

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](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

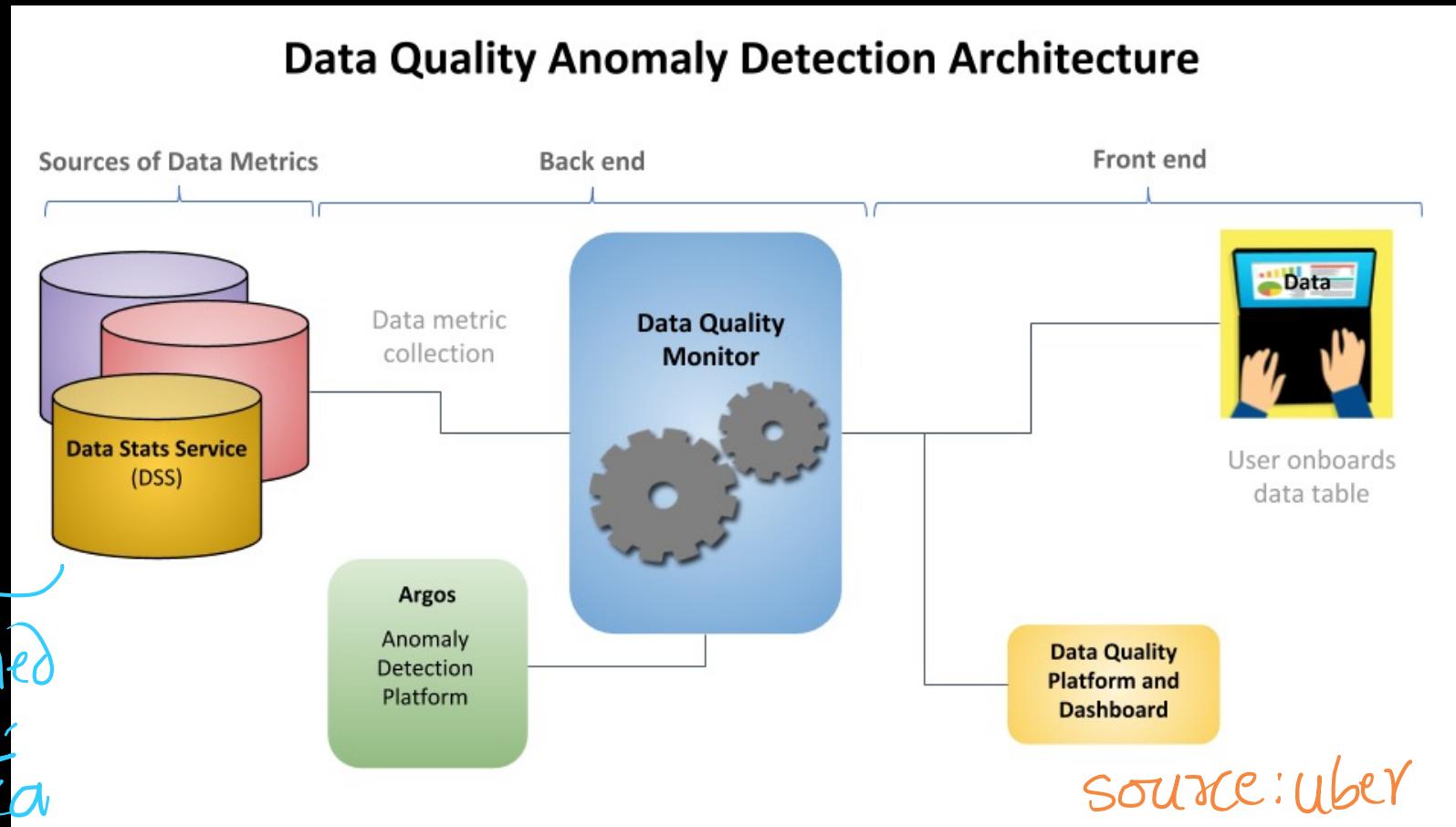
column-level

missing  
values

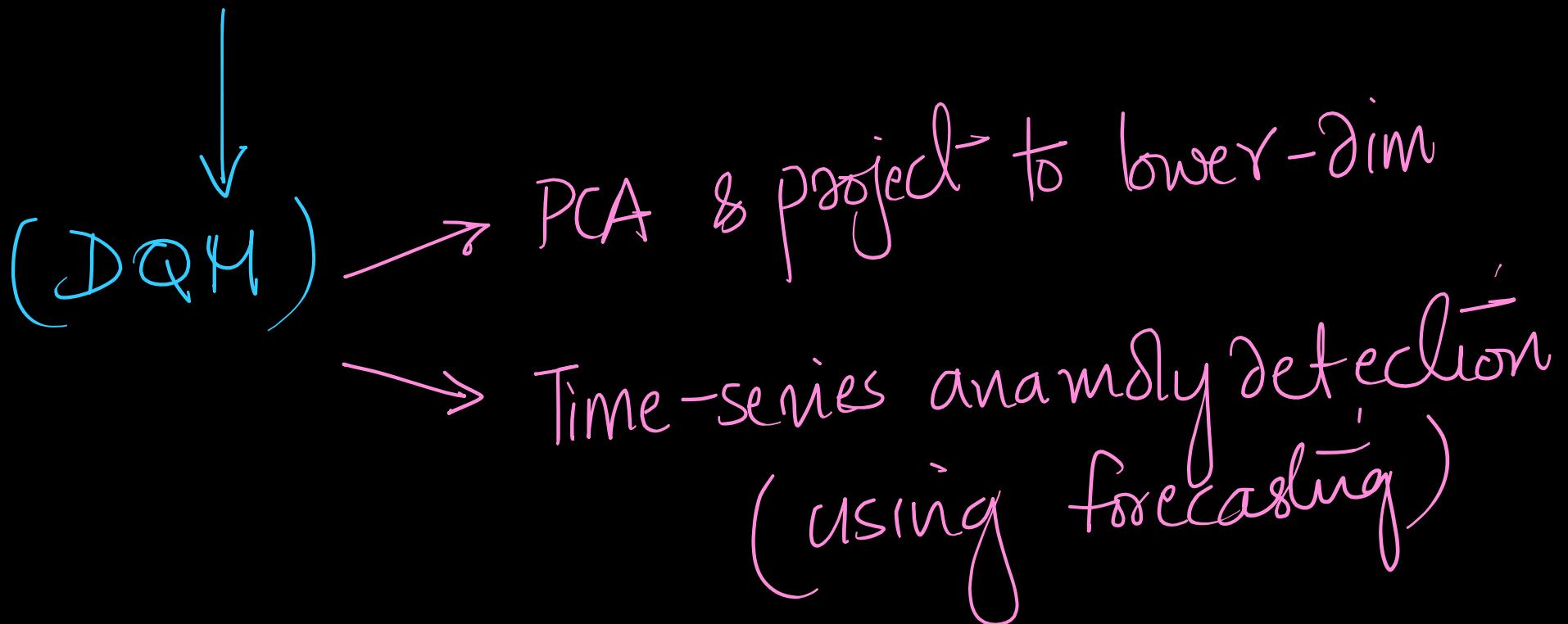
min; max; unique-values

