# Personalizing Netflix with Streaming datasets

Shriya Arora

Senior Data Engineer
Personalization Analytics

@shriyarora

NETFLIX

**What is this talk about ?**

- **Helping you decide if a streaming pipeline fits your ETL problem**

- **If it does, how to make a decision on what streaming solution to pick**

**What is this NOT talk about ?**

- **X streaming engine is the BEST, go use that one!**

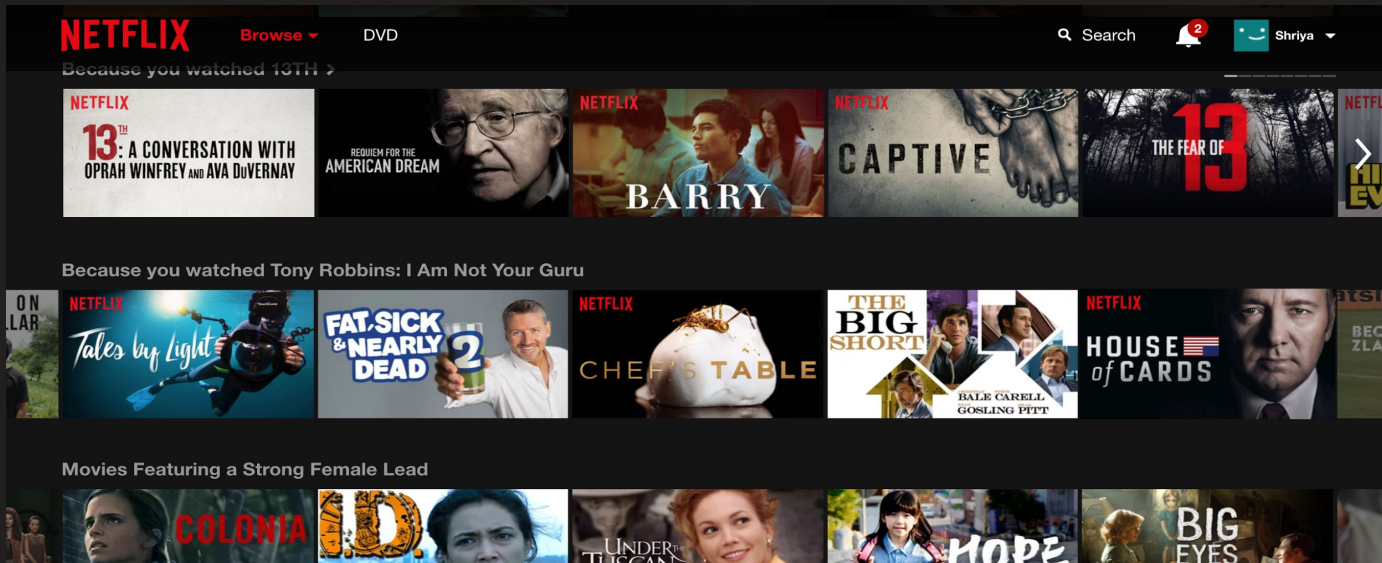- **Batch is dead, must stream everything!**

