



The Data Flywheel

(v1.00)

Week 1: September 4, 2025

Jimmy Lin
David R. Cheriton School of Computer Science
University of Waterloo

These slides are available at <https://lintool.github.io/cs451-2025f/>

This work is licensed under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International
See <https://creativecommons.org/licenses/by-nc-sa/4.0/> for details



Who am I?



What's this course about?



What does it mean to an AI-first company?
What does it mean to a data-driven company?





When people talk about AI these days,
they really mean ML...

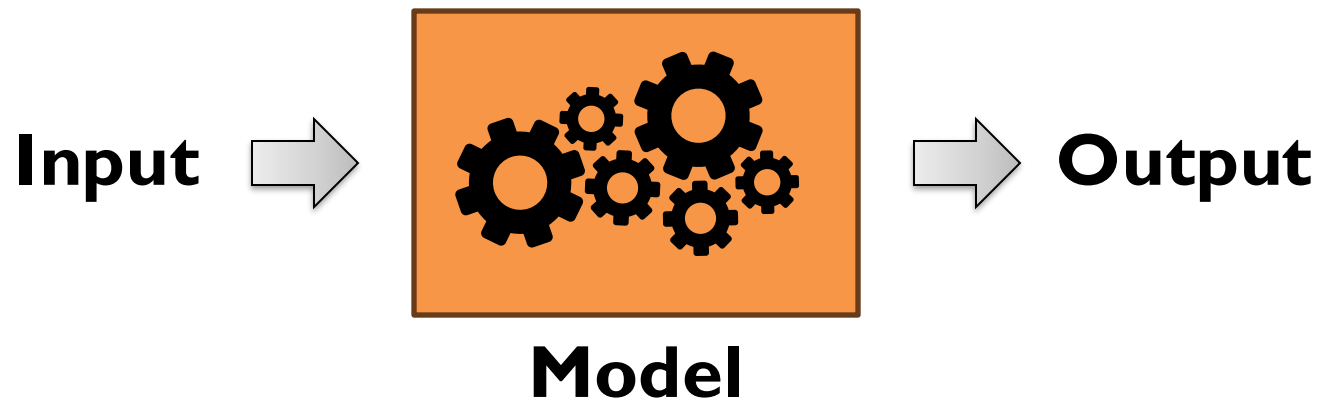




When people talk about ML these days,
they really mean supervised ML...



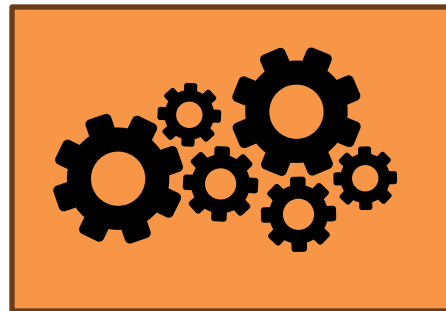
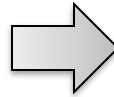
What's supervised machine learning?



(accomplishes *some* task)
(hopefully *well*)

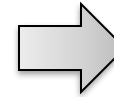


Input



Model

(accomplishes some task)
(hopefully well)



Output

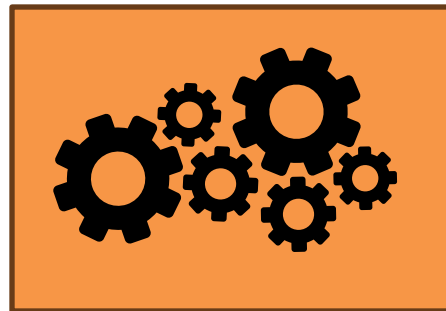
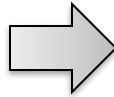
owl

cat

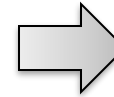


amazing spot for good food & a fun time
🍕🍹 they offer a super unique dine-in
experience with their interactive tables!
also love that they have innovative
weekly feature dishes 😊

Input



Model



Output

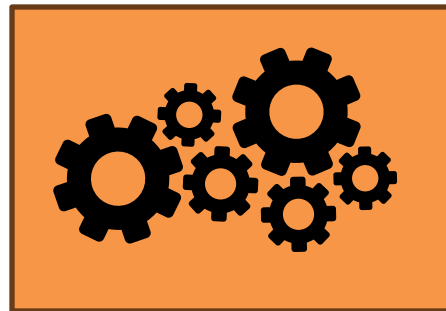
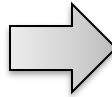


Worst service I have ever experienced!
Waited 25 mins for our waitress to
come to our table once seated, no
apology! Waited another 30 mins for our
drinks, which we had to ask about.

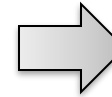
(accomplishes some task)
(hopefully well)

amazing spot for good food & a fun time
🍕🍹 they offer a super unique dine-in
experience with their interactive tables!
also love that they have innovative
weekly feature dishes 😊

Input



Model



Output



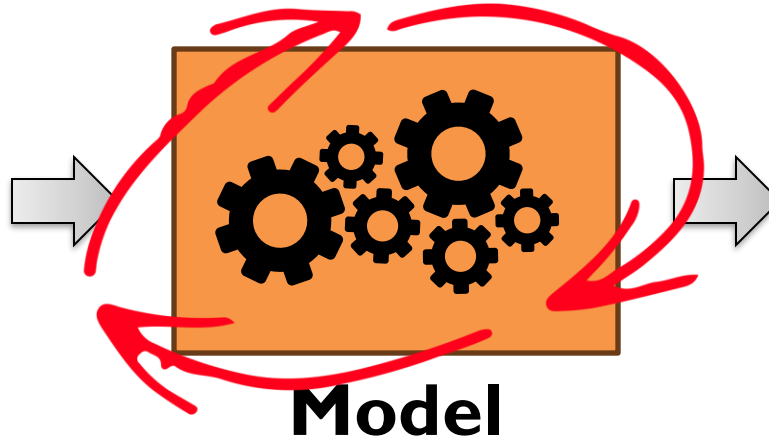
Worst service I have ever experienced!
Waited 25 mins for our waitress to
come to our table once seated, no
apology! Waited another 30 mins for our
drinks, which we had to ask about.

How do you build such a model?

Model *learns* from the data

(review, 👍)
(review, 👎)
(review, 👍)
(review, 👍)
(review, 👎)
(review, 👎)

Input



Model

Output

👍
👎

Machine learning algorithm
adjusts the model *parameters*

How do you build such a model?