→ 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

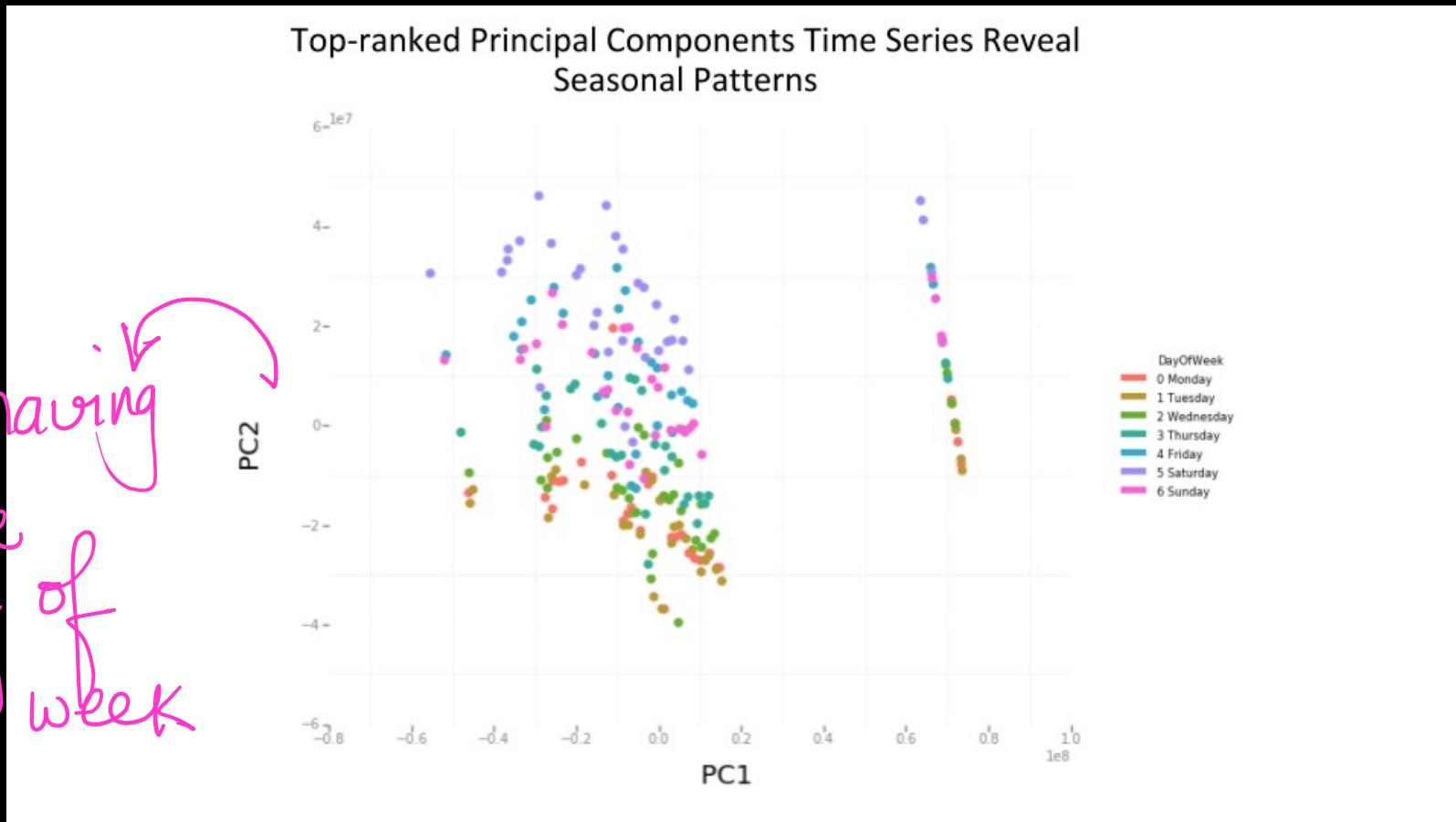


e.g.: trip duration & trip distance

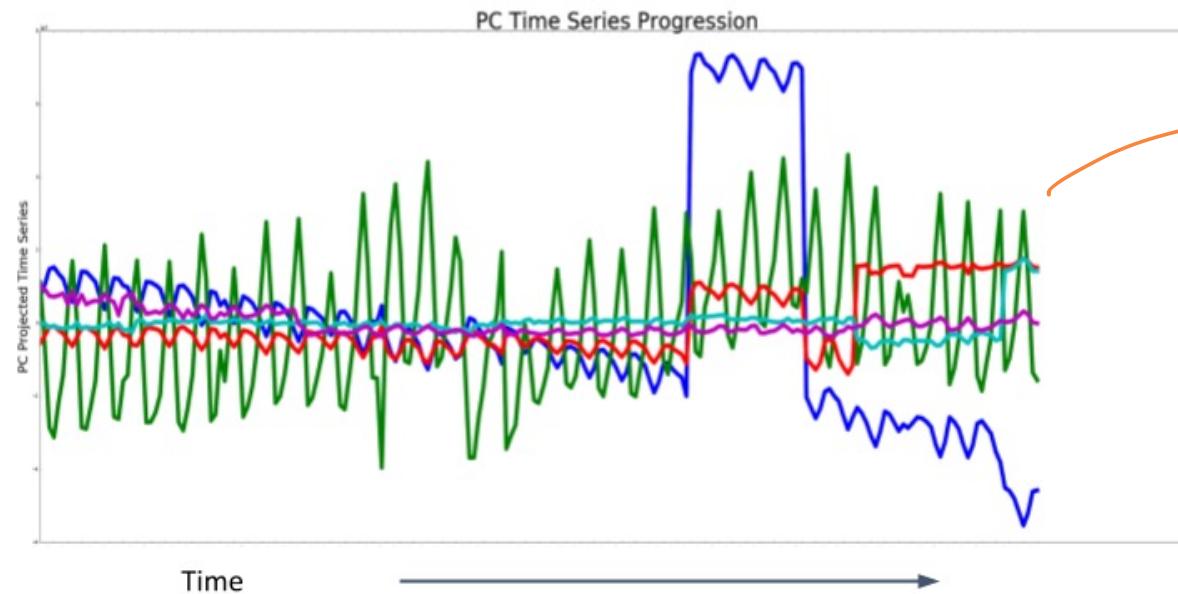
Uber: 90% of variability explained by the top -PC .