NETFLIX

# What is Netflix's Mission?

**Entertaining you by allowing you to stream content anywhere, anytime**

NETFLIX

# What is Netflix's Mission?

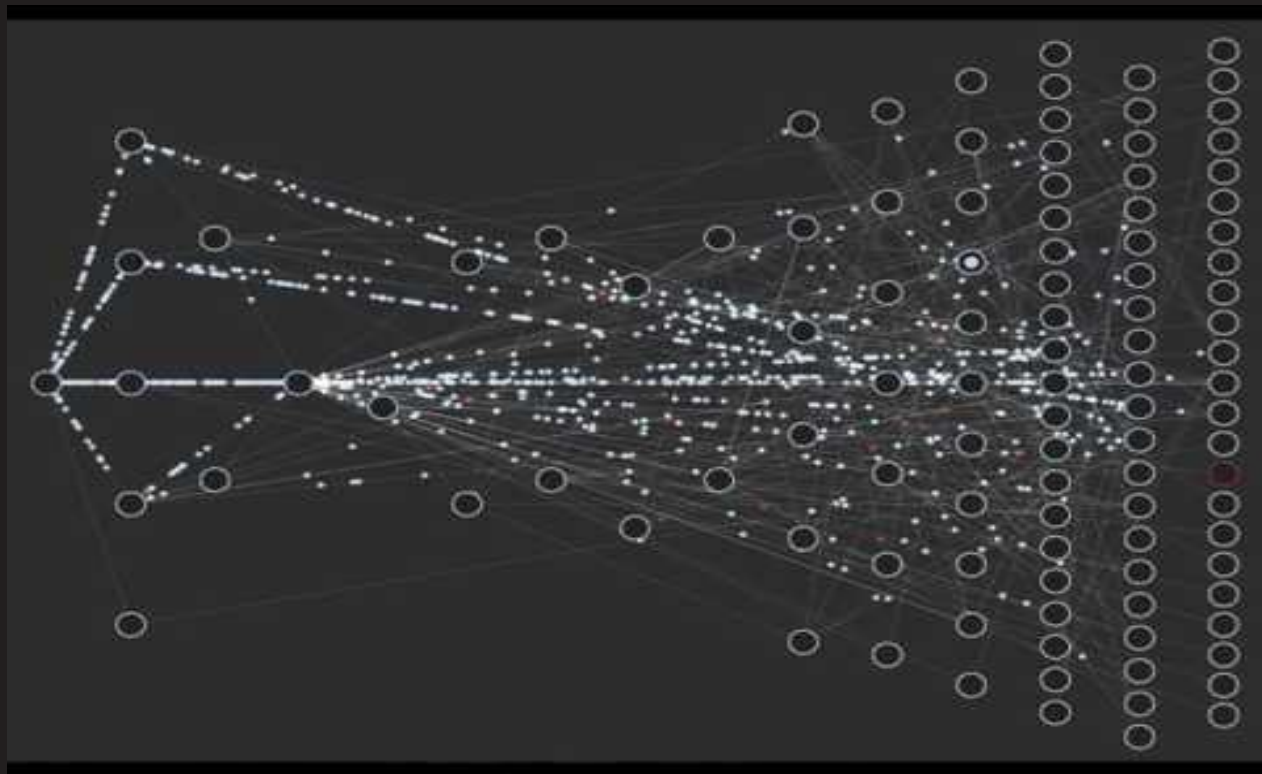**Entertaining you by allowing you to stream personalized content anywhere, anytime**

