# Data Intensive Applications

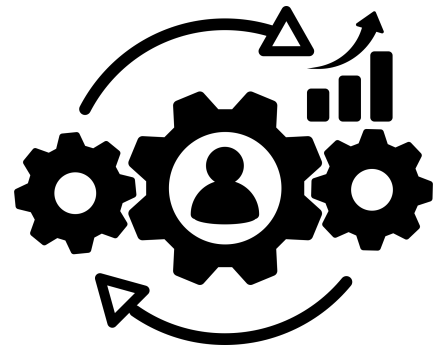## An Introduction

# Tonight's Quest

**What** is a Data Intensive Application - Our definition?

What are some of the **challenges** in building DIA?

**Who** Builds DIA?

How do we **classify** and **assess** the value of a DIA?

**How** are DIA Built?

Data intensive applications combine a large complement of data and compute to augment or replace human decision making for the benefit of optimizing important processes.
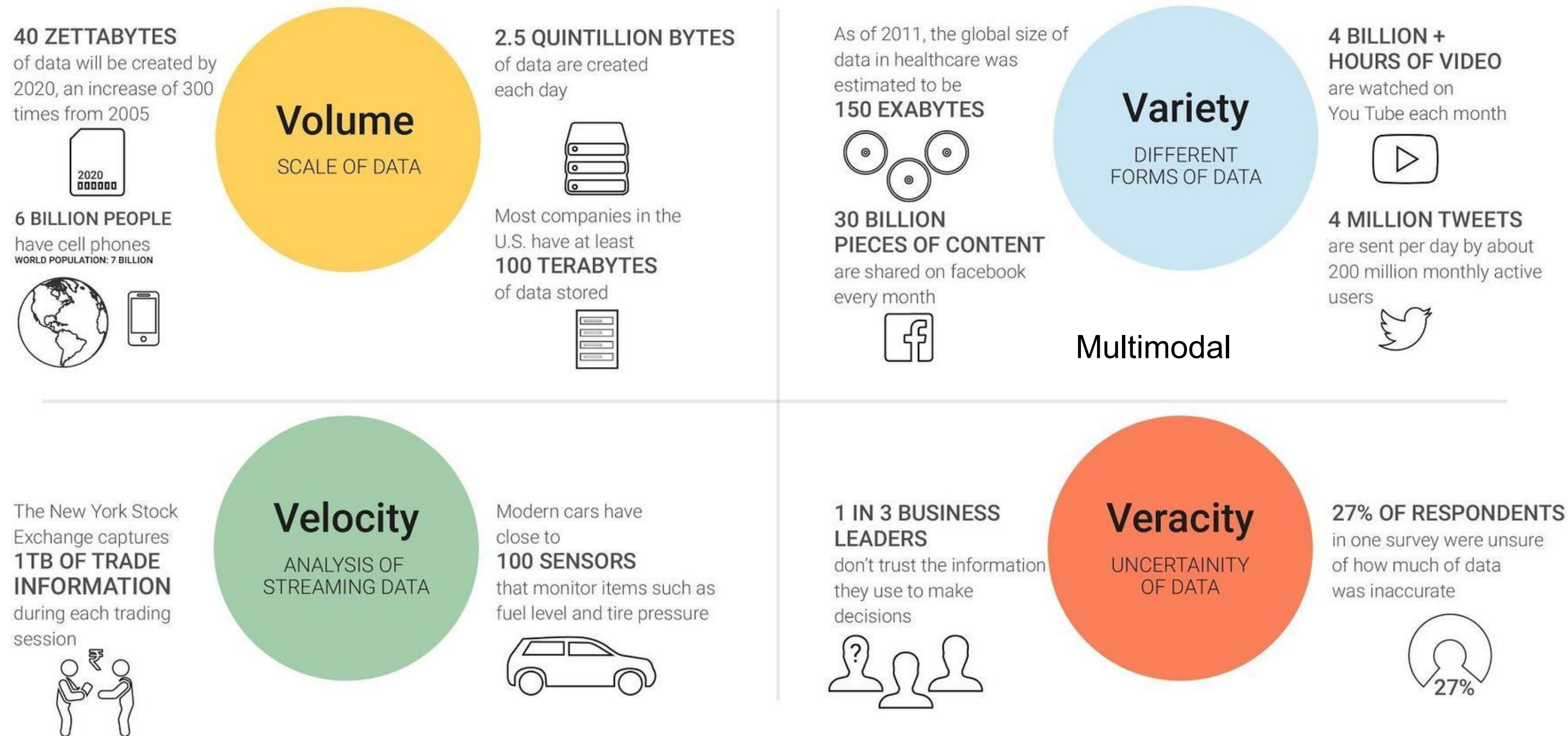
Data intensive applications combine a **_large_** complement of **_data_** and **_compute_** to **_augment_** or **_replace_** human **decision** making for the benefit of **_optimizing_** important **_processes_**.