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

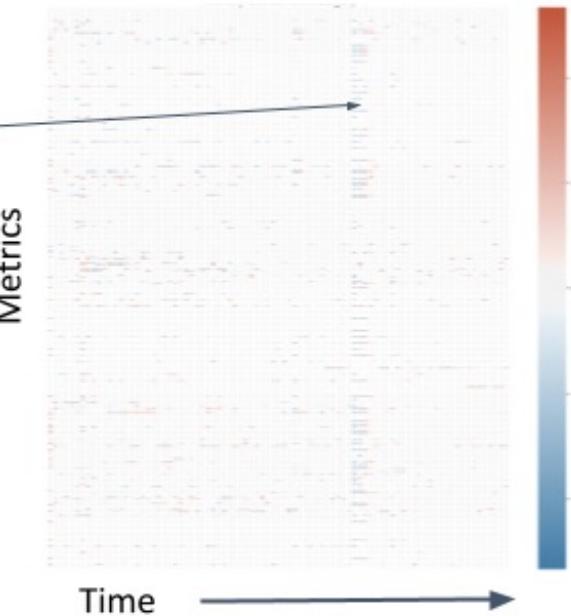
↳ Time-Series of data in a Table

→ using top PC's with equal weightage

→ anomaly detection

Source: Uber

## Anomaly Capture with Table and Metric-Level Scores

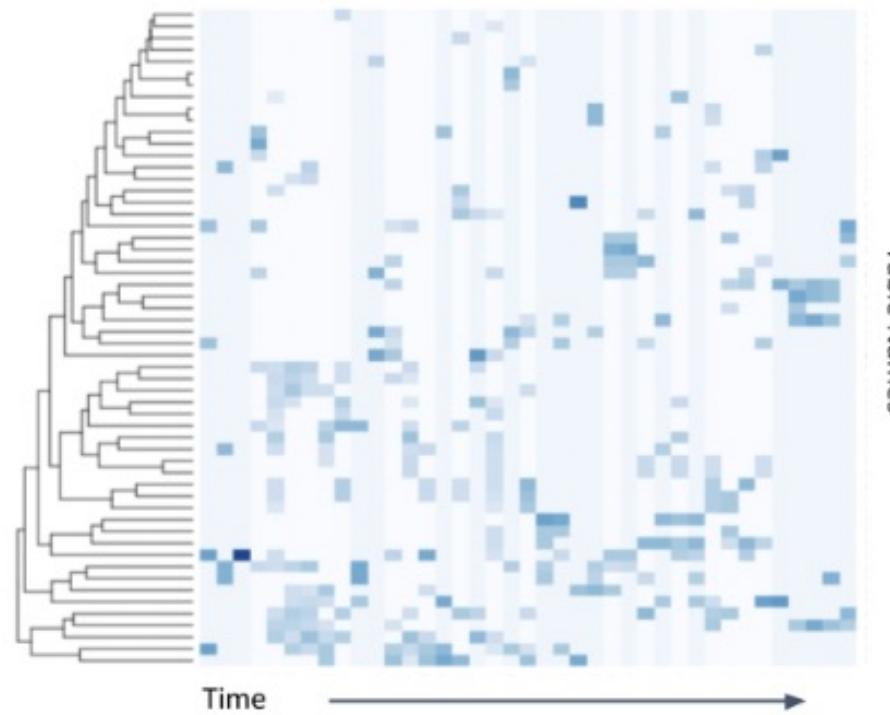
**Table-level Score****Metric-level Score**

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



(slides & video)