# How much data do we process to have a personalized Netflix for everyone?

- 100**M+** active members

- 125**M** hours/ day

- 190 countries with unique catalogs

- 450**B** unique events/day

- 700+ Kafka topics

HOUSE of DATA

Image credit:http://www.bigwisdom.net/

# DEA Personalization at a (very) high level



User watches a
video on Netflix

Data flows through
Netflix Servers

DEA

DATA ENGINEERING & ANALYTICS

NETFLIX

NETFLIX
RESEARCH

NETFLIX

# Data Infrastructure

# Why have data later when you can have it now?

# Business wins
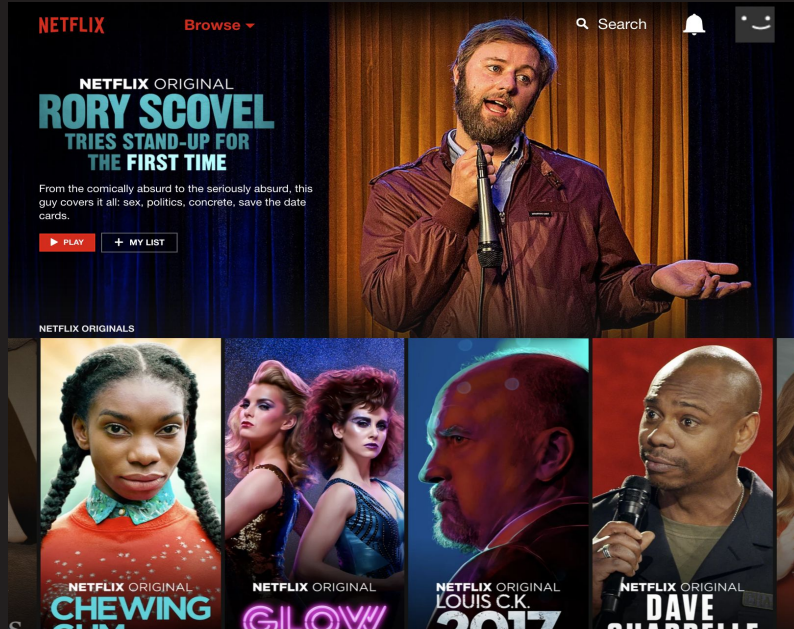
- Algorithms can be trained with the latest data



NETFLIX

# Business wins

- Innovation in marketing of new launches



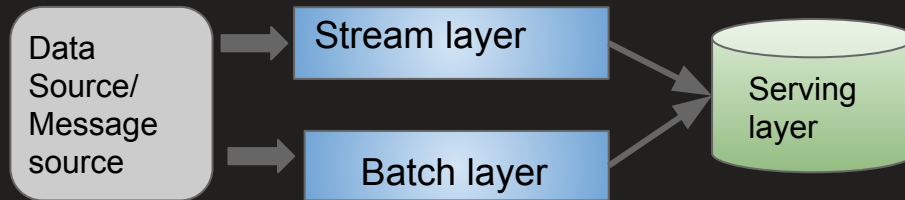- Creates opportunity for news kinds of algorithms

NETFLIX

# Technical wins

- ## Save on storage costs
  - Raw data in its original form has to be persisted

- ## Faster turnaround time on error correction
  - Long-running batch jobs can incur significant delays when they fail

- ## Real-time auditing on key personalization metrics

- ## Integrate with other real-time systems
  - Additional infrastructure is required to make 'online' systems be available offline

NETFLIX

# How to pick a Stream Processing Engine?

Problem Scope/Requirements

- Event-based streaming or micro-batches?

- What features will be the most important for the problem?

- Do you want to implement Lambda?

# How to pick a Stream Processing Engine?

Existing Internal Technologies
- Infrastructure support: What are other teams using?
- ETL eco-system:  Will it fit in with the existing sources and sinks

What's your team's learning curve?
- What do you use for batch?
- What is the most fluent language of the team?

# Our problem: Source of Play / Source of Discovery
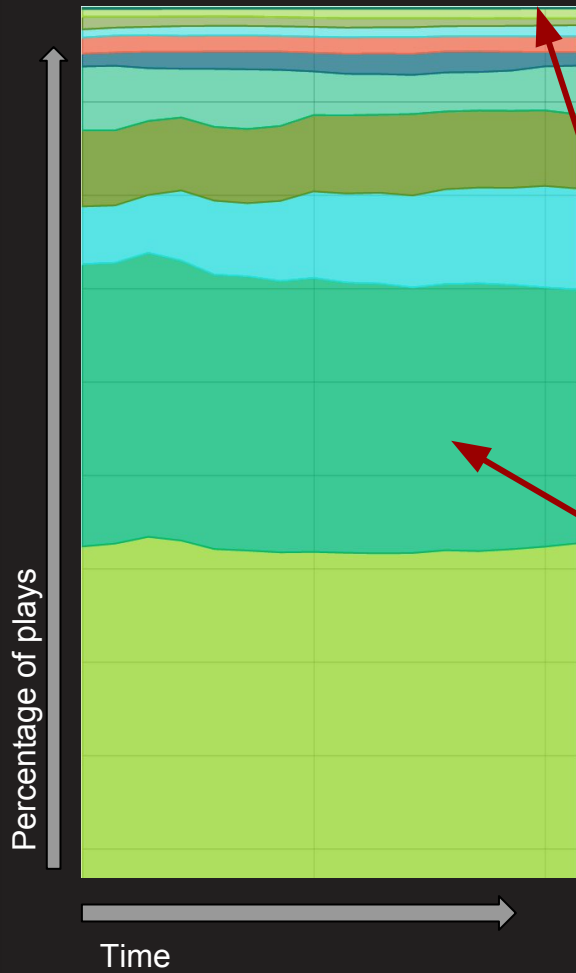


**Anatomy of a Netflix Homepage:**

Billboard

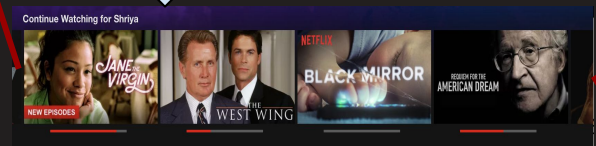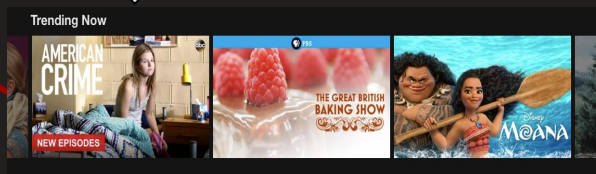Video Rankings (ordering of shows within a row)