UNIVERSITY *of* ROCHESTER

# THE 4 V'S OF BIG DATA

**40 ZETTABYTES**
of data will be created by 2020, an increase of 300 times from 2005

**6 BILLION PEOPLE**
have cell phones
WORLD POPULATION: 7 BILLION

## Volume
SCALE OF DATA

**2.5 QUINTILLION BYTES**
of data are created each day

Most companies in the U.S. have at least
**100 TERABYTES**
of data stored

As of 2011, the global size of data in healthcare was estimated to be
**150 EXABYTES**

**30 BILLION PIECES OF CONTENT**
are shared on facebook every month

## Variety
DIFFERENT FORMS OF DATA

**4 BILLION + HOURS OF VIDEO**
are watched on You Tube each month

**4 MILLION TWEETS**
are sent per day by about 200 million monthly active users

Multimodal

The New York Stock Exchange captures
**1TB OF TRADE INFORMATION**
during each trading session

## Velocity
ANALYSIS OF STREAMING DATA

Modern cars have close to
**100 SENSORS**
that monitor items such as fuel level and tire pressure

**1 IN 3 BUSINESS LEADERS**
don't trust the information they use to make decisions

## Veracity
UNCERTAINITY OF DATA

**27% OF RESPONDENTS**
in one survey were unsure of how much of data was inaccurate

27%

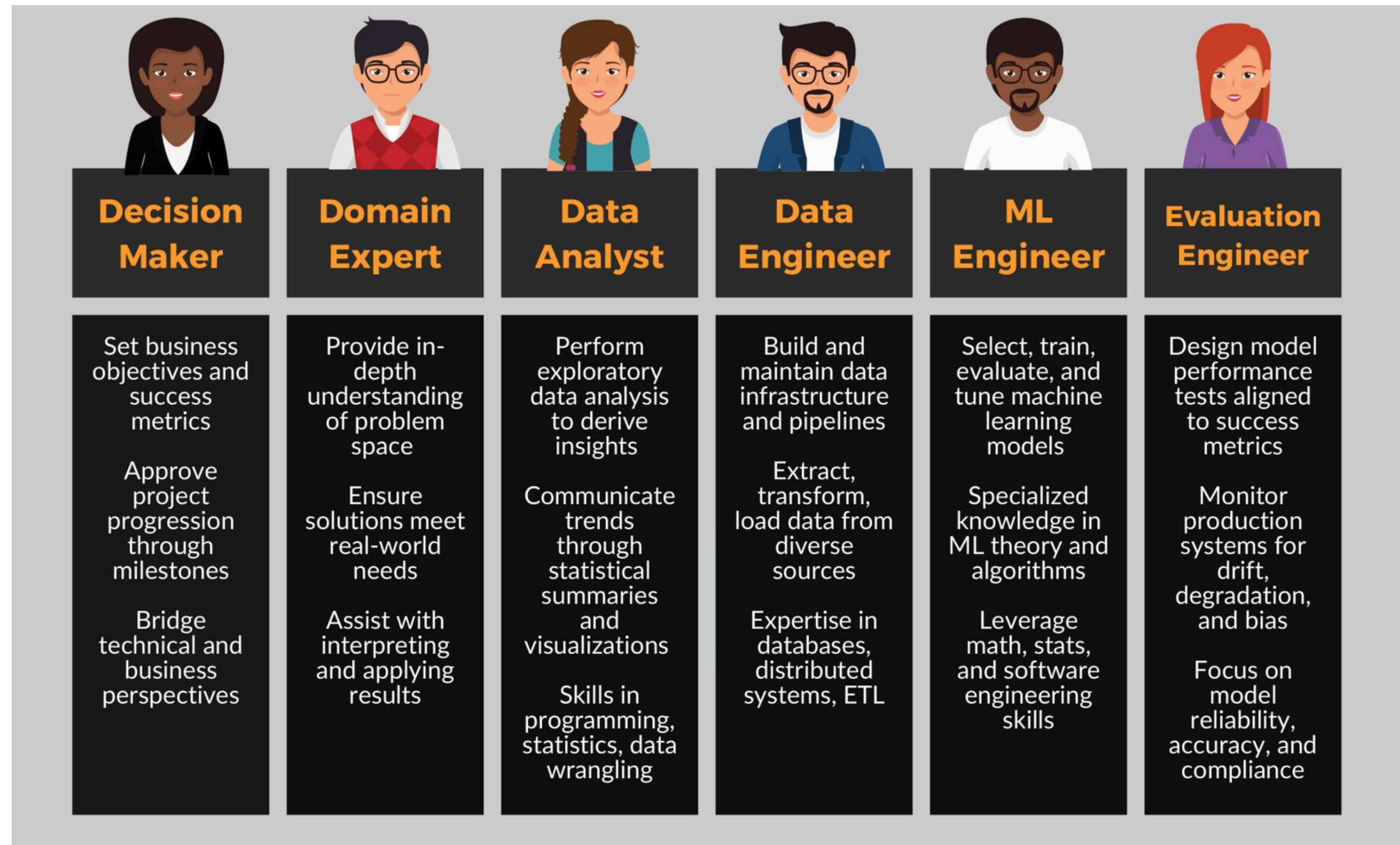Reference : http://www.ibmbigdatahub.com/infographic/four-vs-big-data