More details later in the course...



What's supervised machine learning?



What does it mean to an AI-first company?
What does it mean to a data-driven company?

AI ~means ML

ML ~means supervised ML

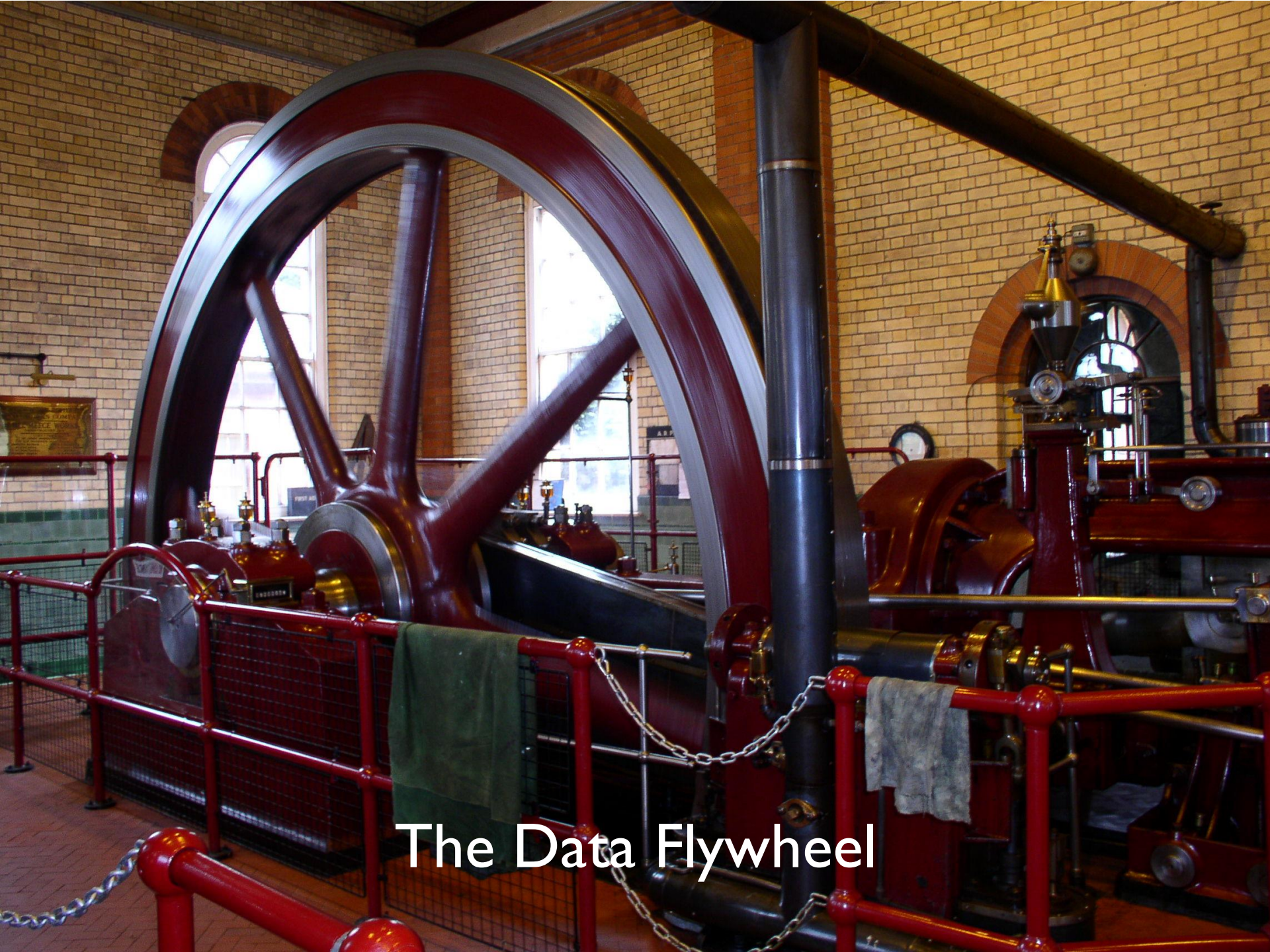
supervised ML requires lots of data

data is what powers AI

AI is what you “do” with data

Where does the data come from?

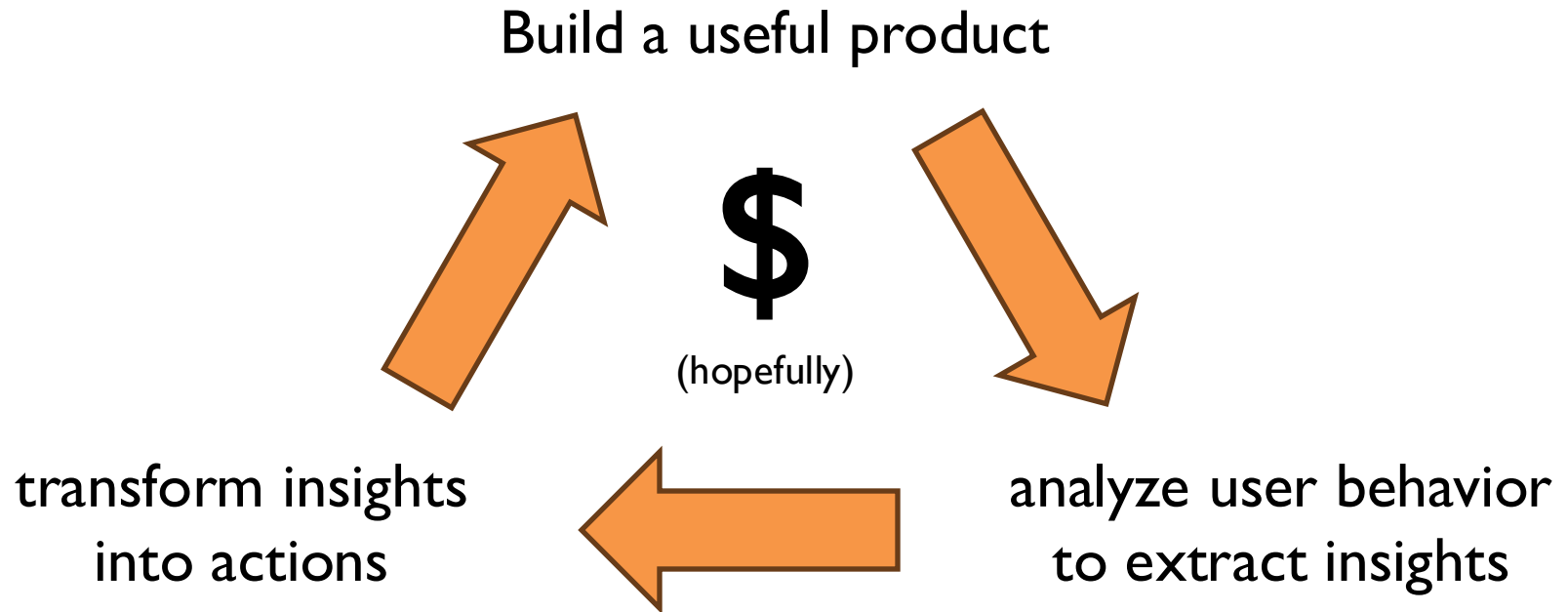
What's this course about?