→ lots of attributes

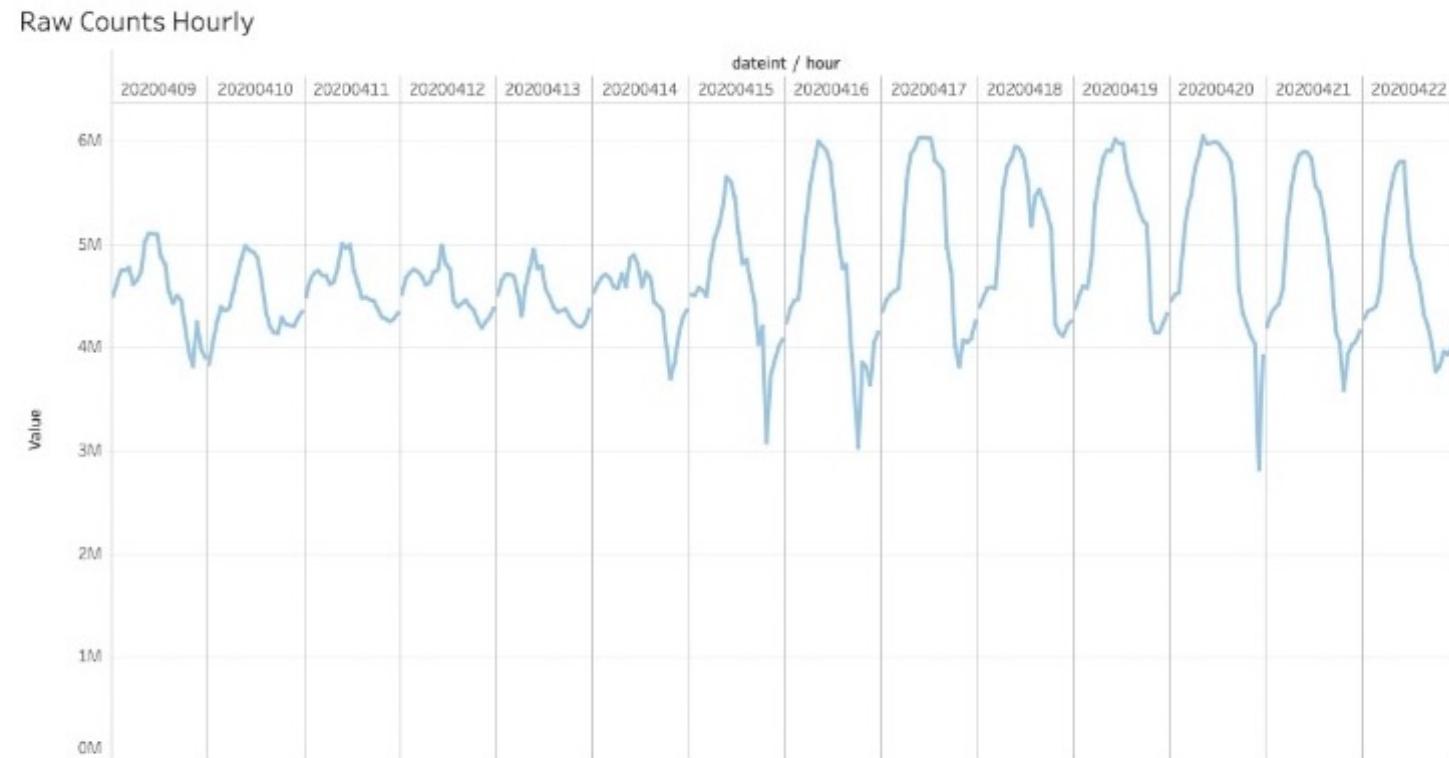
## Netflix PI Historical Fact Store

- 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

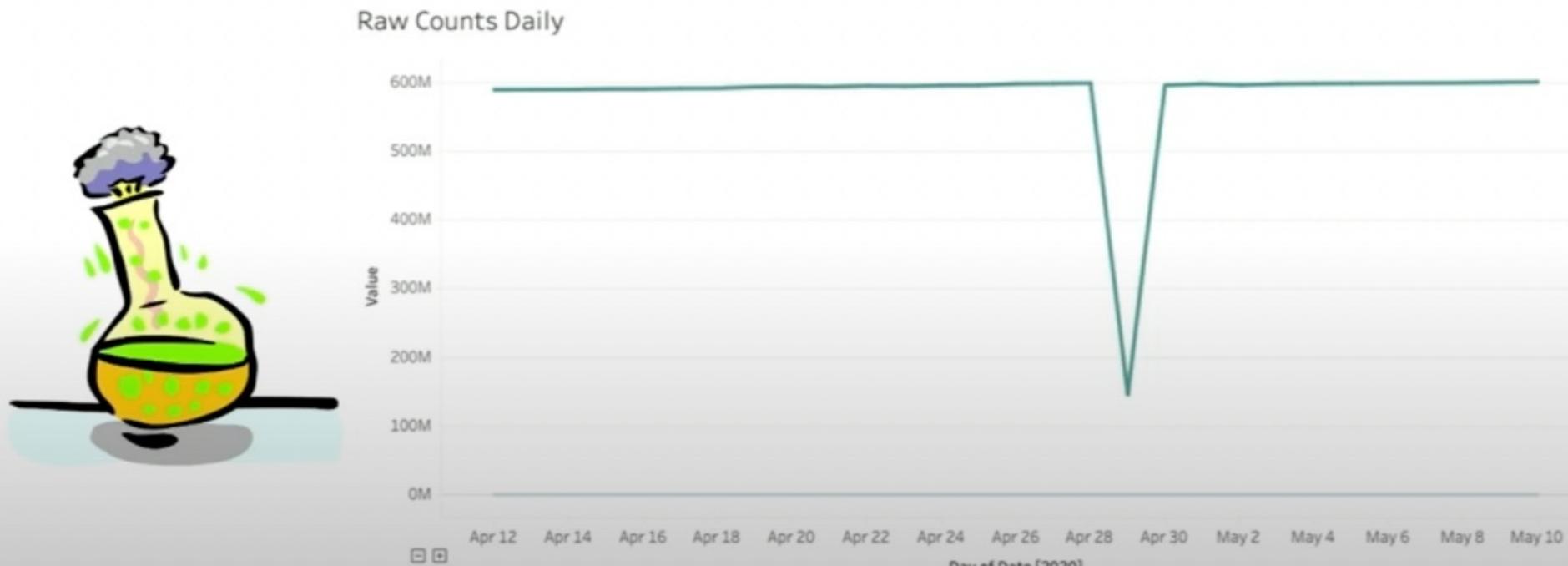
HIVE,  
Spark

e.g.: profile mix-up (lids vs adults) bug.