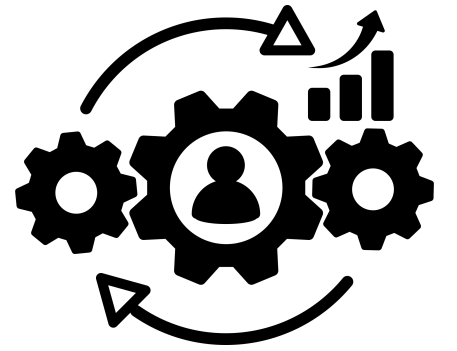**What do we mean by "large" complement of data?**
**Our work in this class is focused on "handling" these data challenges.**
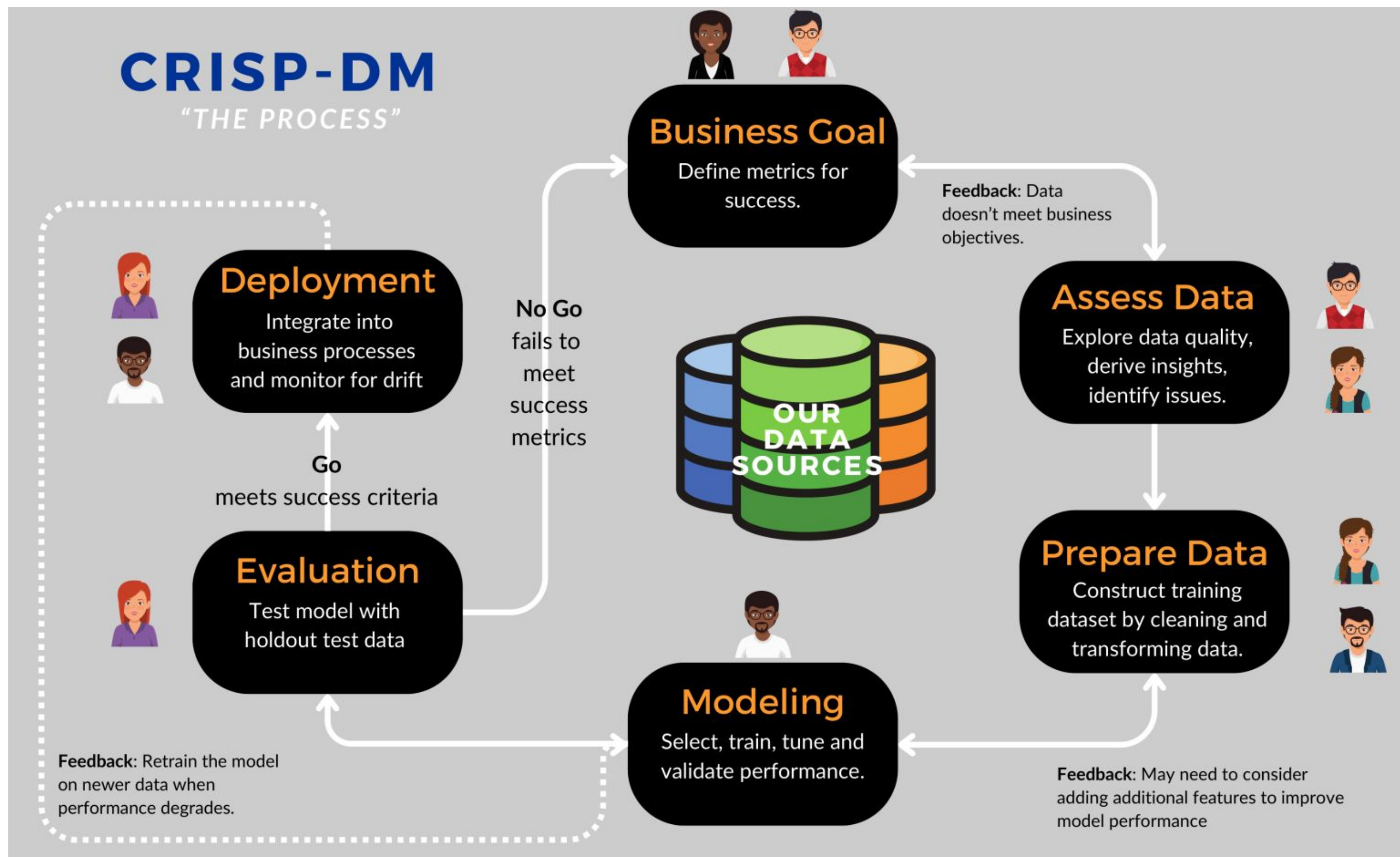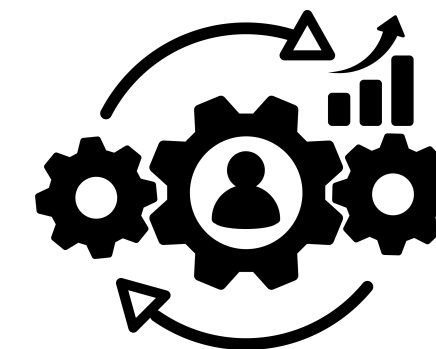
**What do we mean by "large" complement of compute?**

# Building Data Intensive Applications is a team sport.



| Decision Maker | Domain Expert | Data Analyst | Data Engineer | ML Engineer | Evaluation Engineer |
|---|---|---|---|---|---|
| Set business objectives and success metrics | Provide in-depth understanding of problem space | Perform exploratory data analysis to derive insights | Build and maintain data infrastructure and pipelines | Select, train, evaluate, and tune machine learning models | Design model performance tests aligned to success metrics |
| Approve project progression through milestones | Ensure solutions meet real-world needs | Communicate trends through statistical summaries and visualizations | Extract, transform, load data from diverse sources | Specialized knowledge in ML theory and algorithms | Monitor production systems for drift, degradation, and bias |
| Bridge technical and business perspectives | Assist with interpreting and applying results | Skills in programming, statistics, data wrangling | Expertise in databases, distributed systems, ETL | Leverage math, stats, and software engineering skills | Focus on model reliability, accuracy, and compliance |

