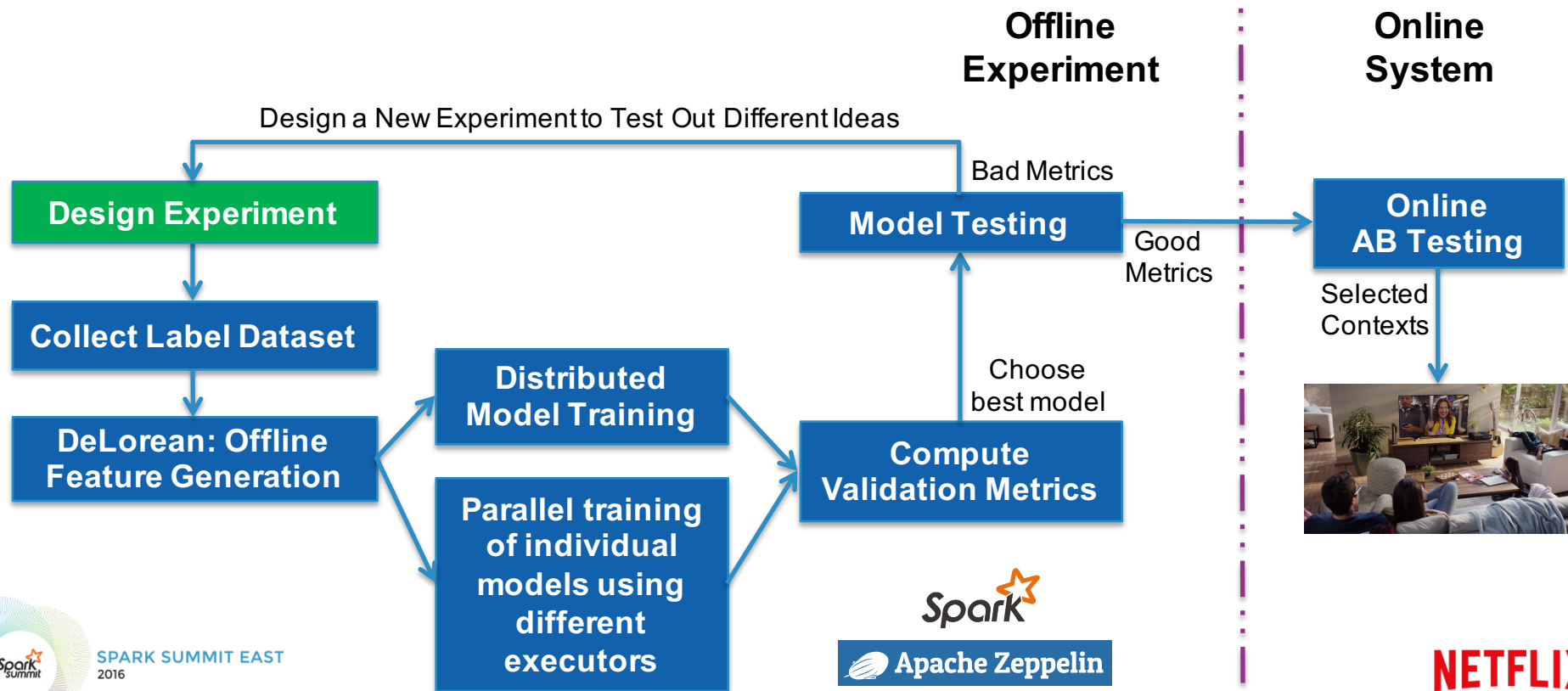     experimentation to A/B test to production

SPARK SUMMIT EAST
2016

# Running Time Travel Experiment

Select the **destination time**

Bring it up to **88 miles** per hour!

NETFLIX

# Running Time Travel Experiment

**Offline Experiment**

**Online System**

Design a New Experiment to Test Out Different Ideas

**Design Experiment**

**Collect Label Dataset**

**DeLorean: Offline Feature Generation**

**Distributed Model Training**

**Parallel training of individual models using different executors**

**Compute Validation Metrics**

Choose best model

**Model Testing**

Bad Metrics

Good Metrics

**Online AB Testing**

Selected Contexts



Spark

Apache Zeppelin

SPARK SUMMIT EAST 2016

NETFLIX

# DeLorean Input Data

- Contexts: The setting for evaluating a set of items (e.g. tuples of member profiles, country, time, device, etc.)

- Items: The elements to be trained on, scored, and/or ranked (e.g. videos, rows, search entities).

- Labels: For supervised learning, this will be the label (target) for each item.

NETFLIX

# Feature Encoders

- Compute features for each item in a given context

- Each type of raw data element has its own data key

- Data map is a map from data keys to data objects in a given context

- Data map is consumed by feature encoder to compute features

# Two type of Data Elements

- Context-dependent data elements

    - Viewing History

    - Mylist

    - ...

- Context-independent data elements

    - Video Metadata

    - Genre Metadata

    - ...