Rows

**Source of Discovery**

**Source of Play**

Continue Watching

Trending now

Percentage of plays

Percentage of plays

Time

Time
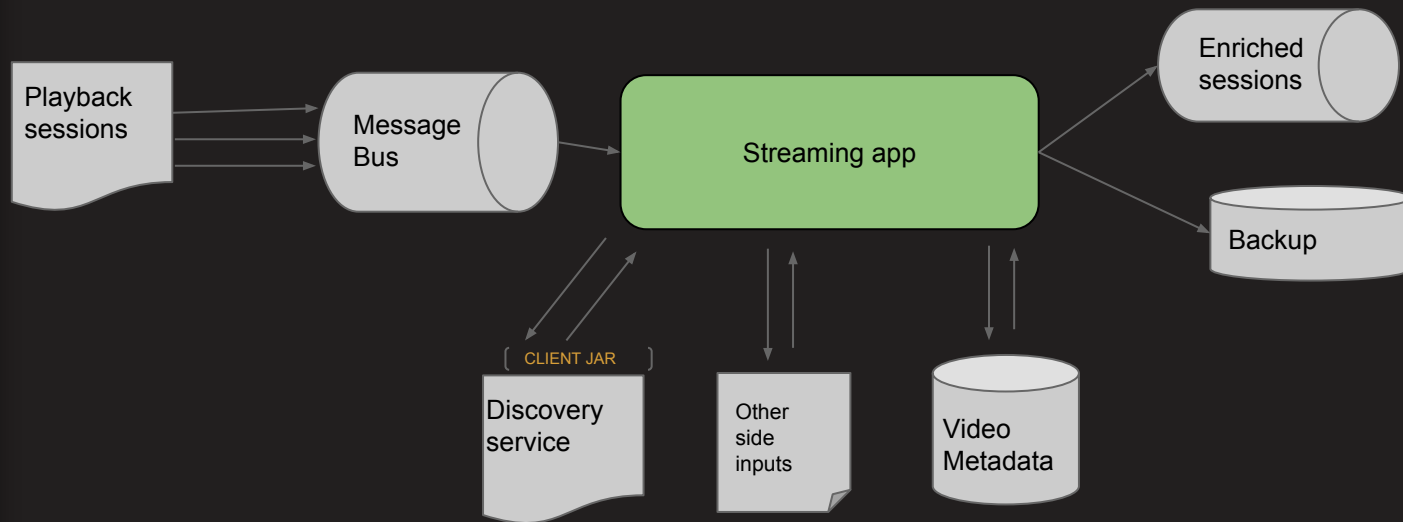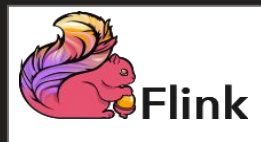
**What we need to solve for Source of Discovery:**

- High throughput
    - ~100M events/day
- Talk to live micro-services via thick clients
- Integrate with the Netflix platform eco-system
- Small State
- Allow for side inputs of slowly changing data