The Data Flywheel

The Data Flywheel

(a virtuous cycle)



Google. Facebook. Twitter. Amazon. Uber.

Business Intelligence

An organization should retain data that result from carrying out its mission and exploit those data to generate insights that benefit the organization, for example, market analysis, strategic planning, decision making, etc.

Duh!?

This is not a new idea...

Case Study

In the 1990s, Wal-Mart found that customers tended to buy diapers and beer together. So they put them next to each other and increased sales of both.*

Why didn't they do it earlier?

Relational databases weren't invented until 1970

Computers were slow and expensive

* BTW, this is completely apocryphal. (But it makes a nice story.)



5 MB hard drive in 1956

Case Study

In the 1990s, Wal-Mart found that customers tended to buy diapers and beer together. So they put them next to each other and increased sales of both.*

So what's changed?

More compute and storage

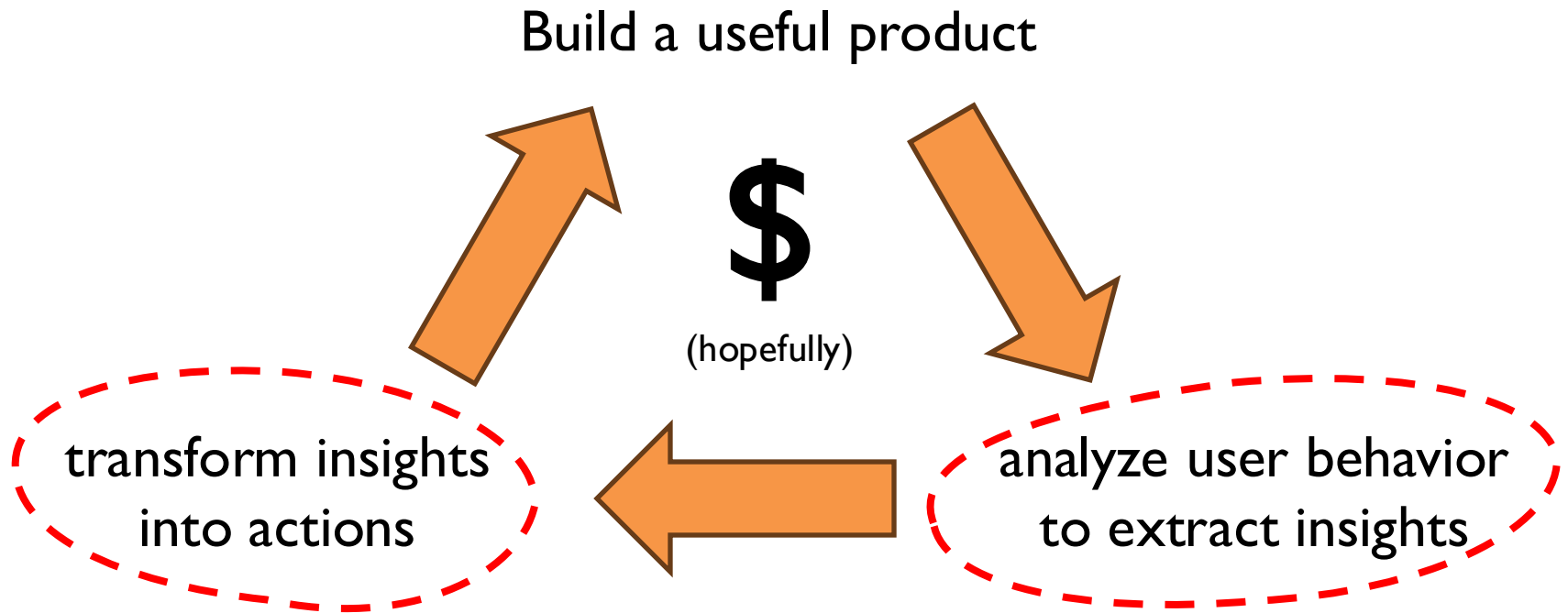
Ability to gather behavioral data

Advent of machine learning

* BTW, this is completely apocryphal. (But it makes a nice story.)

The Data Flywheel

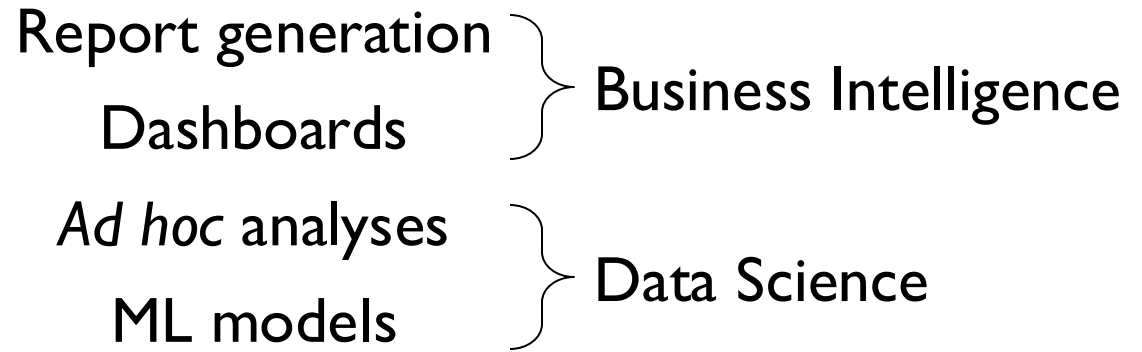
(a virtuous cycle)



Google. Facebook. Twitter. Amazon. Uber.