7

# What we need in a process?

| Structured Approach | Colaboration | Iteration | Predictability | Flexibiliy | Focus |
|---|---|---|---|---|---|
| Defines a clear, methodical sequence of steps to guide the end-to-end data value creation. | Clearly defines roles & responsibility to coordinate work across teams with diverse skills. | Allows folding back to previous steps when needed to refine the solution. | Following consistent standard methodology improves predictability. | The methodology can be tailored to fit different project needs and constraints. | Business goals anchor the process to ensure work stays aligned with delivering value. |

CRISP-DM
"THE PROCESS"

**Business Goal**
Define metrics for success.

Feedback: Data doesn't meet business objectives.

**Deployment**
Integrate into business processes and monitor for drift

**No Go**
fails to meet success metrics

**Assess Data**
Explore data quality, derive insights, identify issues.

OUR DATA SOURCES

**Go**
meets success criteria

**Evaluation**
Test model with holdout test data

**Prepare Data**
Construct training dataset by cleaning and transforming data.

**Modeling**
Select, train, tune and validate performance.

Feedback: Retrain the model on newer data when performance degrades.

Feedback: May need to consider adding additional features to improve model performance

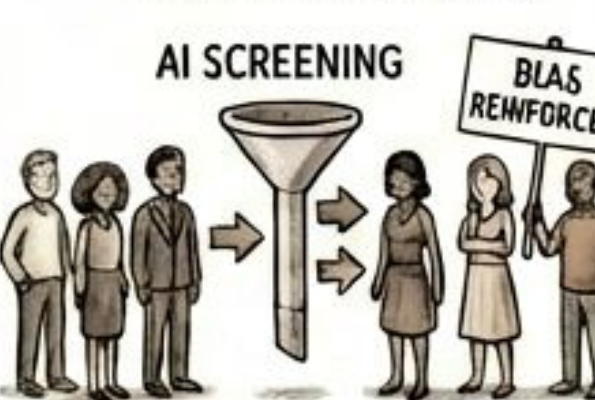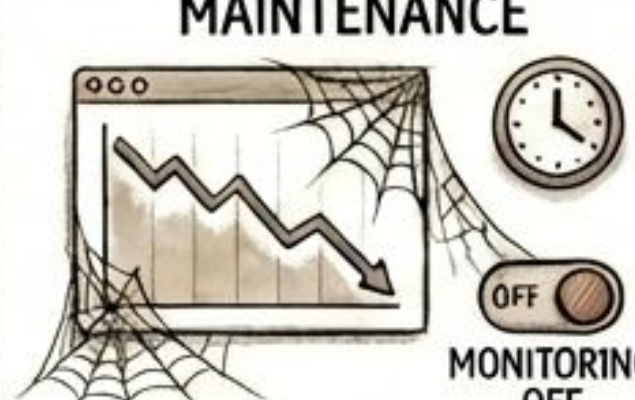# 10 Ways Your Data Project Might Fail

- Lack of Clear **Objectives**

- Insufficient **Data Quality**

- Poor **Data Management**

- Inadequate **Skills** & **Expertise**

- Ignoring **Business Context**

- **Overcomplicating** Solutions

- Poor **Communication** & **Collaboration**

- Ignoring **Ethical Considerations**

- Lack of **Scalability**

- Insufficient **Monitoring** & **Maintenance**

Source: Adapted from Martin Goodson, "Ten Ways Your Data Project Is Going To Fail"

10

# TEN WAYS YOUR DATA PROJECT IS GOING TO FAIL

Source: Martin Goodson's Analysis. *Pitfalls for Data Science Practitioners.*

## 1 LACK OF CLEAR OBJECTIVES

**Example:** Nonprofit aims to 'increase donor engagement' without defining metrics (frequency vs. size vs. retention).
**Result:** Generic reports, dissatisfied stakeholders.