NETFLIX

# Feature Generation



Diagram components:

- **S3** / **Snapshot** (bucket) feeds into:
  - **Label Data**
  - **Data Elements**
- **Label Data** → *Data* →
- **Data Elements** → *Data in POJOs* →
- *Data Keys* ← Data Elements
- **Feature Model (JSON)** → *Required Feature Keys* →
- Central processing (image) outputs:
  - *Data Map* → **Feature Encoders**
  - *Features* ← **Feature Encoders**
  - *Data Keys* ← Feature Encoders
  - *Label Features* → **Model Training**

# Features

- Represented in Spark's DataFrames

- In nested structure to avoid data shuffling in ranking process
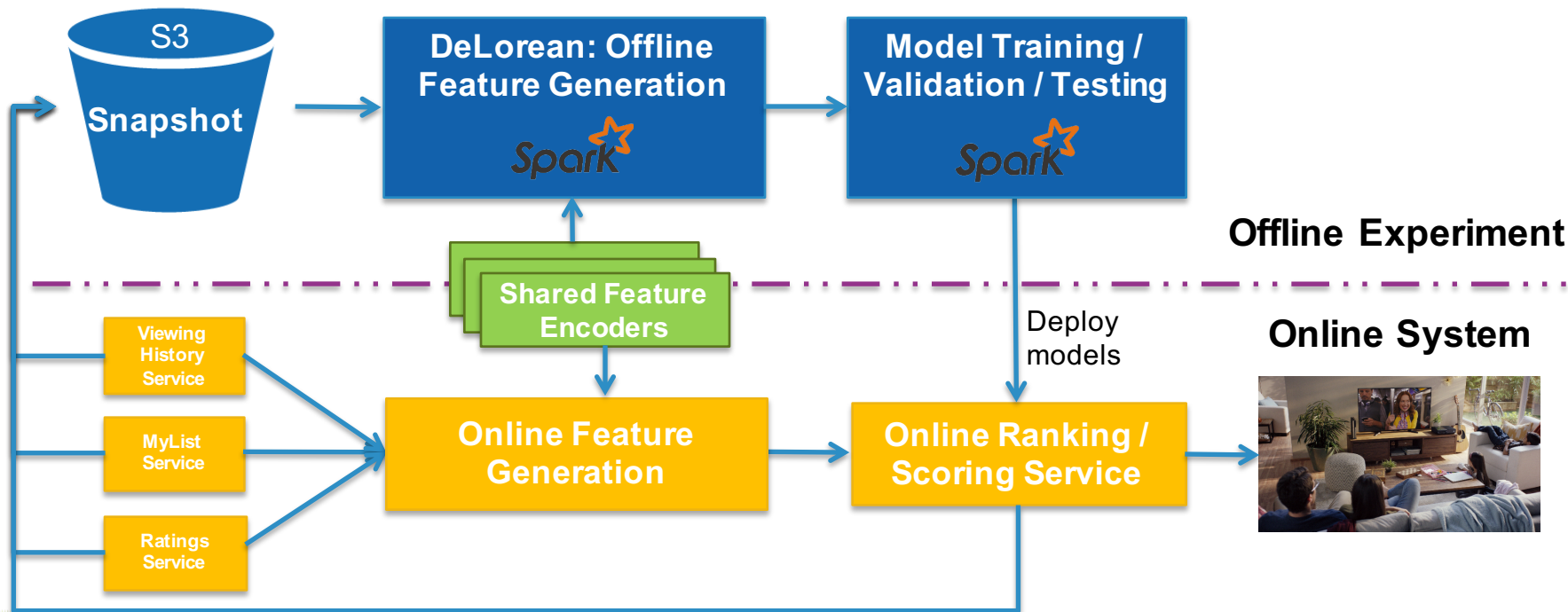
- Stored with Parquet format in S3

NETFLIX

# Features

Context

Item, label,
and features

```
root
 |-- Visitor: long (nullable = false)
 |-- Country: string (nullable = false)
 |-- data: array (nullable = true)
 |    |-- element: struct (containsNull = true)
 |    |    |-- videoId: long (nullable = false)
 |    |    |-- weight: double (nullable = false)
 |    |    |-- label: double (nullable = false)
 |    |    |-- features: struct (nullable = false)
 |    |    |    |-- rating: double (nullable = false)
 |    |    |    |-- unpersonalizedPopularity: double (nullable = false)
 |    |    |    |-- ...
 |    |    |    |-- ...
 |    |    |    |-- ...
```

NETFLIX

# Going Online

# Conclusion

Spark helped us <span style="color:red">significantly reduce</span> the time from an idea to an AB Test

NETFLIX

# Future work

Event Driven Data Snapshots

Time Travel to the Future!!

NETFLIX

# We're hiring!
## (come talk to us)
## **https://jobs.netflix.com/**

## **Tech Blog: http://bit.ly/sparktimetravel**

NETFLIX