Transform Insights into Actions

What does that really mean?



What's this course about?

The *infrastructure* that supports the data flywheel.

data platforms

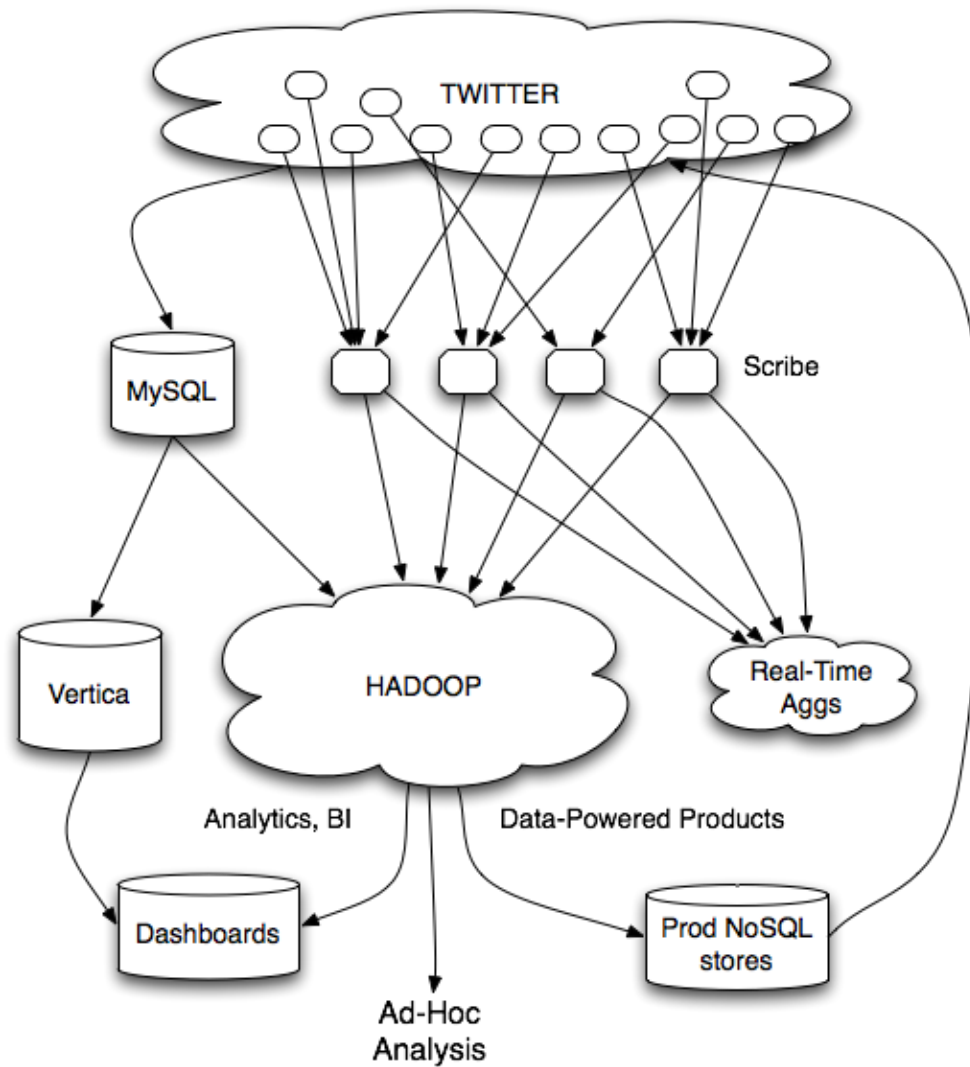
lakehouse = data warehouse + data lake

data engineering

What problems do data platforms solve?

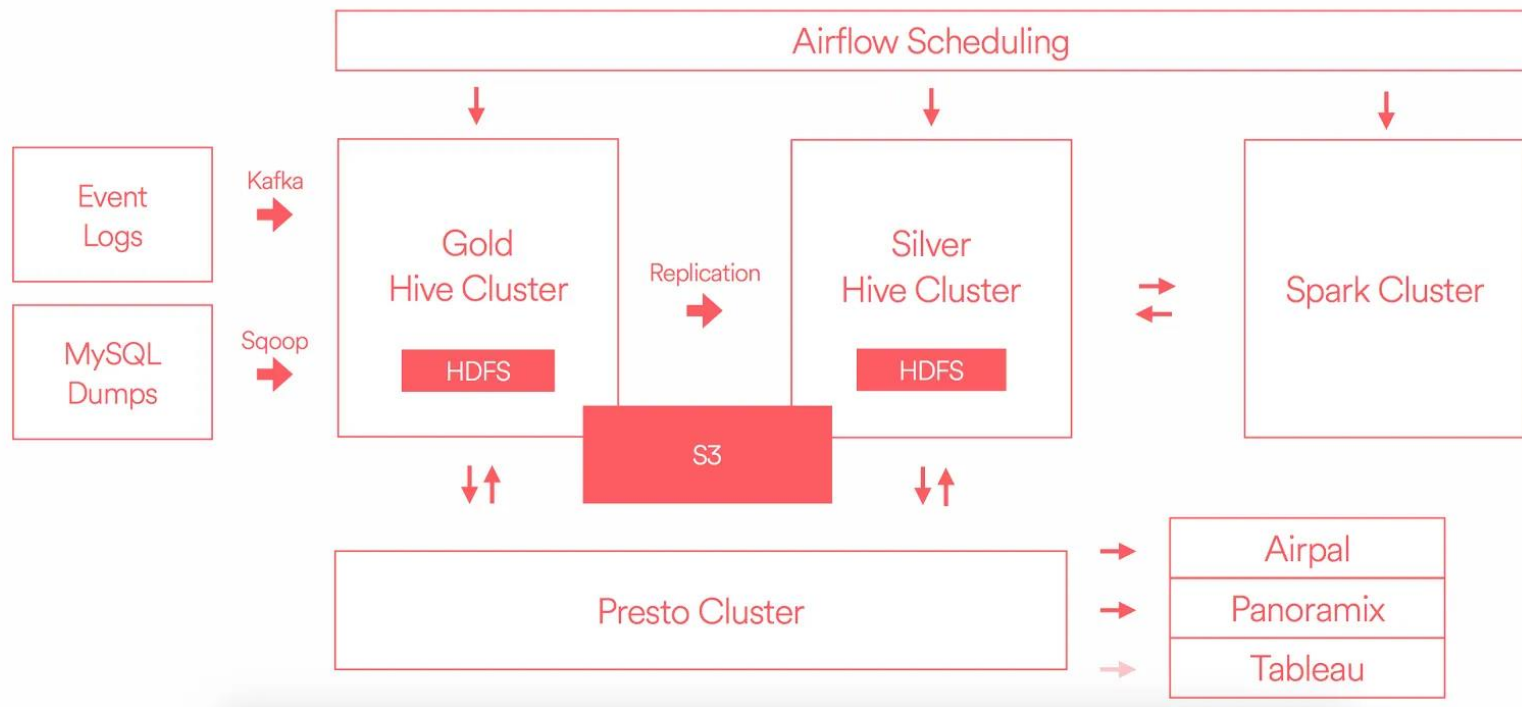
Ingesting, storing, manipulating, maintaining, serving...
the data that supports the data flywheel.