## 2 INSUFFICIENT DATA QUALITY

**Example:** Healthcare predictive model for at-risk patients relies on inconsistent EHRs (missing fields, errors).
**Result:** Unreliable predictions, ineffective project.

## 3 POOR DATA MANAGEMENT

SILOED
DUPLICATED
UNTAGGED

**Example:** Media company has terabytes of unstructured, siloed customer data without governance.
**Result:** Impossible to extract insights for recommendation system.

## 4 INADEQUATE SKILLS & EXPERTISE

**Example:** Retail chain hires data scientist without domain expertise.
**Result:** Complex, impractical model incompatible with inventory workflows.

## 5 IGNORING BUSINESS CONTEXT

FRAUD DETECTED

**Example:** Bank's accurate fraud model ignores manual approval process.
**Result:** Overwhelming false positives, delayed transactions, frustrated staff.

## 6 OVERCOMPLICATING SOLUTIONS

DEEP LEARNING ROUTE OPTIMIZER

HEURISTIC

**Example:** Logistics company uses complex deep learning when simple heuristic could achieve 90% benefit.
**Result:** unsustainable resource & expertise requirements.

## 7 POOR COMMUNICATION & COLLABORATION

DATA TEAM — MARKETING TEAM

**Example:** E-commerce customer segmentation model built without marketing input.
**Result:** Model too technical, misaligned with campaigns, underutilized.

## 8 IGNORING ETHICAL CONSIDERATIONS

AI SCREENING
BIAS REINFORCED

**Example:** Recruitment AI uses historical hiring data, reinforcing bias against underrepresented groups.
**Result:** Public backlash, reputational damage.

## 9 LACK OF SCALABILITY

GROWING DATA VOLUME

**Example:** Startup's personalized email system on a simple database fails as customer base grows.
**Result:** Delays, failures, frustrated customers.

## 10 INSUFFICIENT MONITORING & MAINTENANCE

MONITORING OFF

**Example:** Social media recommendation engine deployed without performance monitoring.
**Result:** Relevance drift due to behavior changes, decreased engagement.

**Course Note:** Success requires balancing technical rigor with clear goals, domain context, and ongoing stewardship. Avoid these traps!

## *CRISP-DM* covers a lot of ground, but must be complemented by strong communication, ethical considerations, and scalability planning.

| CRISP-DM Stage | What Pitfall does it help mitigate? |
| --- | --- |
| Business Goal | 1: Clear Objectives & ROI - Defines success metrics upfront. |
| Assess Data | 2: Data Quality - Identifies and addresses data issues early. |
| Prepare Data | 2: Data Quality - Cleans & transforms data for reliability. 3: Data Management - Supports documentation and governance. |
| Modeling | 4: Expertise - Emphasizes model validation & performance. 6: Overcomplicating - Promotes pragmatic solutions. |
| Evaluation | 10: Monitoring - Tests model with holdout data prior to deployment. |
| Deployment | 3: Data Management - Supports integration. 5: Context - Integrates into business processes.  10: Monitoring - Provides framework for monitoring after integration. |

# DATA + COMPUTE

necessary / not sufficient

# MODELING + OPTIMIZATION

deliver the value

|  | | Human Decision Making | |
|---|---|---|---|
|  | | Augment - Narrow | Replace - General |
| **Optimizing For …** | More | Skin lesion Classification > Correct Diagnosis | Trading Bots > profits |
| | Less | Delivery Route Planning < Fuel use | Self driving car < Accidents |

**Key Reason DIA are important: Data 2 Optimal Decisions**

| | | Decision Frequency | |
|---|---|---|---|
| | | High | Low |
| **Decision Impact** | High | **Prime Value Target** E.g. medical diagnosis and autonomous vehicles | **Mixed Bag** |
| | Low | **May be worth the investment** E.g. route planning and recommendations | **Forget It!** Too Costly |

**Key Reason DIA are important: Data 2 Optimal Decisions**

Save ~100M driving miles and $350–400M per year with ORION route optimization system

| | | Decision Making | |
|---|---|---|---|
| | | Augment | Replace |
| **Optimizing For …** | More | | |
| | Less | **X** Routes - Fuel | |

| | | Decision Frequency | |
|---|---|---|---|
| | | High | Low |
| **Decision Impact** | High | 🟩 | 🟨 |
| | Low | **X** aggregation | 🟥 |

Planners and Drivers frequent decision process (every day). Low impact on a given delivery although aggregation across the high number of deliveries.

