Knowledge from Nothing

# Creating Data Pipelines for User Insights from Zero

JEFFREY LEUNG

# A DAY IN DATA

The exponential growth of data is undisputed, but the numbers behind this explosion – fuelled by internet of things and the use of connected devcies – are hard to comprehend, particularly when looked at in the context of one day

## DEMYSTIFIYING DATA UNITS

From the more familiar 'bit' or 'megabyte', larger units of measurement are more frequently being used to explain the masses of data

| Unit | | Value | Size |
|------|------|-------|------|
| b | bit | 0 or 1 | 1/8 of a byte |
| B | byte | 8 bits | 1 byte |
| KB | kilobyte | 1,000 bytes | 1,000 bytes |
| MB | megabyte | $1,000^2$ bytes | 1,000,000 bytes |
| GB | gigabyte | $1,000^3$ bytes | 1,000,000,000 bytes |
| TB | terabyte | $1,000^4$ bytes | 1,000,000,000,000 bytes |
| PB | petabyte | $1,000^5$ bytes | 1,000,000,000,000,000 bytes |
| EB | exabyte | $1,000^6$ bytes | 1,000,000,000,000,000,000 bytes |
| ZB | zettabyte | $1,000^7$ bytes | 1,000,000,000,000,000,000,000 bytes |
| YB | yottabyte | $1,000^8$ bytes | 1,000,000,000,000,000,000,000,000 bytes |

*A lowercase "b" is used as an abbreviation for bits, while an uppercase "B" represents bytes.

## 463EB
of data will be created every day by 2025
IDC

## 500m
tweets are sent every day
Twitter

## 4PB
of data created by Facebook, including

**350m** photos

**100m** hours of video watch time

Facebook Research

## 95m
photos and videos are shared on Instagram
Instagram Business

## 320bn
emails to be sent each day by 2021

## 306bn
emails to be sent each day by 2020

## 294bn
billion emails are sent
Radicati Group

## 65bn
messages sent over WhatsApp and two billion minutes of voice and video calls made
Facebook

## 3.9bn
people use emails

## 4TB
of data produced by a connected car
Intel

## 28PB
to be generated from wearable devices by 2020
Statista

Searches made a day — 5bn

Searches made a day from Google — 3.5bn

Smart Insights

## ACCUMULATED DIGITAL UNIVERSE OF DATA

**4.4ZB** 2013

**44ZB** 2020

PwC

RACONTEUR

# Why do you(r leaders) want data?

Analyze past performance.

Understand current performance.

Experiment on new features.

Make informed business decisions.

Machine Learning for new insights.

# Outline

# 0. the fundamentals

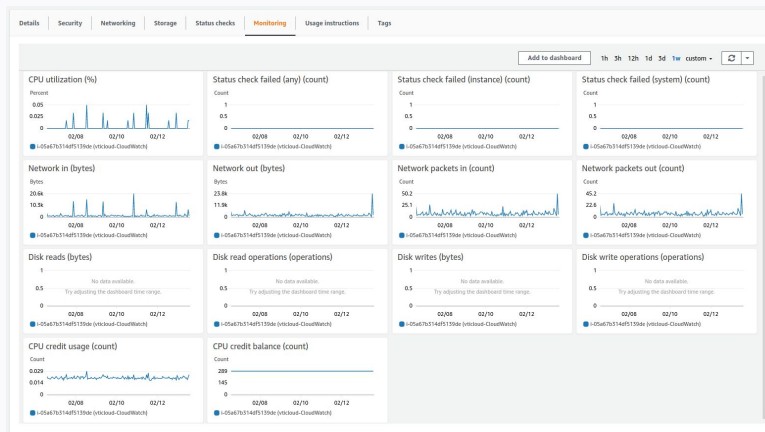# Where does it come from, where does it go?