Some examples...



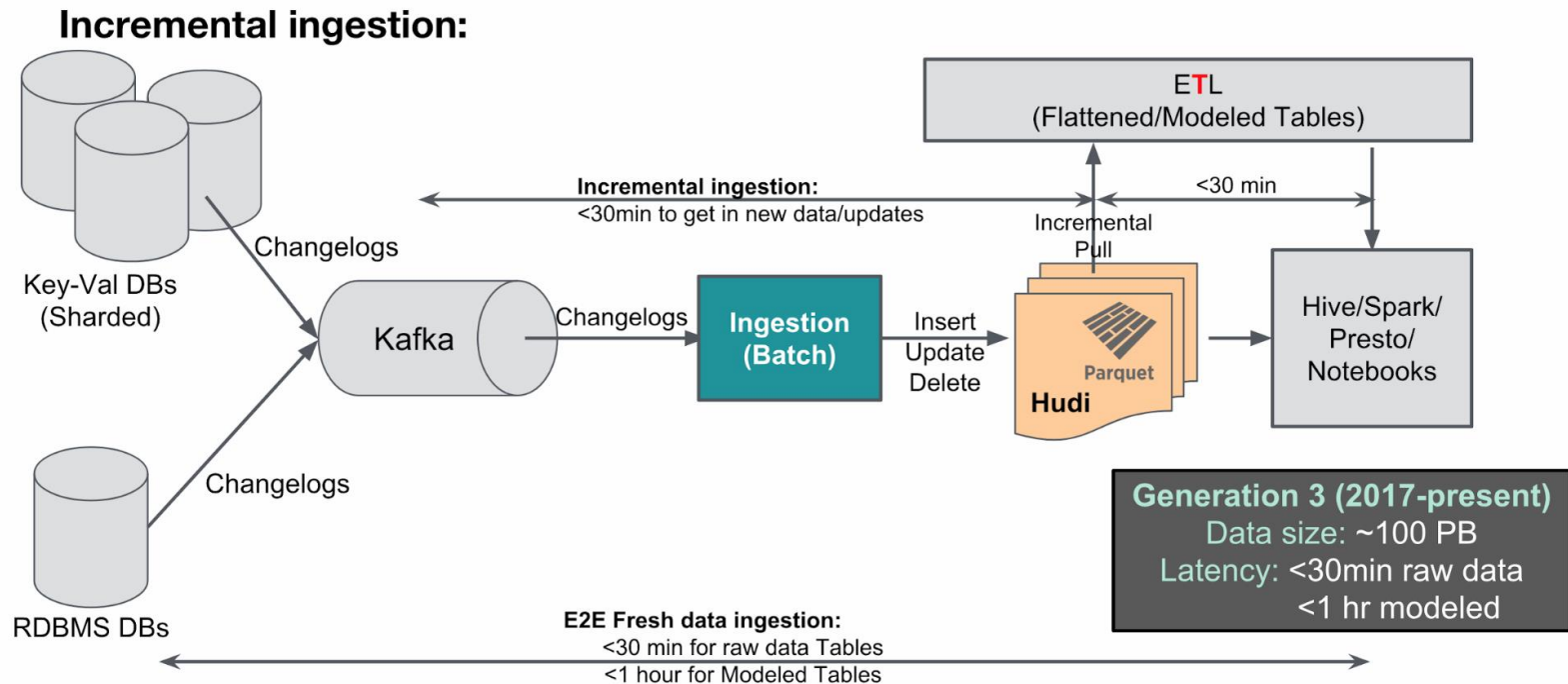
Twitter's data platform (circa 2012)

AIRBNB DATA INFRA



AirBnB's data platform (circa 2016)