**Outcome**
Use Less fuel (KPI) by suggesting routes for each days deliveries.

Example:  UPS Route Optimization

|  |  | Decision Making | |
| --- | --- | --- | --- |
|  |  | Augment | Replace |
| **Optimizing For …** | More | **X** Oil Produced | |
|  | Less | **X** GHG | |

|  |  | Decision Frequency | |
| --- | --- | --- | --- |
|  |  | High | Low |
| **Decision Impact** | High | **X** | |
|  | Low | | |

*__Augmenting__*
Oil field operators with "digital twin" simulations leading to optimizations of operational parameters (flow and pressure) across the large scale operations.

**Outcome**
Less GHG and higher production. (KPI)

Example:  BP's New Oilfield Roughneck Is An Algorithm

Increased booking conversion rate by ~4% by modeling likelihood of host acceptance

|  |  | Decision Making | |
|---|---|---|---|
|  |  | Augment | Replace |
| **Optimizing For …** | More |  | X (rec) |
|  | Less |  |  |

|  |  | Decision Frequency | |
|---|---|---|---|
|  |  | High | Low |
| **Decision Impact** | High | 🟩 | 🟨 |
|  | Low | X (agg) 🟨 | 🟥 |

**Outcome**

More booking revenue (KPI) by recommending listings that are likely to accept the renter. Beware of ethical considerations of algorithmic bias!

**_Replace_**

Replacing travel agent. Low impact on a given listing although aggregation across the high number of listings.

Example:  AirBnB Recommendation - increase booking revenue

**NETFLIX**

Generate more than 80% of content views through ML recommendation system

| | | Decision Making | |
|---|---|---|---|
| | | Augment | Replace |
| **Optimizing For …** | More | | **X** (rec) |
| | Less | | |