NETFLIX

# Source-of-Discovery pipeline: Data Flow

# Source-of-Discovery pipeline: Tech stack
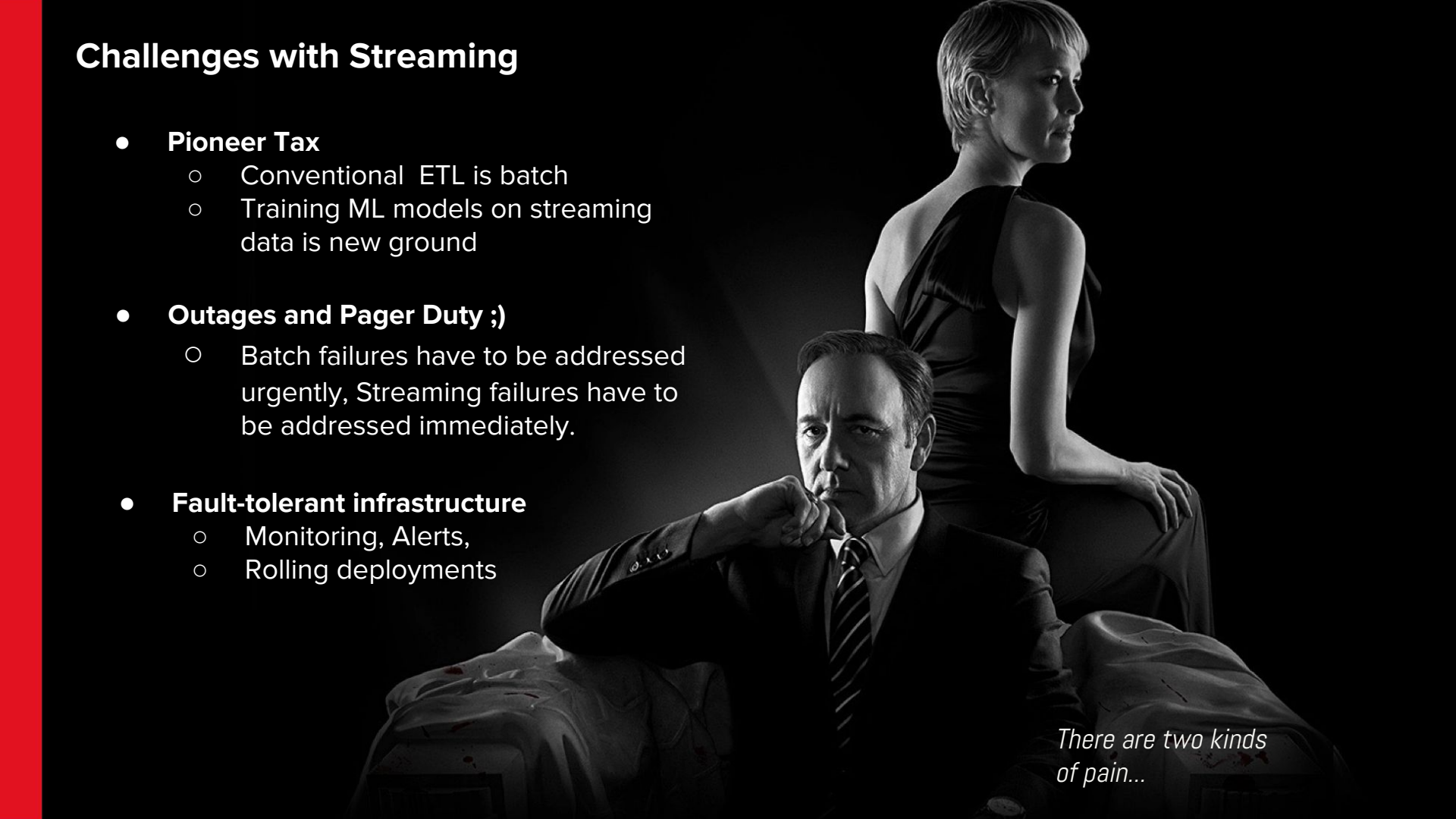
## Getting streaming ETL to work

- Getting Data from Live sources
  - Every event (session) enriched with attributes from past history
  - Making a call to the micro-service via a thick client

- Side inputs
  - Get metadata about shows from the content service
  - Slowly changing data, optimize to call less frequently

- Dependency Isolation
  - Shading jars is fun *(said no one ever)*

NETFLIX

# Getting streaming ETL to work cont..

- Data Recovery
    - Kafka TTLs are aggressive
    - Raw data stored in HDFS for finite time for replay

- Out of order events
    - Late arriving data must be attributed correctly

- Increased Monitoring, Alerts
    - Because recovery is non-trivial, prevent data-loss

NETFLIX

# Challenges with Streaming

- **Pioneer Tax**
  - Conventional ETL is batch
  - Training ML models on streaming data is new ground

- **Outages and Pager Duty ;)**
  - Batch failures have to be addressed urgently, Streaming failures have to be addressed immediately.

- **Fault-tolerant infrastructure**
  - Monitoring, Alerts,
  - Rolling deployments

*There are two kinds of pain...*

# Questions?

Stay in touch!

@**NetflixData**

NETFLIX