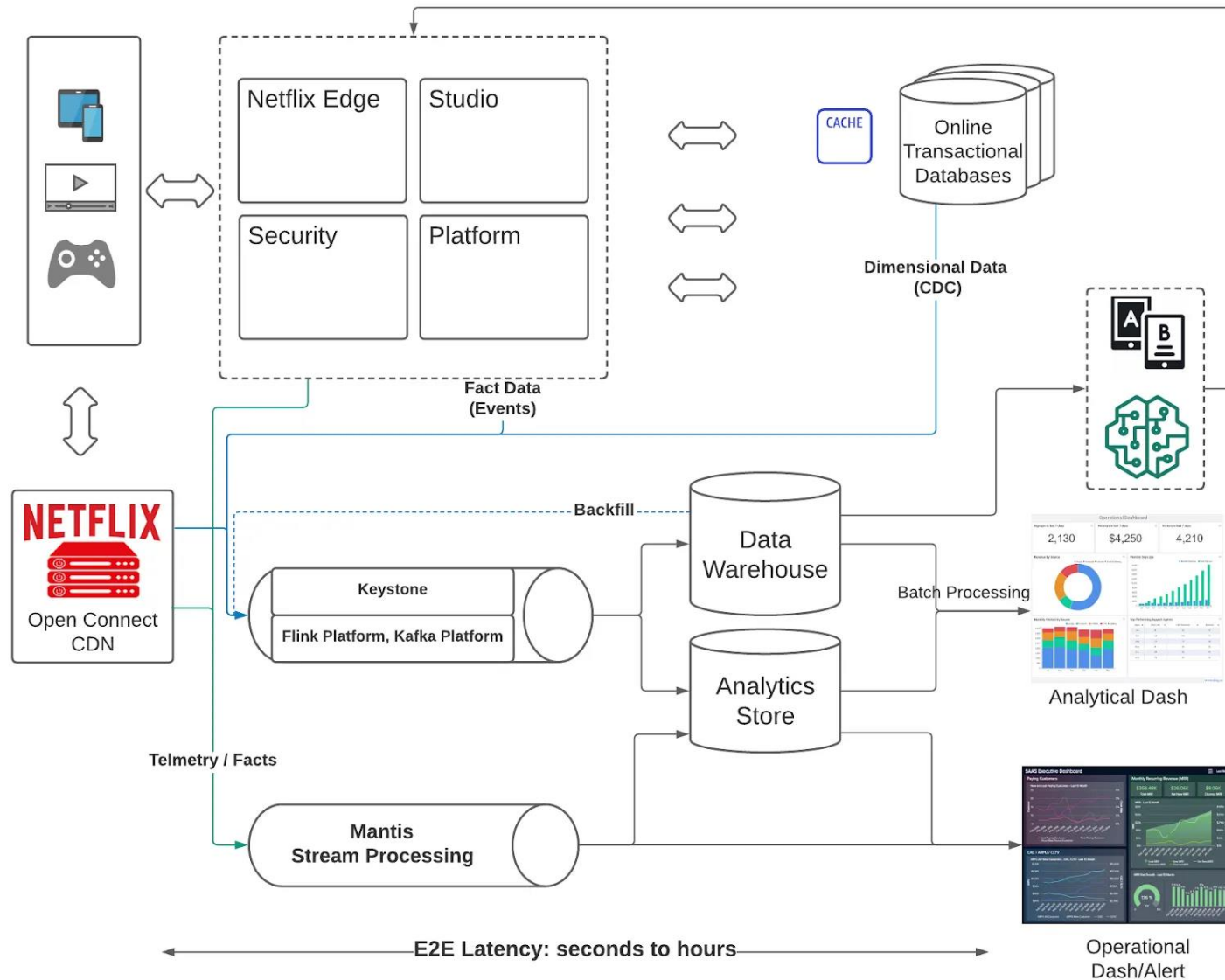
Generation 3 (2017-present) - Let's rebuild for long term



Uber's data platform (circa 2018)

<https://www.uber.com/en-CA/blog/uber-big-data-platform/>

How Stream Processing fit in Netflix (2021)



Netflix's data platform (circa 2021)



Spotify's data platform (circa 2024)

<https://engineering.atspotify.com/2024/5/data-platform-explained-part-ii>

What problem do data platforms solve?

Ingesting, storing, manipulating, maintaining, serving...
the data that supports the data flywheel.

By the end of the course you should be able
to “make sense” of any data platform.

Why is this a difficult problem?

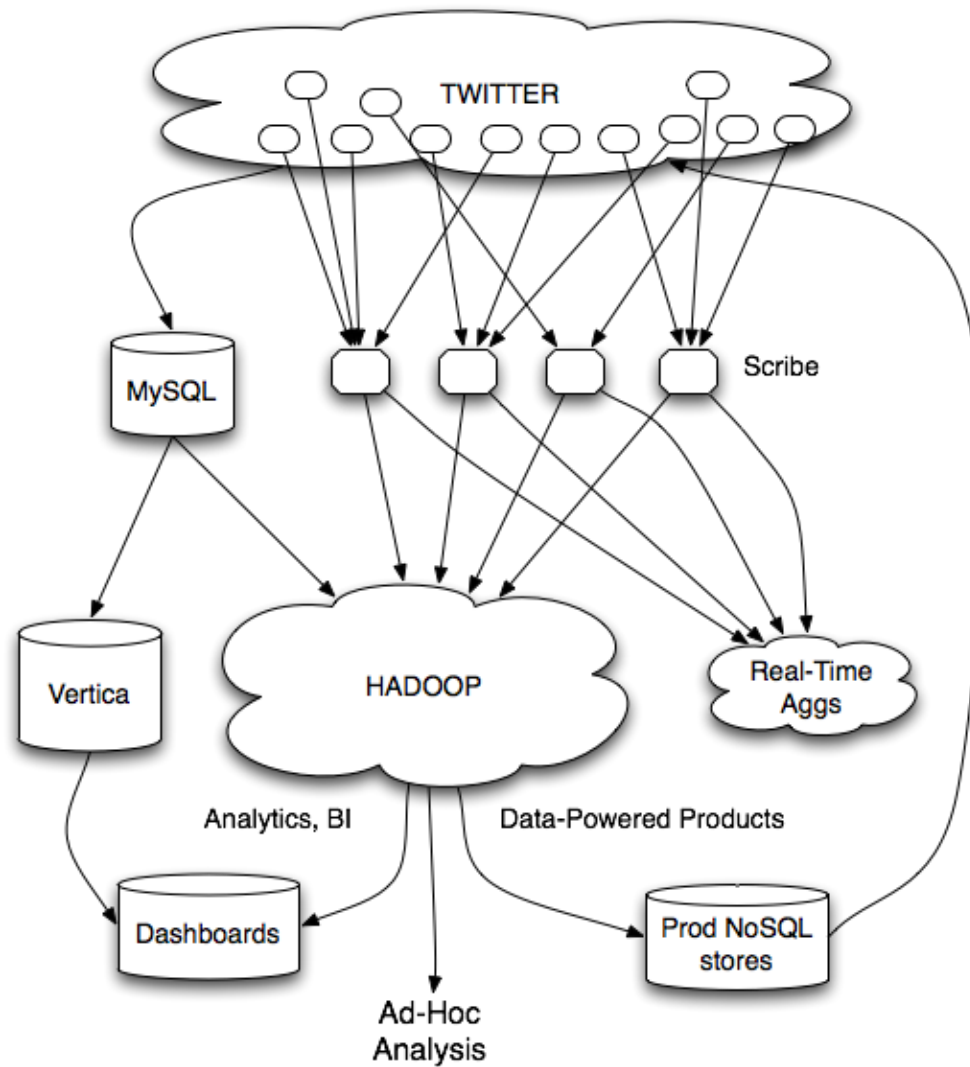
Characteristics of the data: Volume, Velocity, Variety* + Veracity

* Coined by Gartner analyst Doug Laney in 2001.

Why is this a difficult problem?

Characteristics of the data: Volume, Velocity, Variety* + Veracity

* Coined by Gartner analyst Doug Laney in 2001.