# What kind of data?

**User behaviour**
(e.g. signups, interactions)

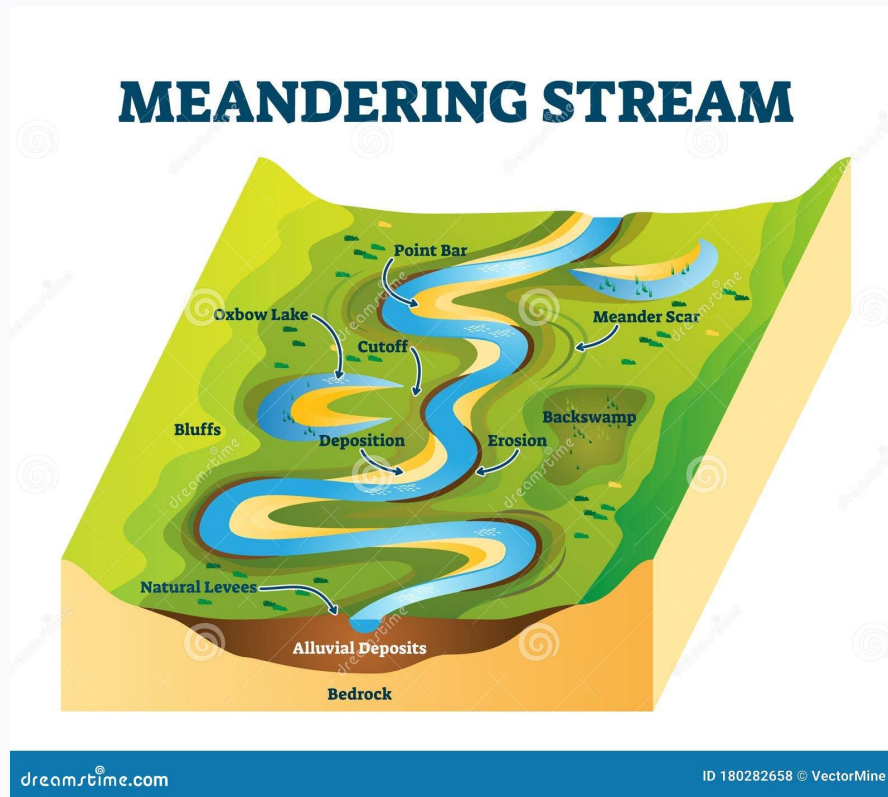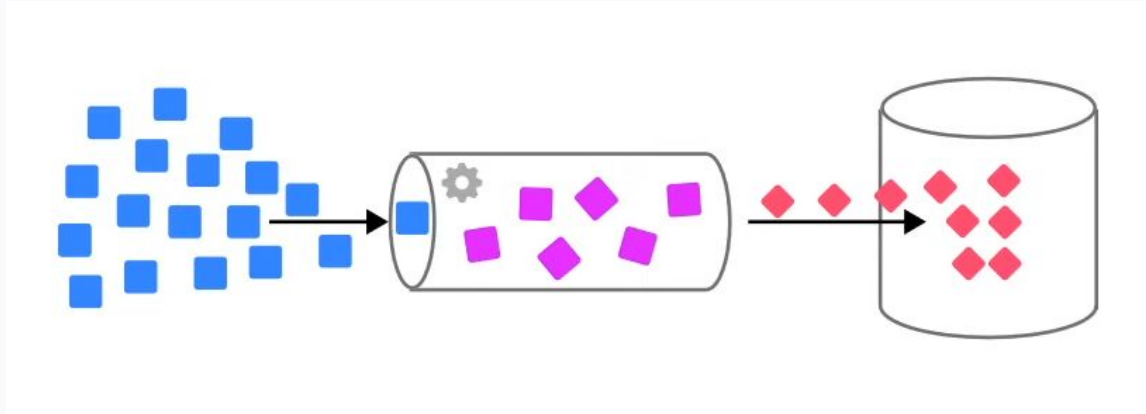**Operational metrics**
(e.g. cpu, memory, latency, network)

# Streaming Data Pipeline



MEANDERING STREAM
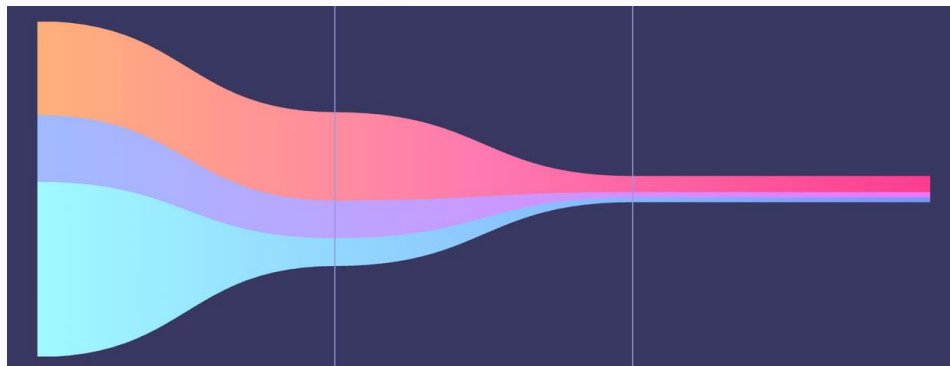
# Streaming Data Pipeline

