

Data Quality & Verification

@

Uber, Netflix & Amazon

Active Learning

↳ Participation & Discussion

Let's solve it together . . .

Sources:

Uber DQM



1. <https://eng.uber.com/monitoring-data-quality-at-scale/> (2020)
2. https://databricks.com/session_na20/an-approach-to-data-quality-for-netflix-personalization-systems (2020)
3. <https://www.amazon.science/publications/automating-large-scale-data-quality-verification> (2018)

Scale @
all companies

~ 1 - 10PB

10K - 100K tables

100s of Millions of events

1000's of metrics

(Q) What platform(s) & techniques
would you use to process internet-scale
data?

Data Quality impacts all downstream tasks



Data Analysis

Biz decision-making

ML & DL models

Infra - costs

Let's map Data Quality to an ML - problem?

Any suggestions?

Data Quality

- Time-series of metrics
(detect anomalies)

KPI

utilization

Table

rows

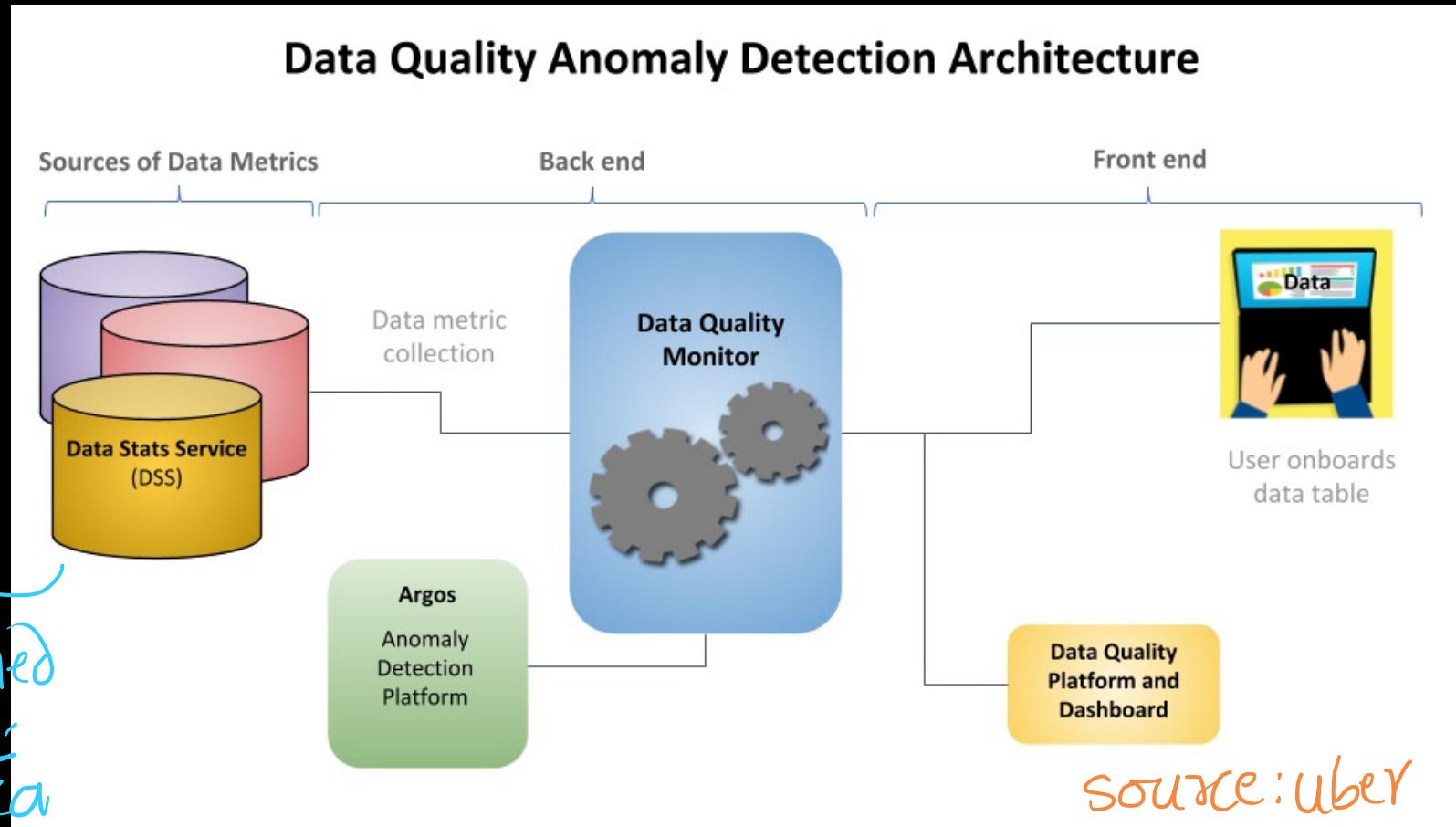
column-level

missing
values

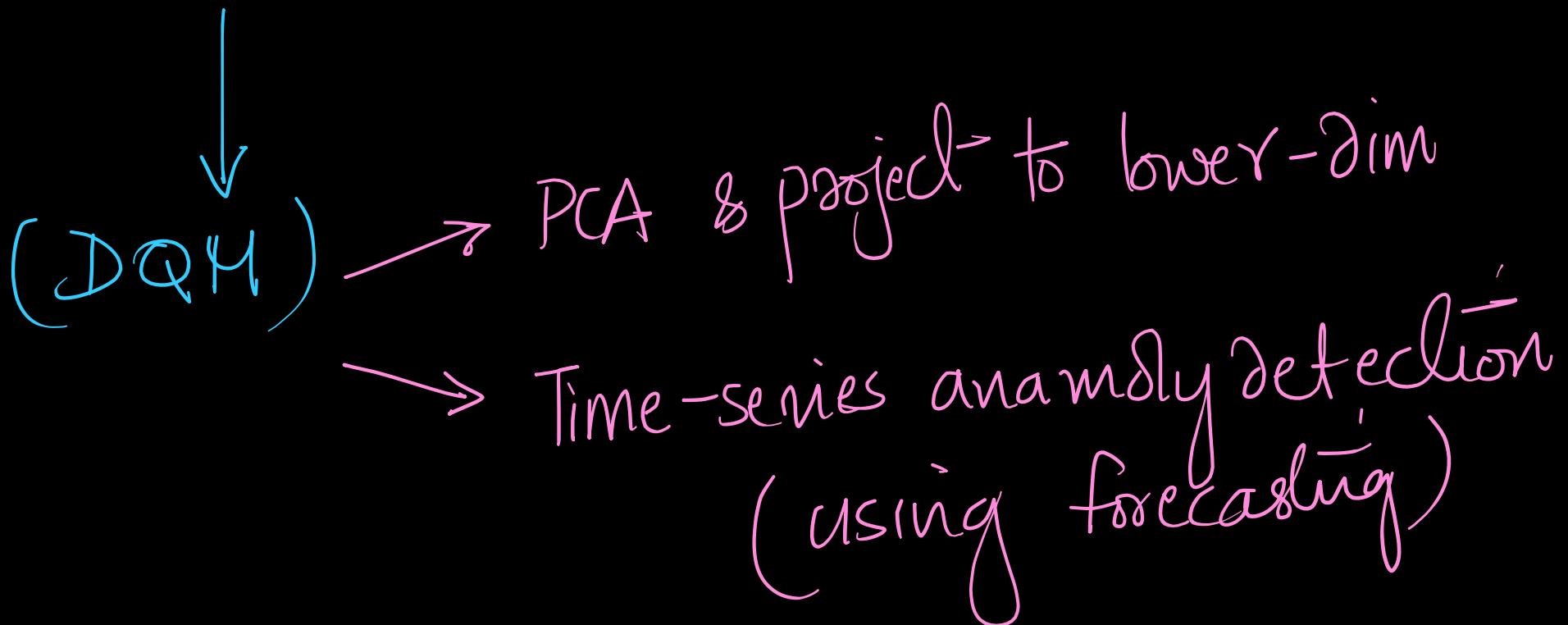
min; max; unique-values

→ Data Anomalies $\mathcal{D} = \{x_i, y_i\}_{i=1}^n$

(Q) What techniques can we employ here ?



DSS: multi-dimensional time series output



(Q) why do you PCA is helpful here?