Twitter's data platform (circa 2012)

circa ~2010

~150 people total

~60 Hadoop nodes

~6 people use analytics stack daily

circa ~2012

~1400 people total

10s of Ks of Hadoop nodes, multiple DCs

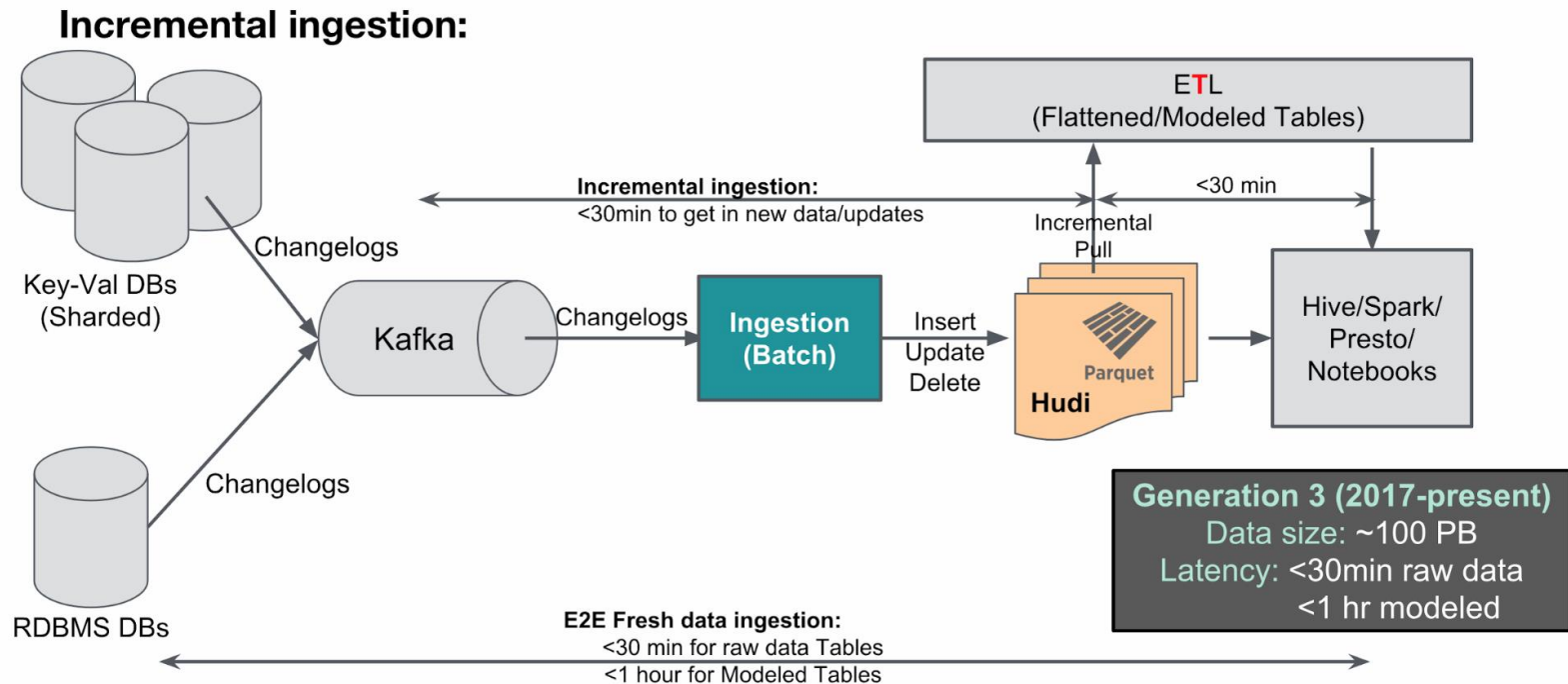
10s of PBs total Hadoop DW capacity

~100 TB ingest daily

dozens of teams use Hadoop daily

10s of Ks of Hadoop jobs daily

Generation 3 (2017-present) - Let's rebuild for long term



Uber's data platform (circa 2018)

<https://www.uber.com/en-CA/blog/uber-big-data-platform/>

Why is this a difficult problem?

Characteristics of the data: Volume, Velocity, Variety* + Veracity

Want more?

<https://lintool.github.io/my-data-is-bigger-than-your-data/>

* Coined by Gartner analyst Doug Laney in 2001.

Why is this a difficult problem?

Characteristics of the data: Volume, Velocity, Variety* + Veracity

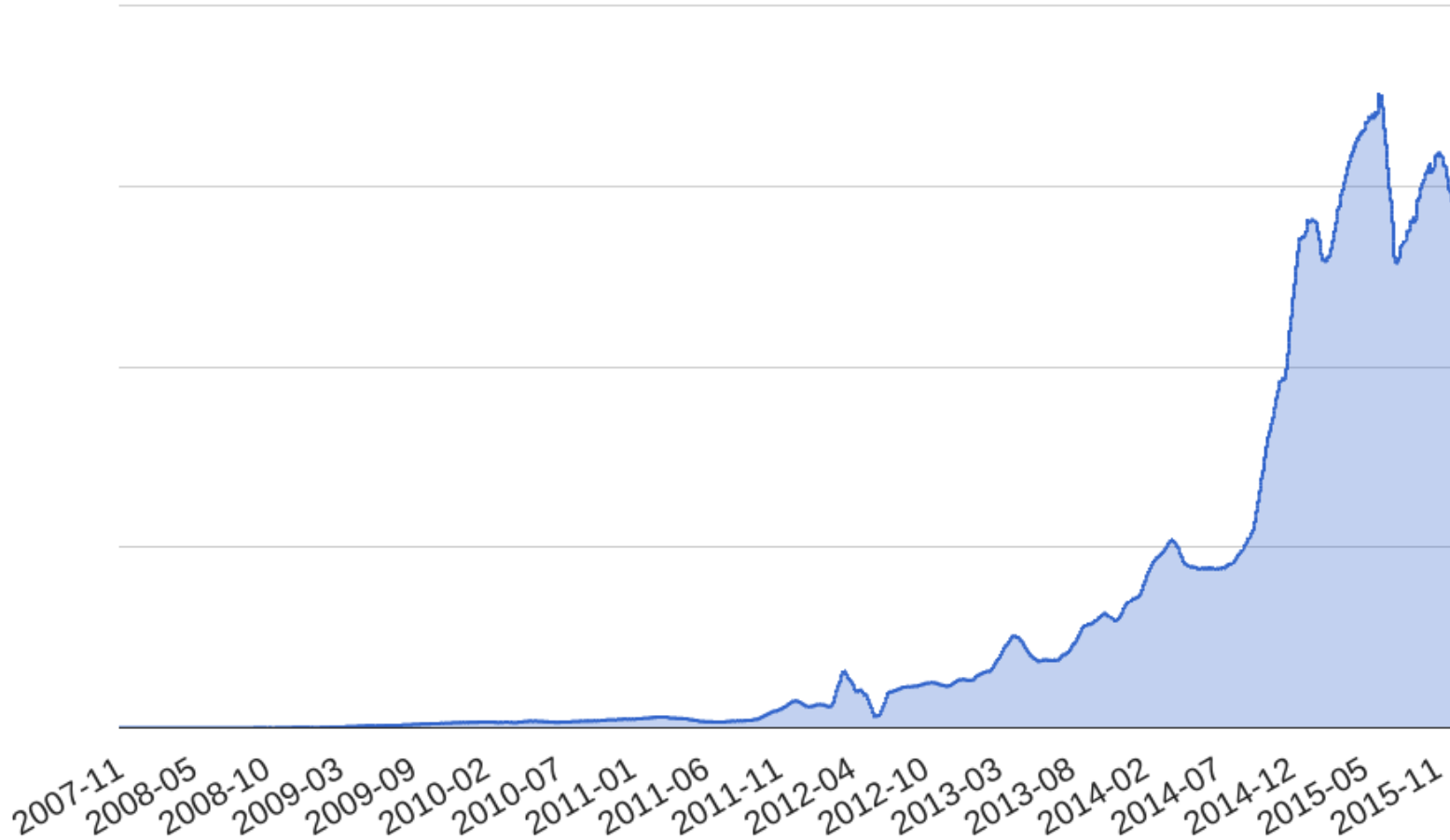
* Coined by Gartner analyst Doug Laney in 2001.



