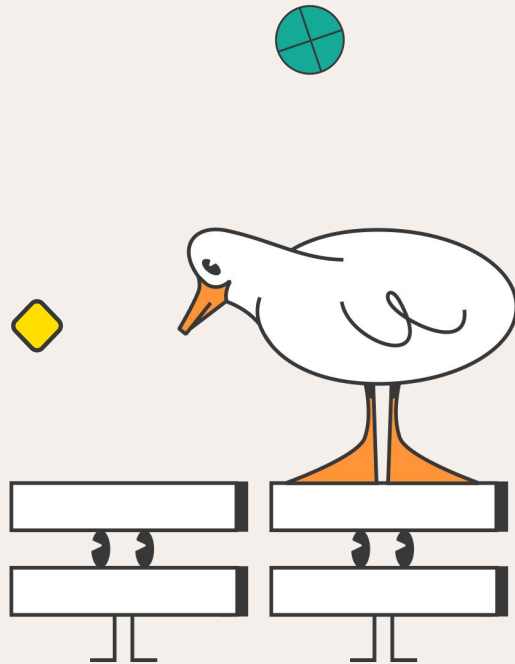# BIG DATA IS DEAD

## AND WHAT THE DUCK YOU CAN DO ABOUT IT

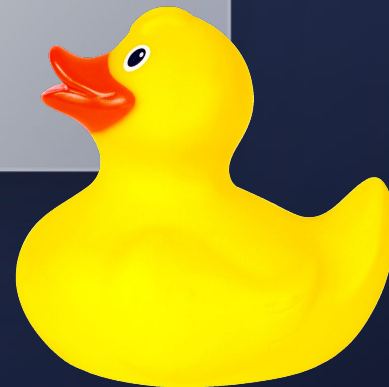JORDAN TIGANI

Chief Duck Herder @ MotherDuck

**Pete Soderling**
@petesoder

@thetinot I'll give @jrdntgn a keynote slot at Data Council if he wears the suit 🦆

1:49 PM · 28 Feb, 2023

2 replies    5 likes

# WHO AM I?

- BigQuery Engineer
- BigQuery Eng Director
- BigQuery PM Director
- MemSQL/SingleStore CPO
- MotherDuck CEO

- Lots of talking about Big Data

Photo above: Me warning about the perils of big data in 2012.

**MotherDuck**

**MotherDuck**

# ABOUT THIS TALK

## TELLING A STORY WITH GRAPHS

- Represent shapes rather than real data
- Based on experience with real data
- Only intended to be directionally accurate
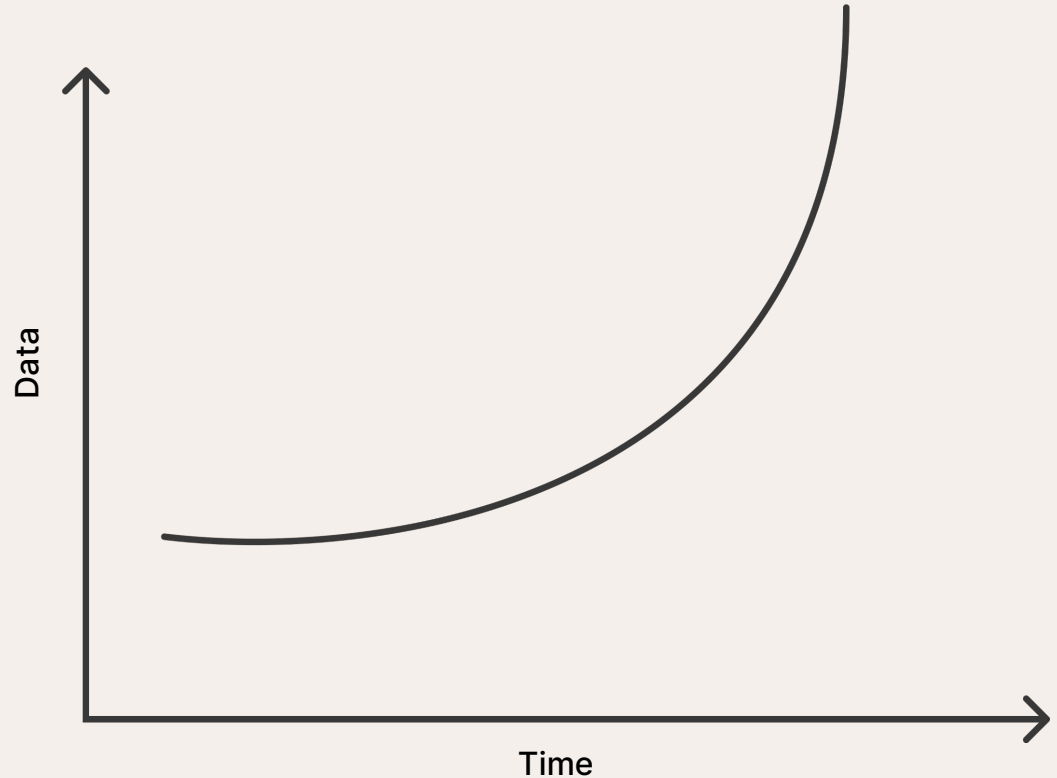
## WHERE THIS DATA COMES FROM