many metrics are highly correlated

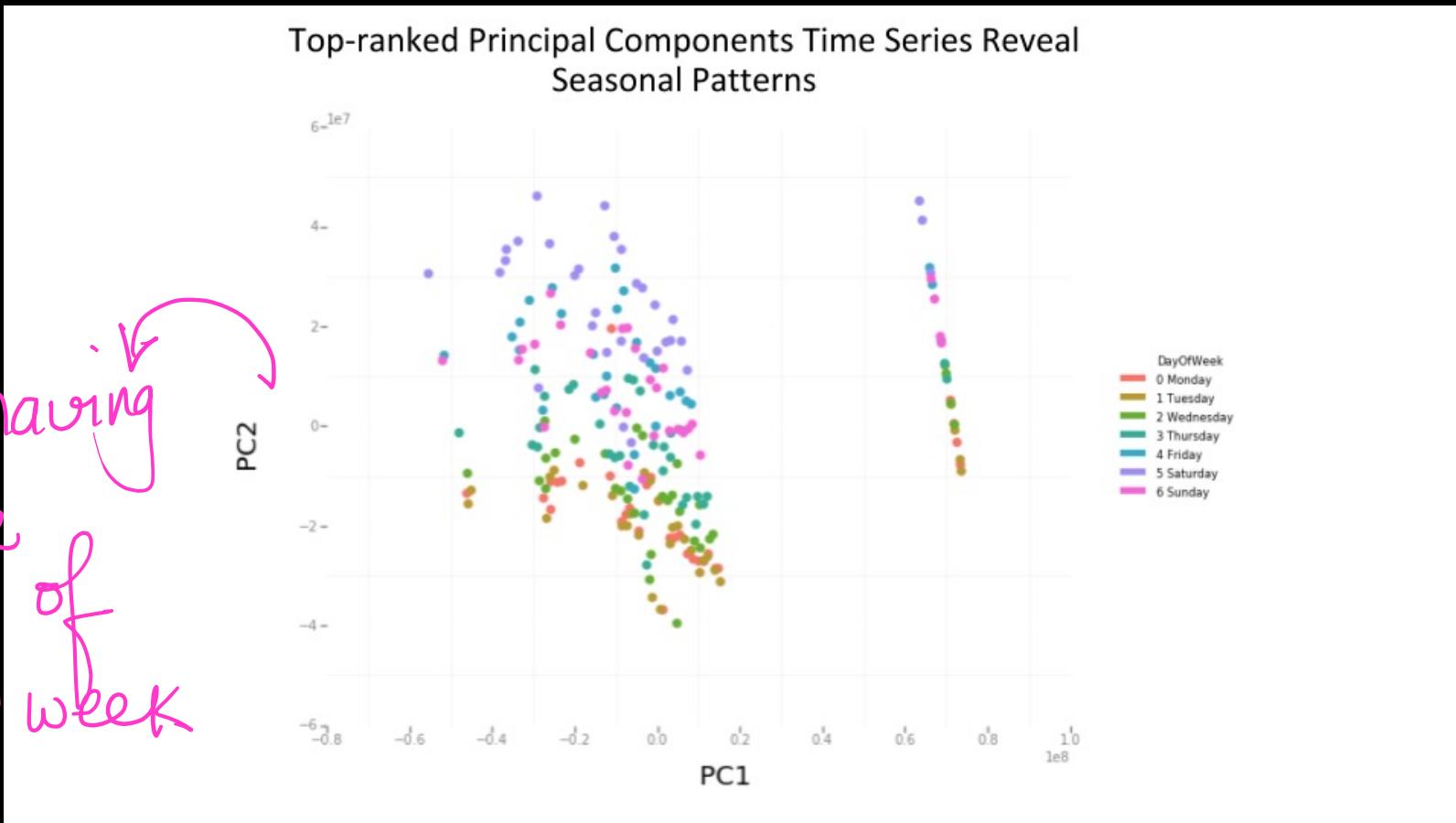


e.g.: trip duration & trip distance

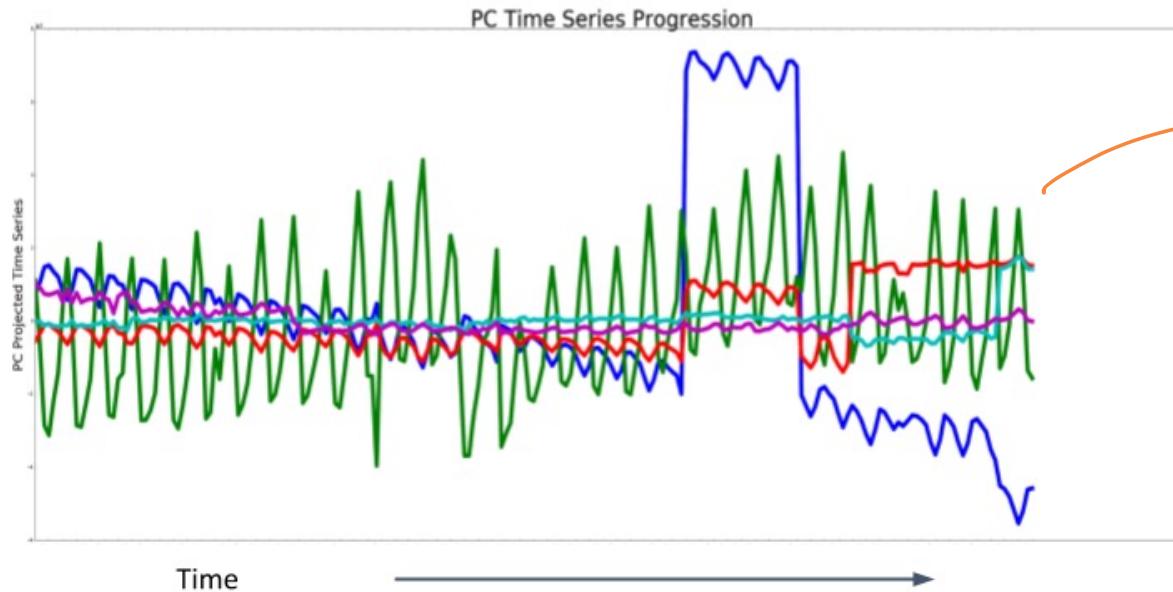
Uber: 90% of variability explained by the top -PC .

source: Uber

behaving
like
day
of
week



Principal Component Time Series Shows State Changes and Seasonalities



→ Top 5-PCS
across
time

(Q) Why not perform time-series anomaly detection on saw metrics?

Analyzing top-5 PCs

vs

Analyzing 100's of metrics

Time-series anomaly detection

- via forecasting

- Holt - winters model

<https://otexts.com/fpp2/holt-winters.html>

exponential moving average

Seasonality

Trends

(Q) Why exponential smoothing?

recent data is given more weightage



fast growing startup with sudden changes

Seasonality → weekends; evening-hours;
holidays - - -

Trend → Typically to capture broader
growth over time.