## Bad Data Example 1: Drift

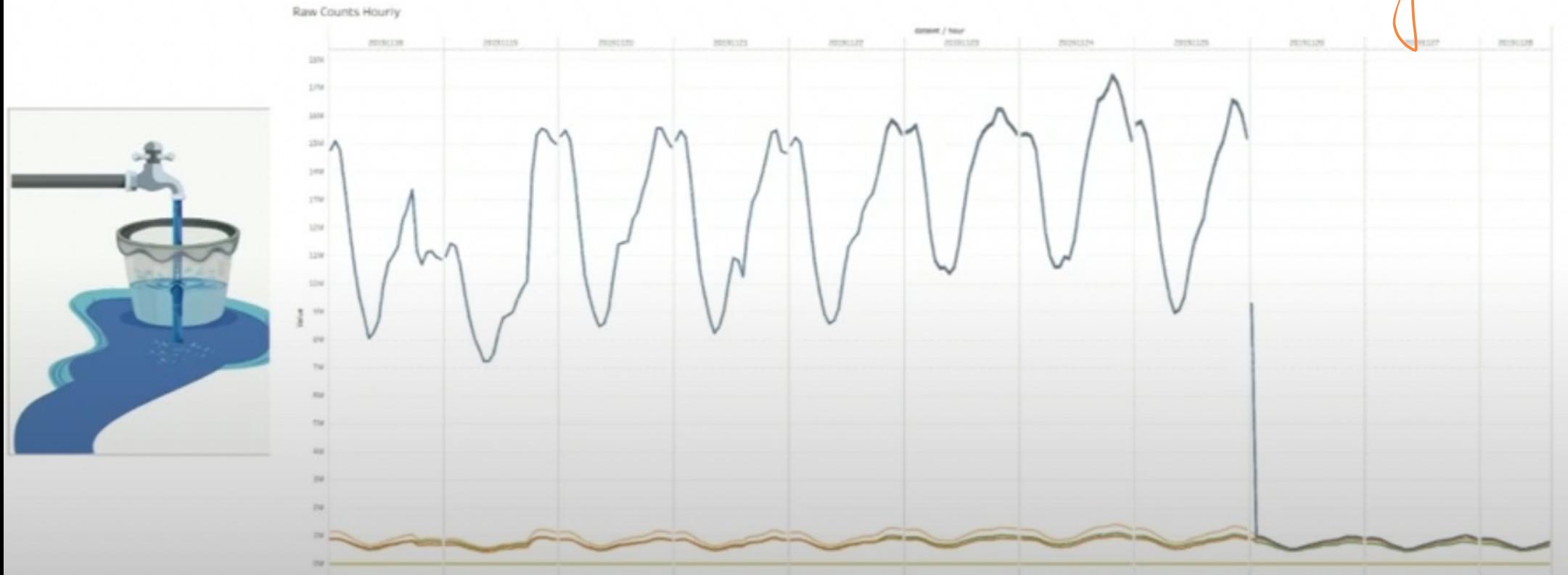


## Bad Data Example 2: Drastic changes



## Bad Data Example 3: Under Utilization

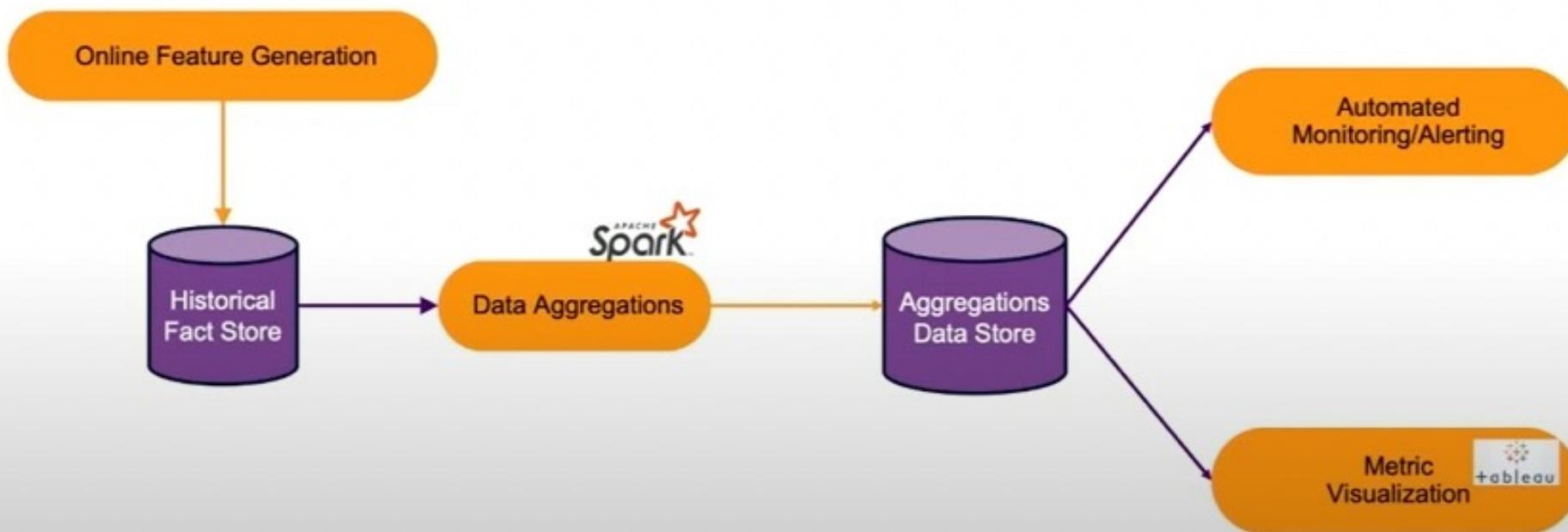
Very important  
but often ignored



## COST-SAVINGS

Source: Netflix.com  
#Dataviz #DataVizGuru

# Data Quality Architecture



## Example Data

ID	Member ID	Video ID	Timestamp (seconds)	Play Duration (seconds)	$\log_{10}(\text{Play Duration})$
1	10001	101	1577000001	360	2.556
2	10002	102	1577000002	780	2.892
3	10002	101	1577000003	NULL	2.556
4	10003	102	1577000004	300	2.477
5	10004	103	1577000005	1000000000	9.0
6	10005	105	1577000006	100	2.0
7	10006	104	1577000007	NULL	2.556

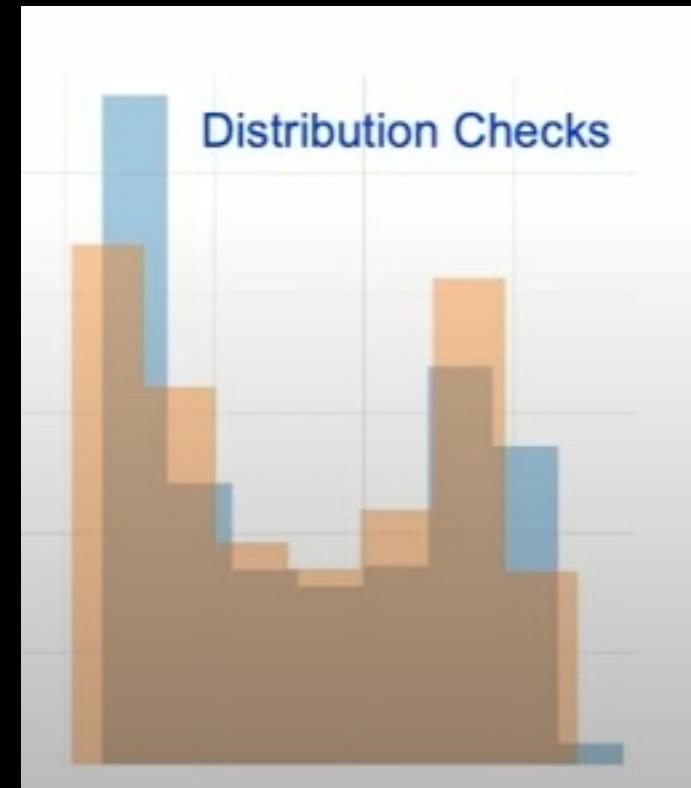
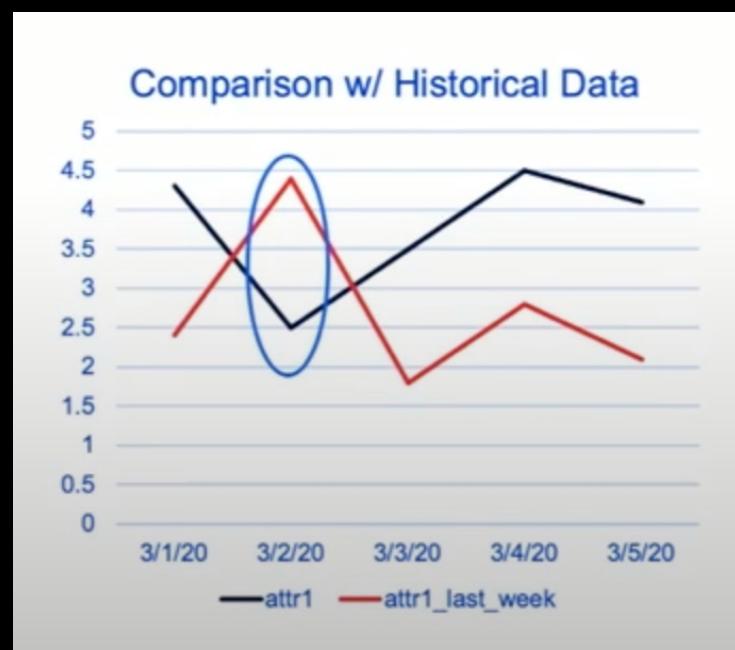
Impulsed

# Aggregations

Date (DD/MM/YYYY)	Number Of Null Plays	Median Play Duration (seconds)	99 <sup>th</sup> 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