# Batch Data Pipeline

# Data Pipeline Architecture

# The Dataflow



## 01. ingestion and storage

pseudonymization
anonymization
replayability
verification from source

## 02. cleaning and filtering

sanity checks
standardization
elimination of invalid data

## 03. aggregating and analyzing

averaging
determining patterns
alerting
p50/p99
diffs across days, weeks,
months, or years
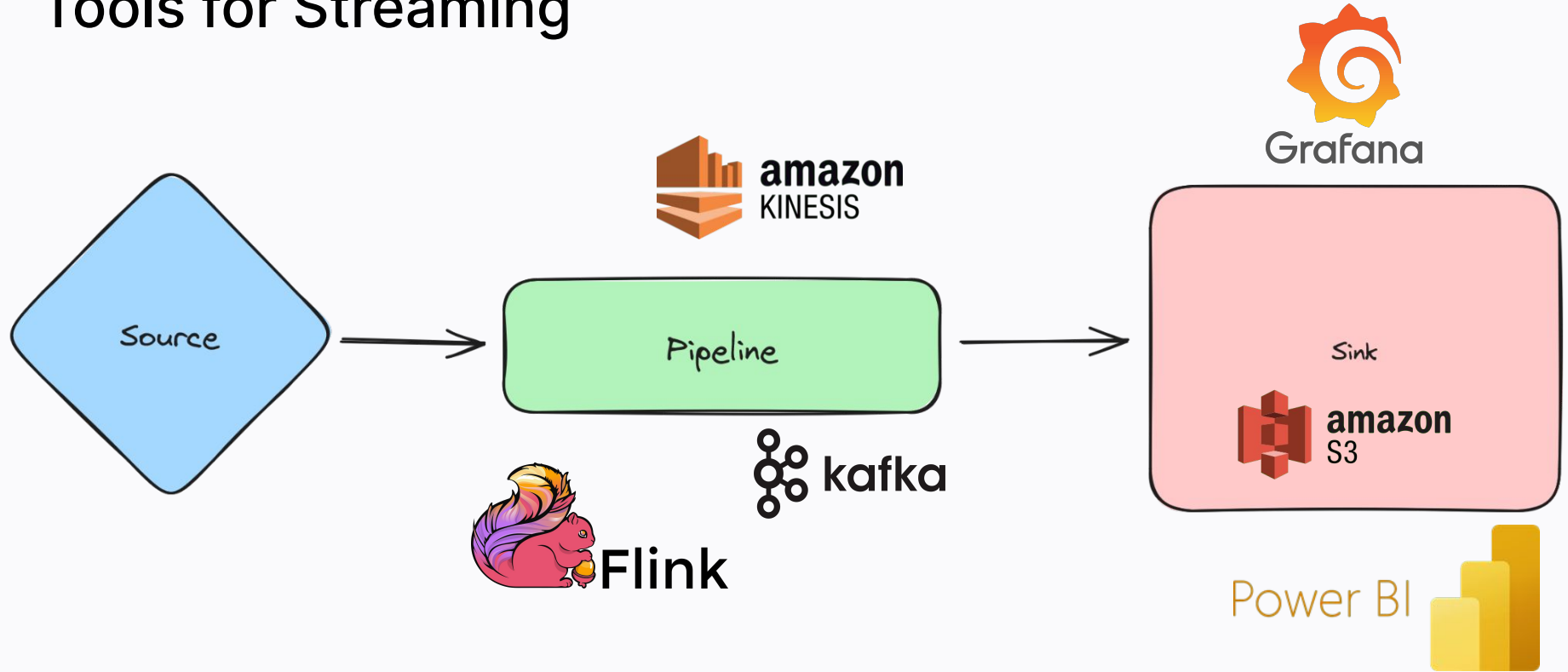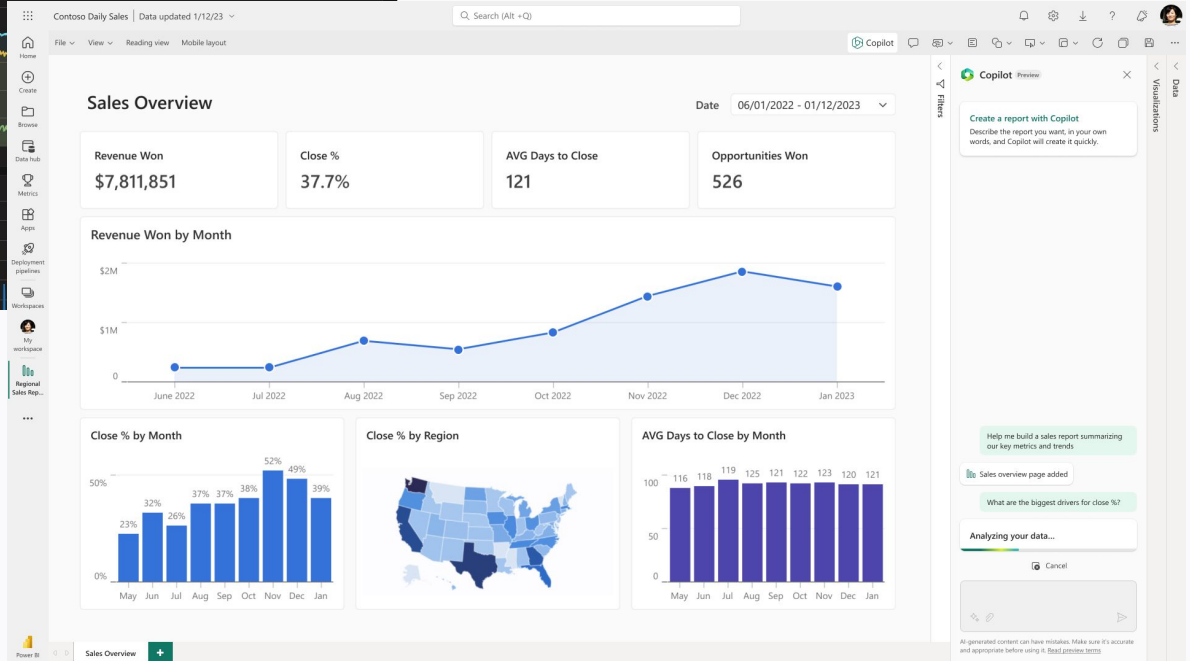graphing time series data
predictive analysis