Lesson: Use the 'simplest' model that works.

easy to interpret

do NOT need DL everywhere

Table-level scores

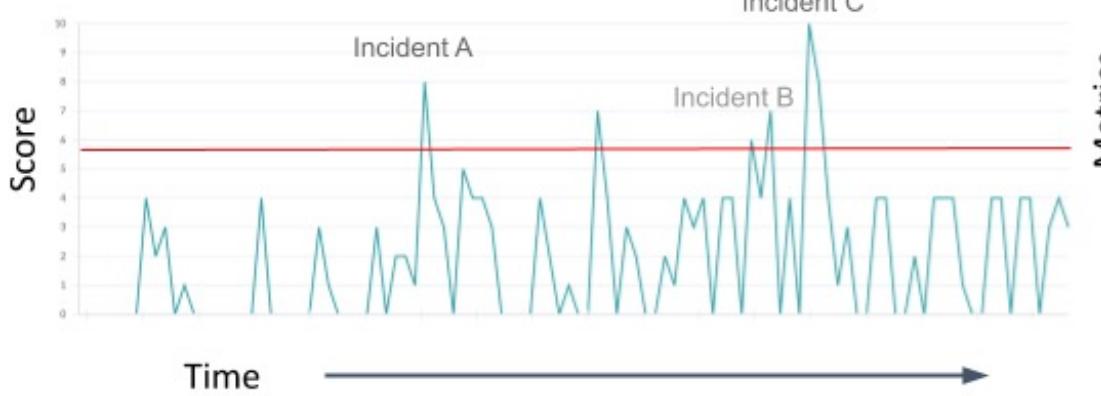
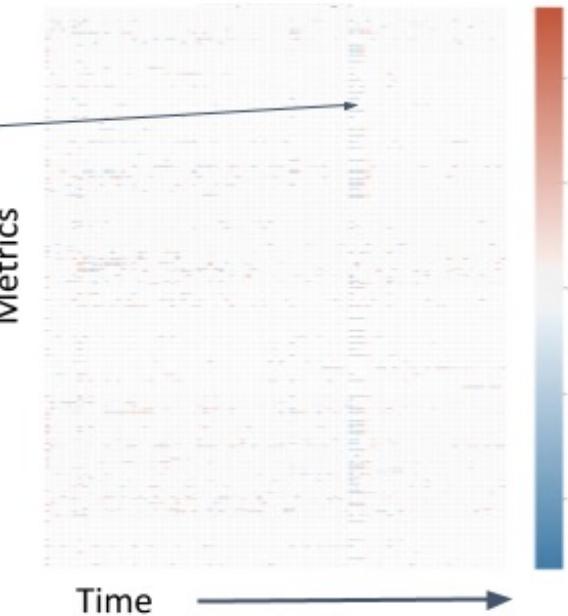
↳ Time-series of data in a Table

→ using top PC's with equal weightage

→ anomaly detection

Source: Uber

Anomaly Capture with Table and Metric-Level Scores

Table-level Score**Metric-level Score**

DQM: PySpark & Hive

↳ (what better platform than this?)

Vedica (columnar NoSQL datastore)

Clustering of Table Quality Scores Reconstructs Table Lineage Pattern

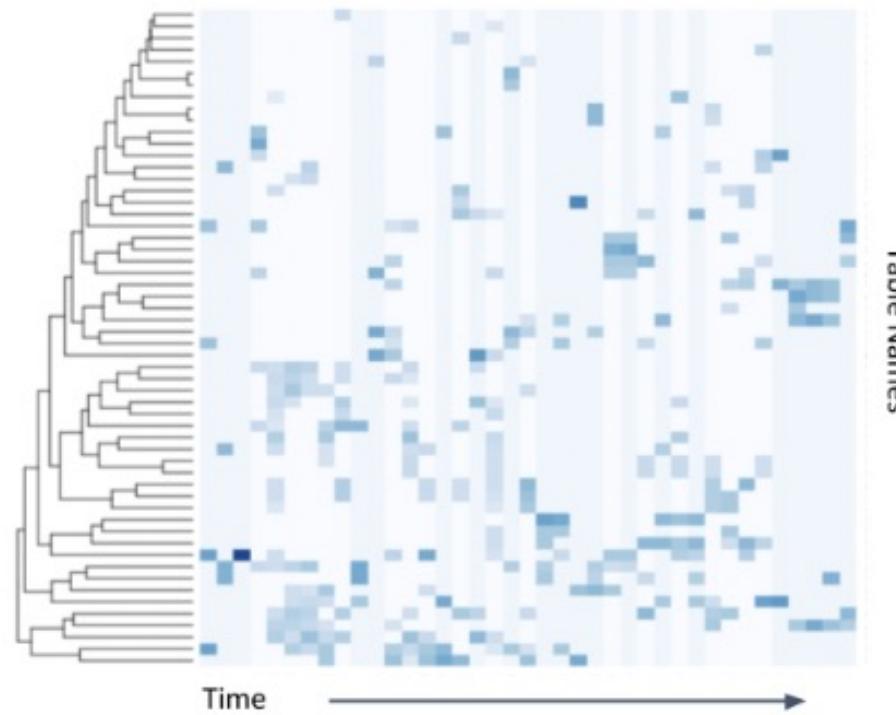


Figure 5. One of our next steps towards making alerts more intelligent is leveraging data table lineage information. In fact, we have observed strong correlation between data table quality and lineage as the clustering of table-level quality scores over time can reconstruct table ancestry. This is validated in practice as we see related tables have common root causes when they degrade in data quality.