Spotify's data platform (circa 2024)

<https://engineering.atspotify.com/2024/5/data-platform-explained-part-ii>

When it comes to scalability, Spotify's Data Collection platform collects more than 1 trillion events per day.
(circa 2024)



Amount of Spotify's Delivered Events over time (circa 2016)

Why is this a difficult problem?

Characteristics of the data: Volume, Velocity, Variety* + Veracity

Want more?

<https://lintool.github.io/my-data-is-bigger-than-your-data/>

* Coined by Gartner analyst Doug Laney in 2001.

Why is this a difficult problem?

Characteristics of the data: Volume, Velocity, Variety* + Veracity

Previously – relational

Today – semi-structured, unstructured, graph, multimodal, etc.

* Coined by Gartner analyst Doug Laney in 2001.

Why is this a difficult problem?

Characteristics of the data: Volume, Velocity, Variety* + Veracity

GI / GO

* Coined by Gartner analyst Doug Laney in 2001.

Why is this a difficult problem?

Characteristics of the data: Volume, Velocity, Variety* + Veracity

By the end of the course you should be able understand
and appreciate the challenges here.

* Coined by Gartner analyst Doug Laney in 2001.

Data Engineering “Truths”

(That are rarely acknowledged explicitly)

Data engineering is practical

No solutions, only tradeoffs (but best practices exist)

Data platforms...

already exist in some form

are complex, messy, and evolving

combination of technologies and processes

subjected to technical, organizational, and business constraints

What's this course about?



The Data Flywheel

What's this course *not* about?

machine learning

data science

distributing processing

building database applications

Course Mechanics

There are two separate sections, beware!
Difference?

Back to this...

Data Engineering “Truths”

(That are rarely acknowledged explicitly)

Data engineering is practical

No solutions, only tradeoffs (but best practices exist)

Data platforms...

already exist in some form

are complex, messy, and evolving

combination of technologies and processes

subjected to technical, organizational, and business constraints

Hence...

My Teaching Philosophy

Focus on concepts and intuition

Focus on reasoning about tradeoffs

Less emphasis on detailed algorithms
(don't confuse with non-technical)

Less emphasis on specific pieces of software
(how concepts are operationalized – in passing)
(how to actually use x – on your own)

Course Materials

All materials will be posted on the course homepage

You are responsible for keeping up to date.

“I didn’t know” or “I didn’t see it” won’t be accepted as an excuse.

Readings will be assigned from textbooks and papers

You are responsible for the readings.

Assigned by week... when to do it?

Generative AI

tl;dr – Allowed. Encouraged.
(if used appropriately)

What's appropriate use?
(Let me know when you figure it out...)

Okay to use in same way as a search engine.
Not okay to wholesale copy and paste “do my assignment”.

Okay to use as a learning tool.
Not okay to use as a crutch to speed run without learning.

Course Grade

six assignments (CS 45I/65I)

GitHub to manage logistics of submissions

midterm/final exam (CS 45I/65I)

final project (CS 65I only)

Back to the beginning...



What does it mean to an AI-first company?
What does it mean to a data-driven company?

AI ~means ML

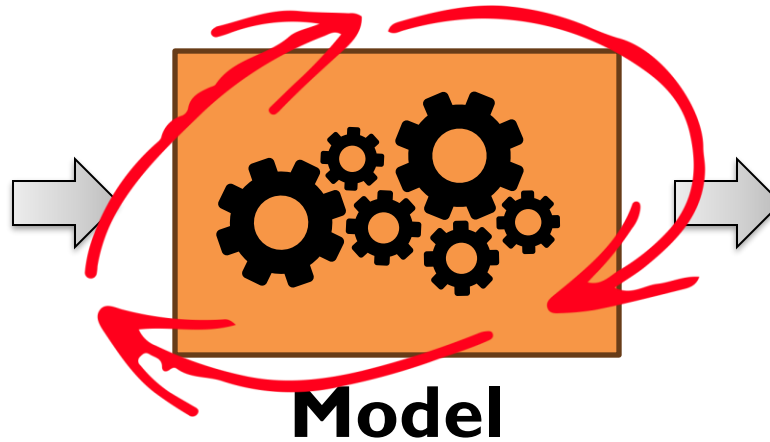
ML ~means supervised ML

supervised ML requires lots of data

Model *learns* from the data

(review, 👍)
(review, 👎)
(review, 👍)
(review, 👍)
(review, 👎)
(review, 👎)

Input



Output

👍
👎

Machine learning algorithm
adjusts the model *parameters*

Where does the data *actually* come from?

富嶽三十六景 神奈川
浪裏

Stay tuned...