# 1. the tools

# Tools for Streaming

# Tools for Batch Processing

# Data Warehouses

# Something Else: Mixpanel

# 2. the **why**

Who wants this data?

How do they want to consume it?

What decisions do they want to make with it?

What is this data intended to measure?

How do we measure what we actually want data on?

# Align on expected output.

# What **business value** does this data provide?

Is this the best way to provide it?

# Data is your **Product**.

Understand the requirements.

Validate the user need.

Design an MVP.

# Present **clearly**.

Use graphs and charts.

Compare for context.

p50, p90, p99.

# 3. data org maturity

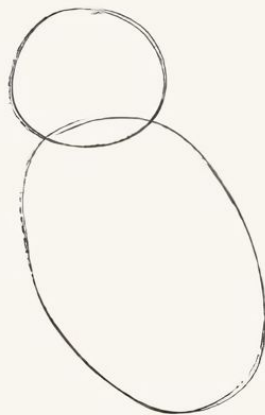# How to draw an Owl.

*"A fun and creative guide for beginners"*

Fig 1. Draw two circles

Fig 2. Draw the rest of the damn Owl

# Follow these people



**Zach Wilson**
Staff Data Engineer
Prev. Facebook, Netflix, Airbnb



**Benjamin Rogojan**
Seattle Data Guy
Prev. Facebook