Source : Uber

Source: Netflix.com

Source ↴

https://databricks.com/session_na20/an-approach-to-data-quality-for-netflix-personalization-systems



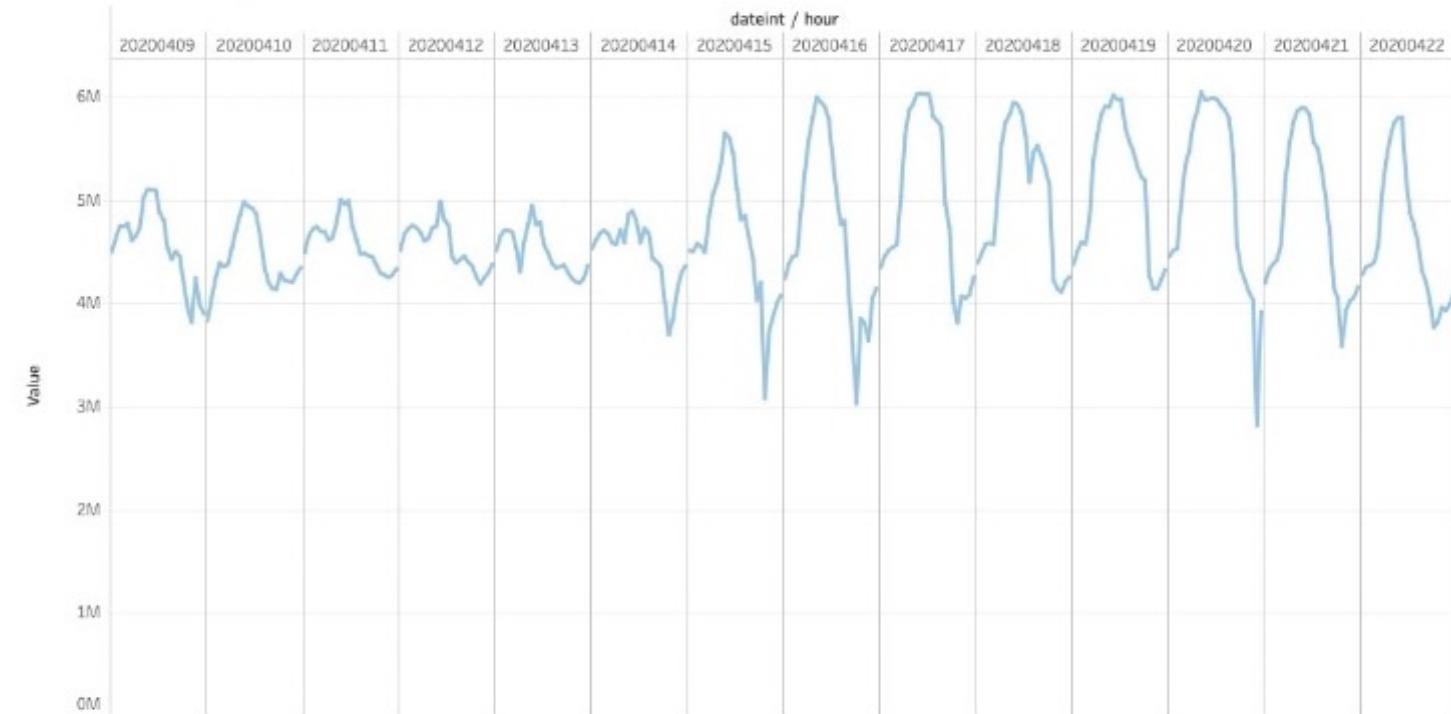
(slides & video)

Netflix PI Historical Fact Store

- Manages more than 10PB of data
- Manages hundreds of attributes
- More than 1 Billion rows flow through the ETL pipelines per day
- This data is used by several machine learned models and algorithms that enable Personalization