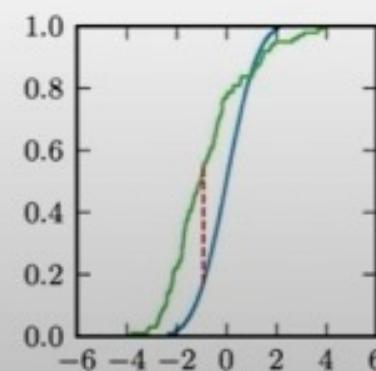
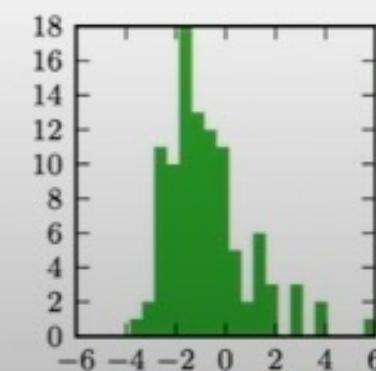
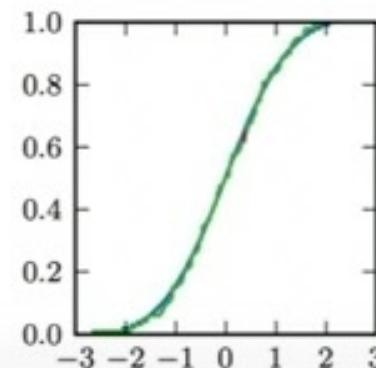
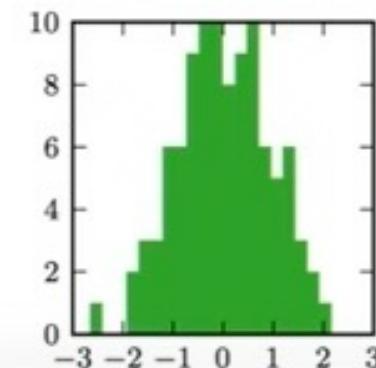
(Q) Any suggestions?

# Distribution Checks - Statistical Test

Kolmogorov-Smirnov statistical test<sup>1</sup>

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

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



<sup>1</sup>Source: [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 $\delta 25$

↳ distribution matching

↳ Robust Anomaly Detec $\underline{t}$

Netflix RAD:

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

↳ Robust PCA (repeatedly calculate SVD  
and apply thresholds)

<https://arxiv.org/pdf/0912.3599.pdf> (original Research Paper  
2009)

(Q) If you were to design Robust-PCA using regular PCA, how would you do it?

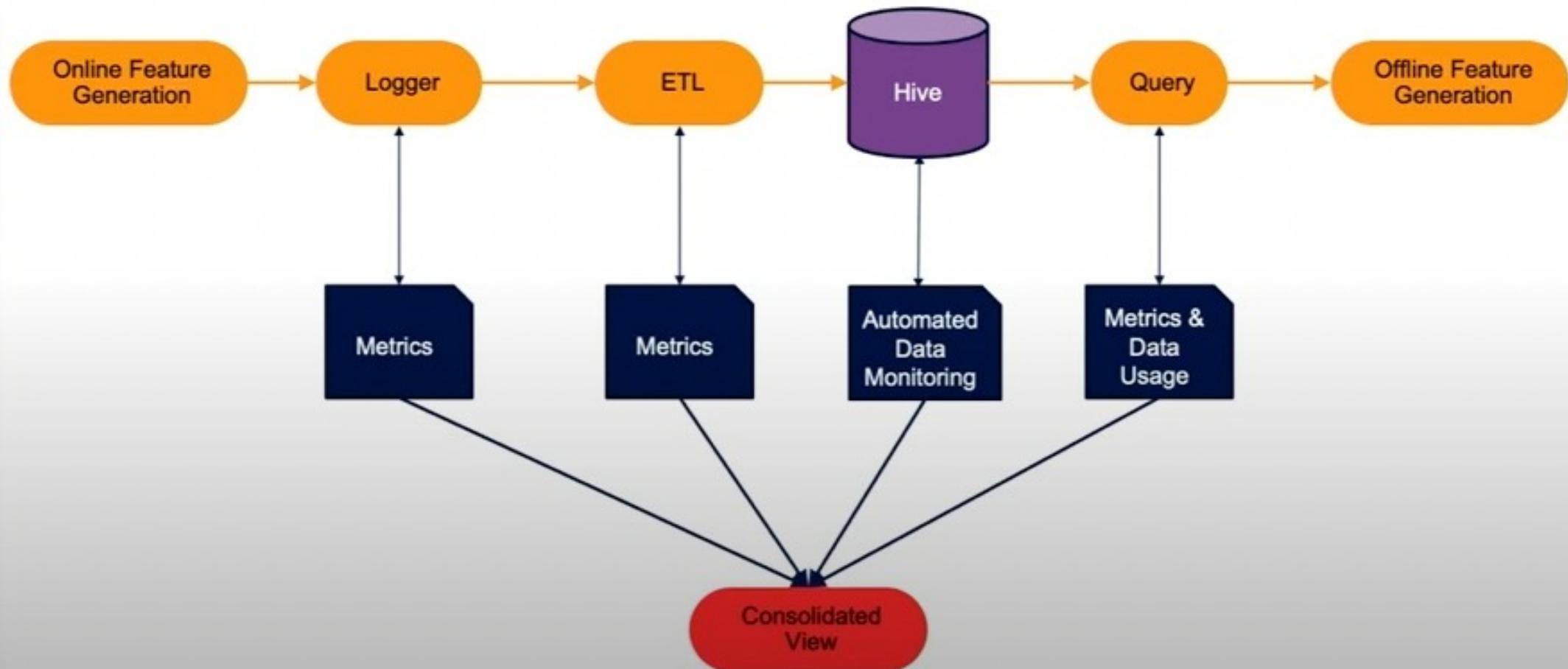
$$\min_u \mathbf{U}^T \mathbf{S} \mathbf{U}$$

$$\text{s.t. } \mathbf{U}^T \mathbf{U} = \mathbf{I}$$

# Plain old visualizations



# Data Quality Checkpoints

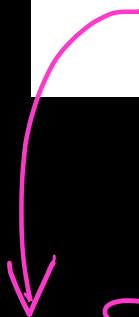


# Automating Large-Scale Data Quality Verification

Sebastian Schelter, Dustin Lange, Philipp Schmidt, Meltem Celikel, Felix Biessmann  
Amazon Research  
`{sseb,langed,phschmid,celikelm,biessmann}@amazon.com`

Andreas Grafberger<sup>\*</sup>  
University of Augsburg  
`andreas.grafberger@student.uni-augsburg.de`

<https://www.amazon.science/publications/automating-large-scale-data-quality-verification>



Spark-based

↑  
SOURCE

VLDB 2018

'Unit tests' for data

- user-defined constraints
- pre-defined tests (optimal implementation)

Examples below:

**Table 1: Constraints available for composing user-defined data quality checks.**

constraint	arguments	semantic
dimension <i>completeness</i>		
isComplete	column	check that there are no missing values in a column
hasCompleteness	column, udf	custom validation of the fraction of missing values in a column
dimension <i>consistency</i>		
isUnique	column	check that there are no duplicates in a column
hasUniqueness	column, udf	custom validation of the unique value ratio in a column
hasDistinctness	column, udf	custom validation of the unique row ratio in a column
isInRange	column, value range	validation of the fraction of values that are in a valid range
hasConsistentType	column	validation of the largest fraction of values that have the same type
isNonNegative	column	validation whether all values in a numeric column are non-negative
isLessThan	column pair	validation whether values in the 1st column are always less than in the 2nd column
satisfies	predicate	validation whether all rows match predicate
satisfiesIf	predicate pair	validation whether all rows matching 1st predicate also match 2nd predicate
hasPredictability	column, column(s), udf	user-defined validation of the predictability of a column
statistics (can be used to verify dimension <i>consistency</i> )		
hasSize	udf	custom validation of the number of records
hasTypeConsistency	column, udf	custom validation of the maximum fraction of values of the same data type
hasCountDistinct	column	custom validation of the number of distinct non-null values in a column
hasApproxCountDistinct	column, udf	custom validation of the approx. number of distinct non-null values
hasMin	column, udf	custom validation of a column's minimum value
hasMax	column, udf	custom validation of a column's maximum value
hasMean	column, udf	custom validation of a column's mean value
hasStandardDeviation	column, udf	custom validation of a column's standard deviation
hasApproxQuantile	column, quantile, udf	custom validation of a particular quantile of a column (approx.)
hasEntropy	column, udf	custom validation of a column's entropy
hasMutualInformation	column pair, udf	custom validation of a column pair's mutual information
hasHistogramValues	column, udf	custom validation of column histogram
hasCorrelation	column pair, udf	custom validation of a column pair's correlation
time		
hasNoAnomalies	metric, detector	validation of anomalies in time series of metric values

e.g.

```
1 val numTitles = callRestService(...)  
2 val maxExpectedPhoneRatio = computeRatio(...)  
3  
4 var checks = Array()  
5  
6 checks += Check(Level.Error)  
7   .isComplete("customerId", "title",  
8     "impressionStart", "impressionEnd",  
9     "deviceType", "priority")  
10  .isUnique("customerId", "countryResidence",  
11    "deviceType", "title")  
12  .hasCountDistinct("title", _ <= numTitles)  
13  .hasHistogramValues("deviceType",  
14    _.ratio("phone") <= maxExpectedPhoneRatio)  
15  
16 checks += Check(Level.Error)  
17   .isNonNegative("count")  
18   .isLessThan("impressionStart", "impressionEnd")  
19   .isInRange("priority", ("hi", "lo"))  
20  
21 checks += Check(Level.Warning, on="delta")  
22   .hasNoAnomalies(Size, OnlineNormal(stdDevs=3))  
23 checks += Check(Level.Error, on="delta")  
24   .hasNoAnomalies(Size, OnlineNormal(stdDevs=4))  
25  
26 checks += Check(Level.Warning)  
27   .hasPredictability("countryResidence",  
28     ("zipCode", "cityResidence"), precision=0.99)  
29  
30 Verification.run(data, checks)
```

Listing 1: Example for declarative data quality constraint definitions using our API.

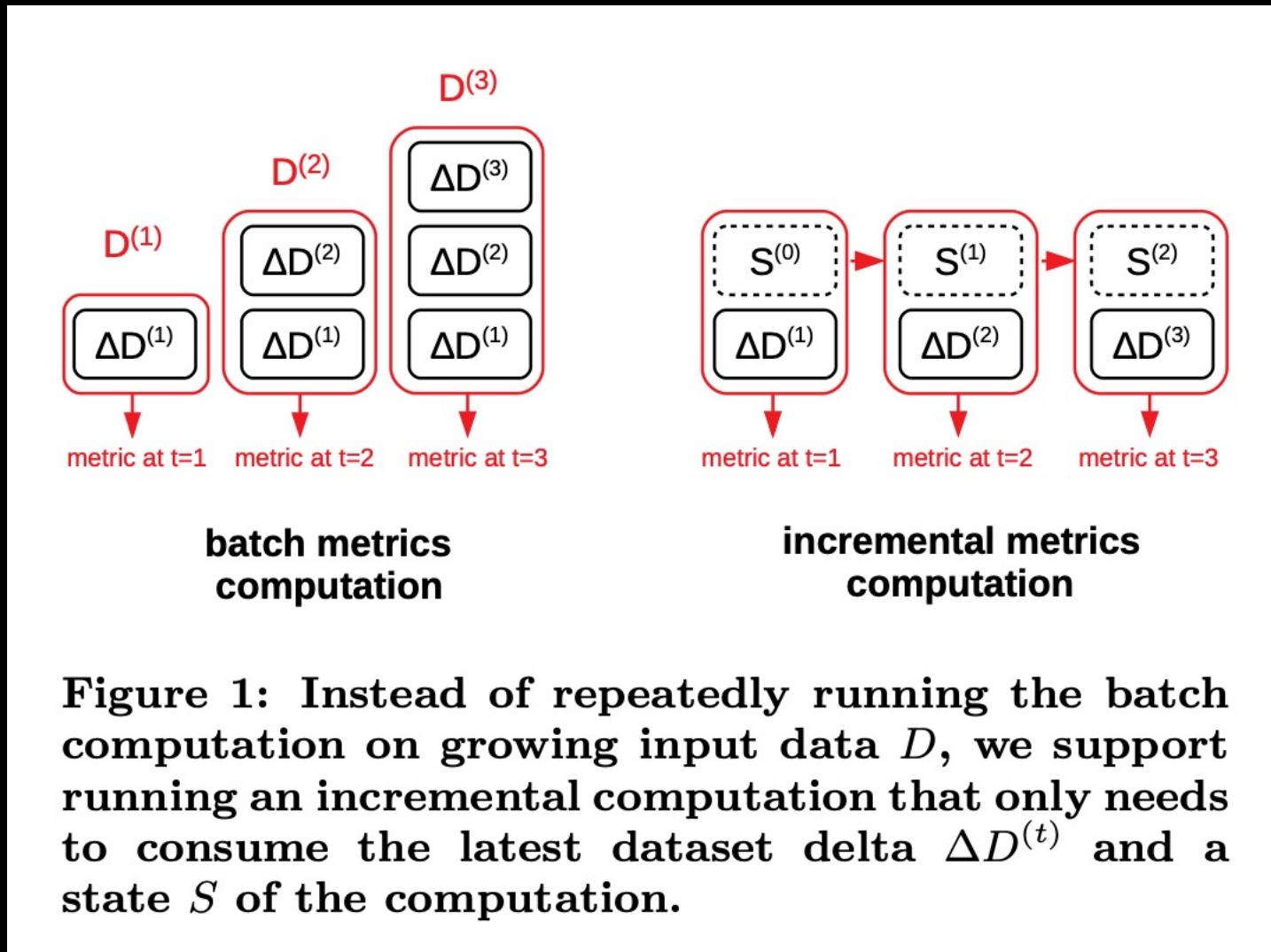
API for  
'unit tests'

Why use the API-based software design?