# Get that data integrity

Compliance (GDPR, DSAR, COPPA, DMA) - anonymization, retention, deletion

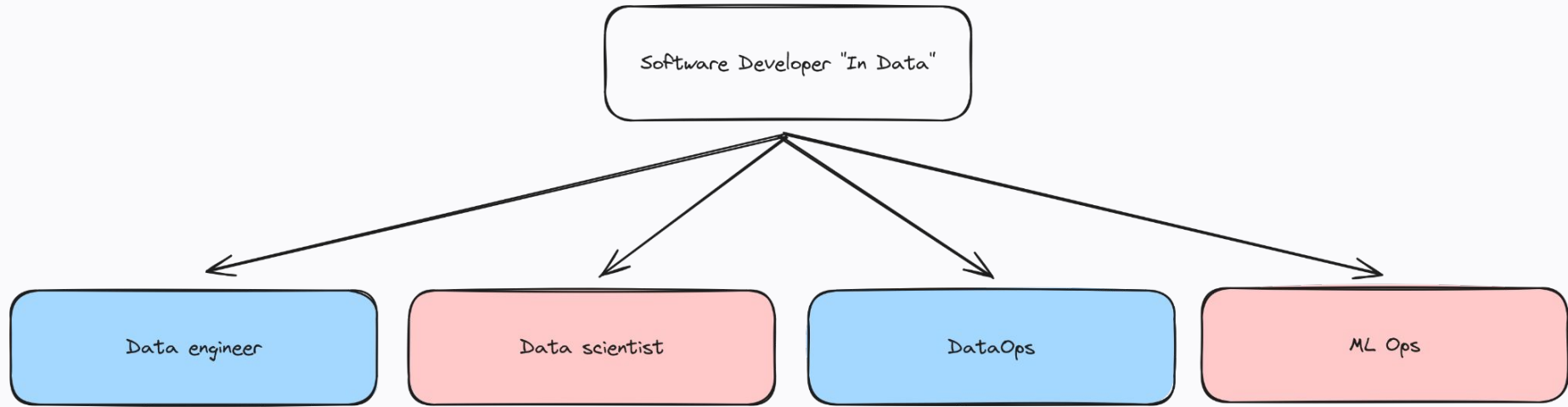Quality verification (e.g. Great Expectations, deequ)

Documentation (e.g. Amundsen)

Monitoring and alerting

you tried

# Evolution of a data-driven organization

# Evolution of a data-driven organization

Experimentation.

Long-range performance analytics.

Deep learning for understanding/prediction.
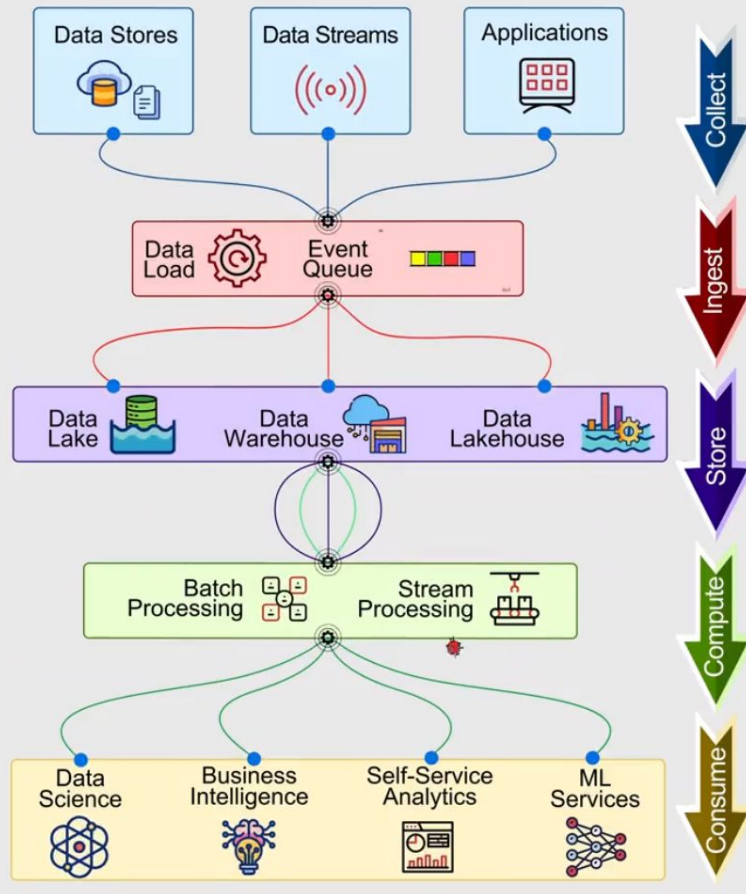
# Questions?



Jeffrey Leung

GitHub: **jleung51**

# Data Pipeline Overview

ByteByteGo

Data Stores · Data Streams · Applications

Collect

Data Load · Event Queue

Ingest

Data Lake · Data Warehouse · Data Lakehouse

Store

Batch Processing · Stream Processing

Compute

Data Science · Business Intelligence · Self-Service Analytics · ML Services

Consume

# The Dataflow

01. ingestion and storage

retention: **short**

02. cleaning and filtering

retention: **medium**

03. aggregating and analyzing

retention: **long**