Bad Data Example 1: Drift

Raw Counts Hourly



SPARK+AI SUMMIT

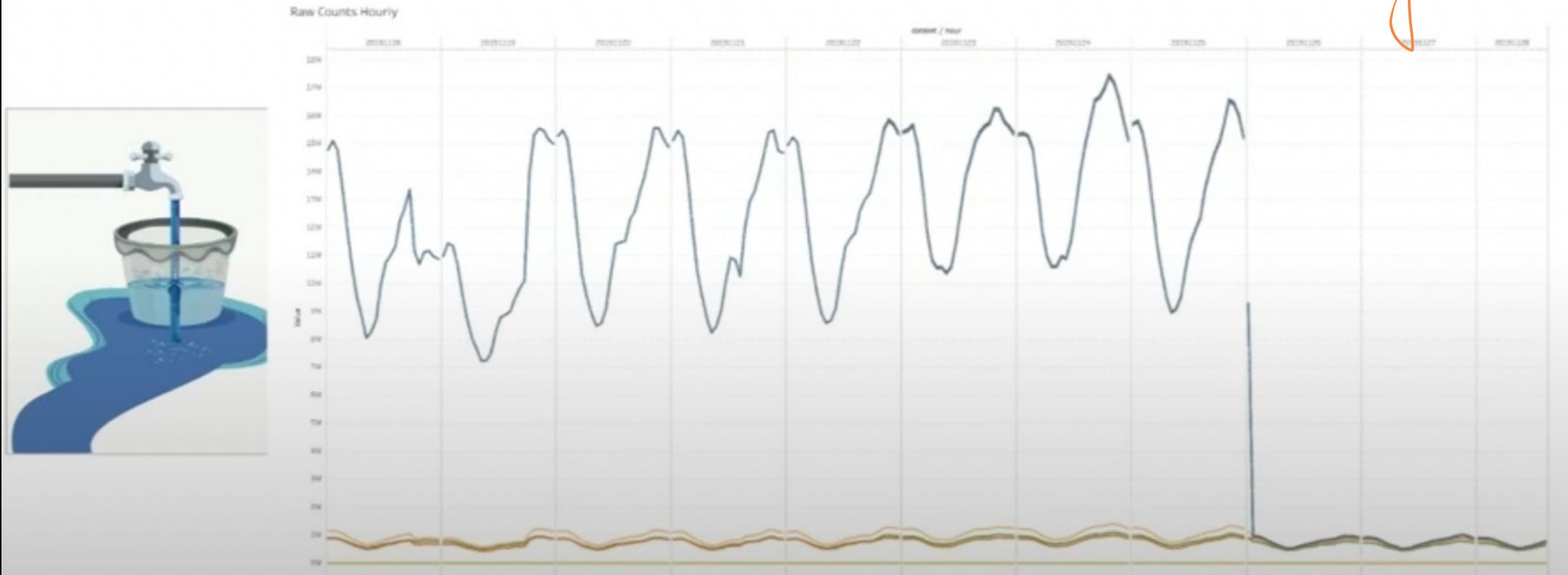
#Databricks #SparkAI

Bad Data Example 2: Drastic changes

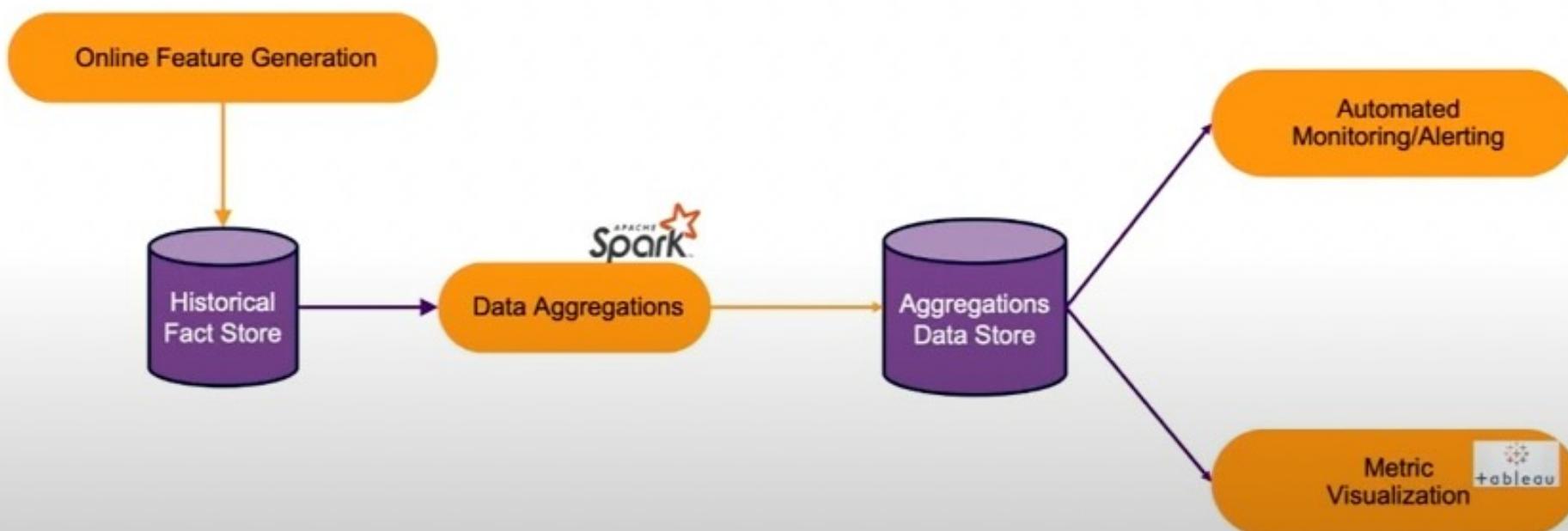


Bad Data Example 3: Under Utilization

Very important
but often ignored



Data Quality Architecture

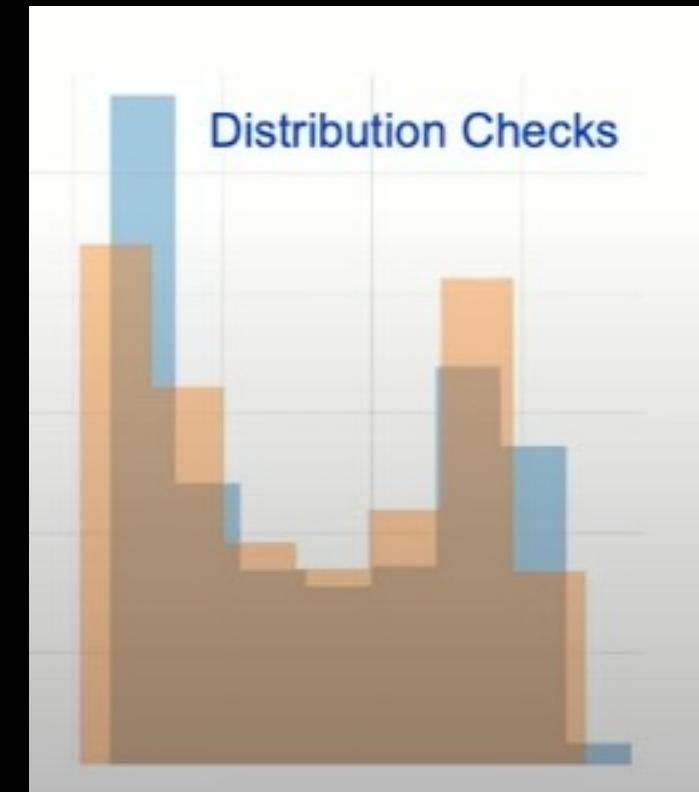


Aggregations