## Implementation of the API-functions:

- Incremental computation e.g.: mean
- Approximate values: e.g.: Quantiles, Count Distinct
- Predictability: build a simple model & use standard metrics

## Incremental Computation

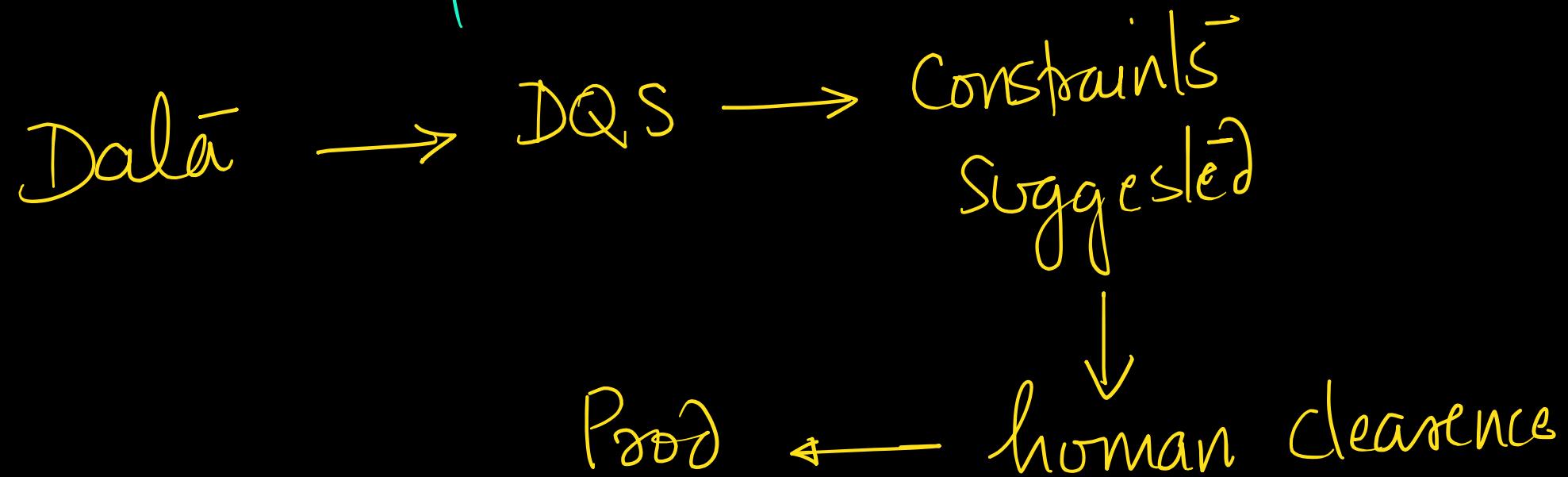


# Constraint Suggestions

- heuristics-based
- single-column profiling
  - ① datatype ; completeness ; #distinct-values ; size,
  - ② Summary stats for Numerical fields;
  - ③ distribution-fitting

## Constraint Suggestion

- human in the loop
- increases adoption



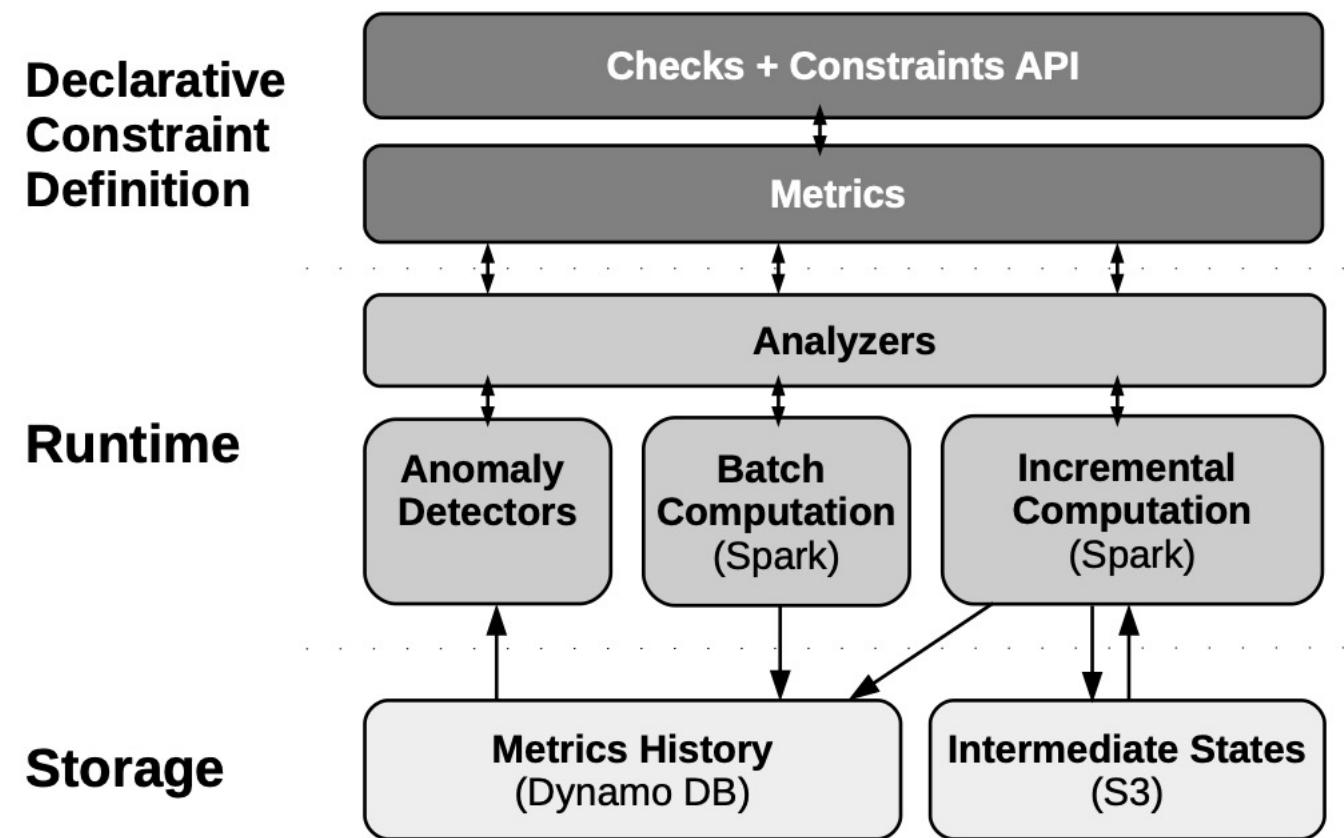
(Q) Can you suggest "Predictability"  
constraints also?

(not discussed in the research paper)

## Anomaly Detection :

- historical time-series
- BYOT : Bring your own technique

# Architecture :



**Figure 2:** System architecture: users declaratively define checks and constraints to combine with their verification code. The system identifies the required metrics and efficiently computes them in the runtime layer. The history of metrics and intermediate states of incremental computations are maintained in AWS storage services.

Section 5 in the research paper

→ experiments on public datasets

## Key-take aways

- Spark & HIVE
- Time-Series anomalies
- Simple & Robust algorithms/techniques
- You can do this

Idea for an open-source project

DQM based on Spark

(very similar to Amazon's  
System)