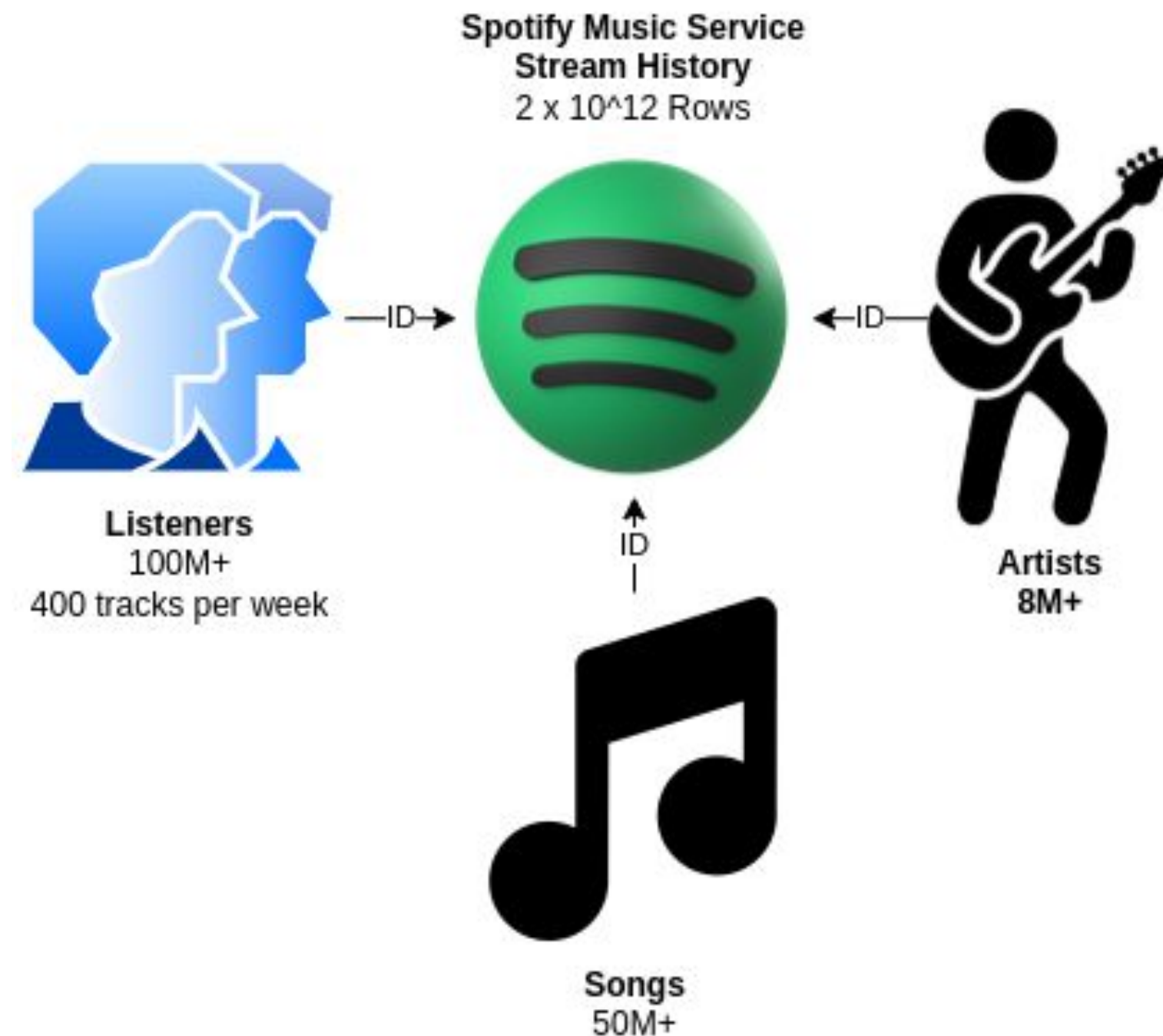| | | Decision Frequency | |
|---|---|---|---|
| | | High | Low |
| **Decision Impact** | High | <span style="background-color:green"> </span> | <span style="background-color:yellow"> </span> |
| | Low | <span style="background-color:yellow">**X** (agg)</span> | <span style="background-color:red"> </span> |

### *Replacing*
Search results based on user entered criteria. Low impact on a given search although aggregation across the high number of searches.

### **Outcome**
Increase engagement - watch time (KPI) by recommending movies.

Example:  Netflix Recommendation - 80% of views

# A DATA & COMPUTE STORY

## How Spotify ran the largest Google Dataflow job ever for Wrapped 2019



Spotify Music Service
Stream History
2 x 10^12 Rows

Listeners
100M+
400 tracks per week

Artists
8M+

Songs
50M+

```
SELECT
    L.UserName,
    A.ArtistName,
    S.Track,
    COUNT(1) as ListenCount
FROM Stream_History SH
JOIN Artists A on A.ID = SH.AID
JOIN Songs S on S.ID = SH.SID
JOIN Listeners L on L.ID = SH.LID
GROUP BY L.UserName, A.ArtistName, S.Track
ORDER BY ListenCount DESC
```

# Spotify Wrapped is a data-intensive application designed to:

| | | Decision Making | |
|---|---|---|---|
| | | Augment | Replace |
| **Optimizing For …** | More | X (engagement) | |
| | Less | | |

| | | Decision Frequency | |
|---|---|---|---|
| | | High | Low |
| **Decision Impact** | High | 🟩 | 🟨 X (agg) |
| | Low | 🟨 | 🟥 |

**Outcome**

Strengthen emotional connections to the platform.
Drive business value through engagement, retention, and brand promotion.

**_Augement_**

Increases the users self-awareness of their listening habits Analytics for the people.

30

# That's a Wrap!

Get going on your Labs!