Date (DD/MM/YYYY)	Number Of Null Plays	Median Play Duration (seconds)	99 th Percentile Play Duration (seconds)
1/1/2020	3	360	450
1/2/2020	20	270	780
1/3/2020	1	150	570

Automated Monitoring:

↳ distribution - matching
of attributes



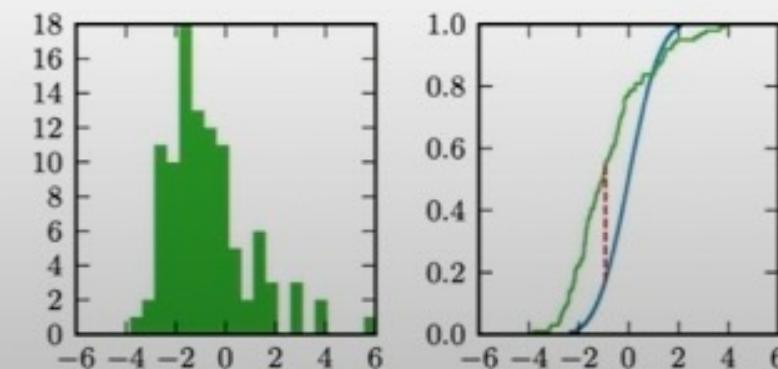
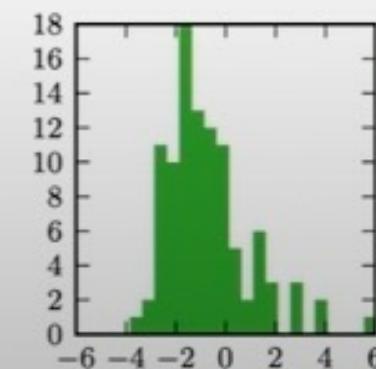
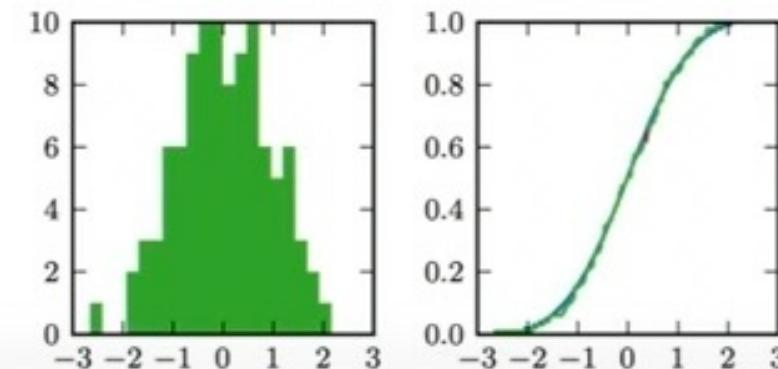
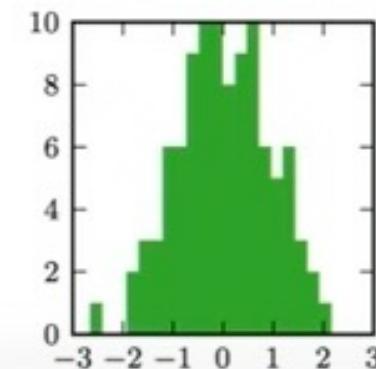
(Q) Any suggestions?

Distribution Checks - Statistical Test

Kolmogorov-Smirnov statistical test¹

$$F_n(a) = \frac{1}{n} \sum_i \mathbb{I}(x_i < a)$$

$$D_n = \max_x |F_n^1(x) - F_n^2(x)|$$



¹Source: [MIT 6.S085](#)

Audits

```
from pyspark.sql import SQLContext
import AuditRunner, Audit

df = SQLContext.sql('SELECT * FROM ml.aggregations WHERE dateint = 20200303')

# Audit Pseudo code
auditRunner = AuditRunner().add(
    Audit(df, 'num_null_plays', { 'upper_bound': 10 }, { 'alert': 'pager' }),
    Audit(df, 'num_null_plays', { 'normal_distribution': True }, { 'alert': 'dev_notification' }),
    Audit(df, 'avg_play_duration', { 'query_anomaly_detection': 'rad' }, { 'alert': 'dev_notification' })
)
# ...
```

Thresh δ 25

↳ distribution matching

↳ Robust Anomaly Detec \underline{t}

Netflix RAD:

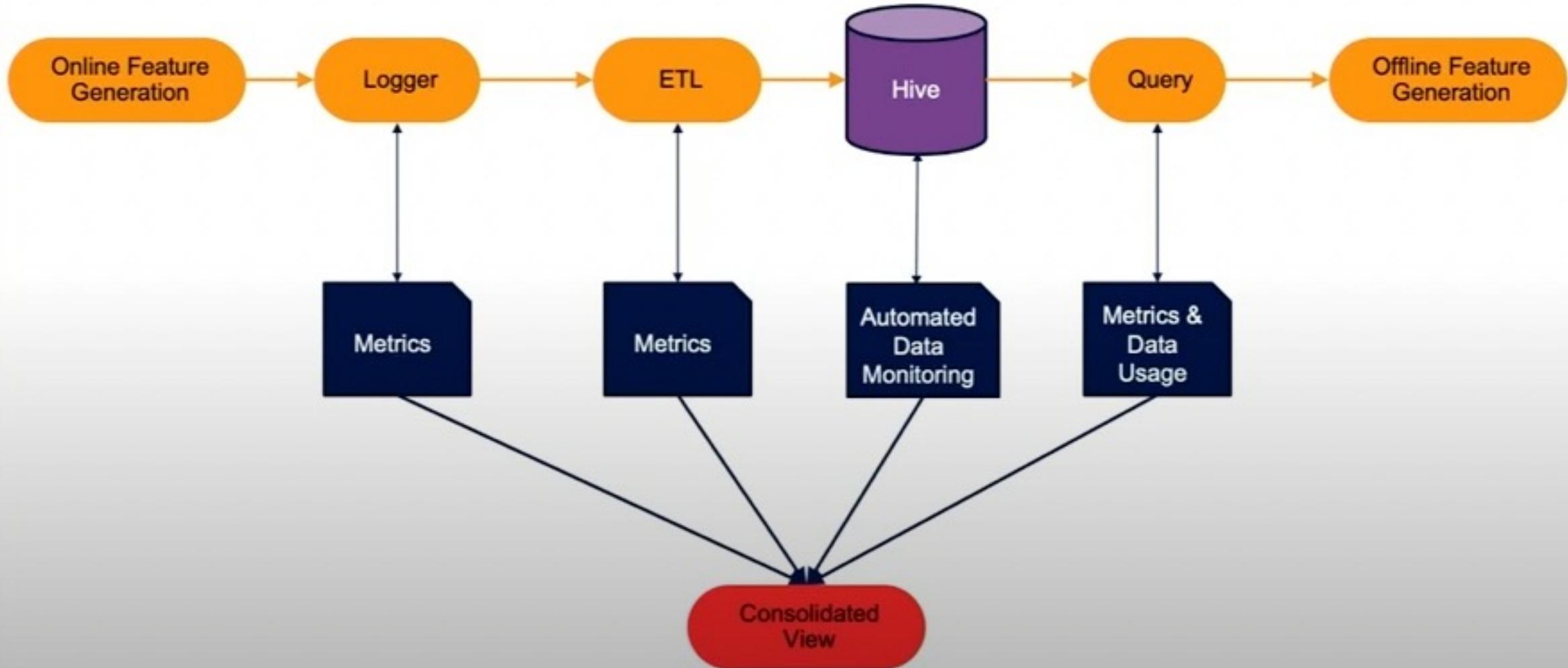
<https://netflixtechblog.com/rad-outlier-detection-on-big-data-d6b0494371cc>

↳ Robust PCA (repeatedly calculate SVD
and apply thresholds)

Plain old visualizations



Data Quality Checkpoints



Next Session : Amazon's Data Quality System