- Query Logs
- Deal Post-Mortems
- Benchmark results
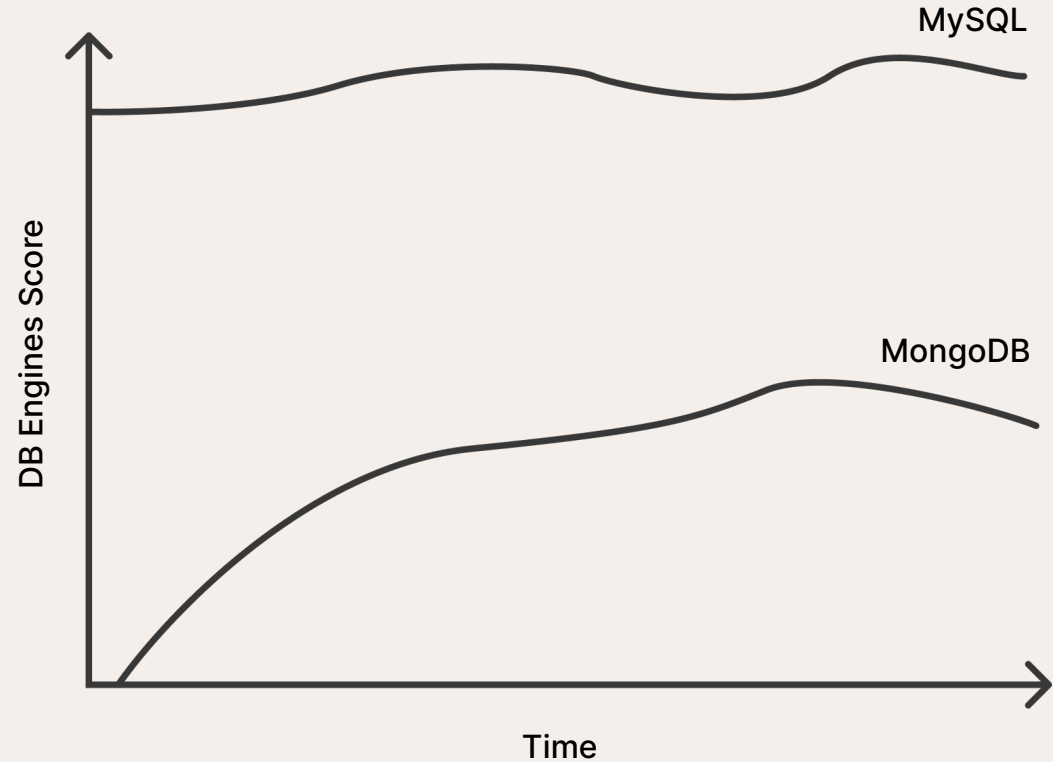- Bugs/Support Tickets
- Customer Conversations
- Service Logs

POPULARITY CONTEST:

BIG DATA IS NOT TAKING OVER

MotherDuck

DB Engines Score

MySQL

MongoDB

Time

# MOST PEOPLE DON'T HAVE THAT MUCH DATA

**MotherDuck**

Data Size (Log)

Customer Number

SEPARATION OF STORAGE AND COMPUTE FAVORS STORAGE

MotherDuck

Compute Needs

Data Size (Log)

WORKLOADS USE LESS DATA THAN YOU THINK

MotherDuck

Percentile

99

90

100MB

10GB

Query Workload Data Size

MOST OF
YOUR DATA
JUST SITS
THERE

MotherDuck

Bytes Processed

Data Age

# SCALE-UP IS NOT A DIRTY WORD

MotherDuck

445 core
24T RAM

**Memory
Optimized**

64 core
256G RAM

**Standard
Instance**

1 core
2G RAM

2006

2023

Time

DREMEL
IN A BOX

**Dremel: Interactive Analysis of Web-Scale Datasets**

Sergey Melnik, Andrey Gubarev, Jing Jing Long, Geoffrey Romer, Shiva Shivakumar, Matt Tolton, Theo Vassilakis
Google, Inc.

Dremel Performance:

87B Rows Scanned
512 GB processed
**3000 nodes**
~20 Seconds

MotherDuck

Cached In Memory: ✔
Cached On SSD: ✔
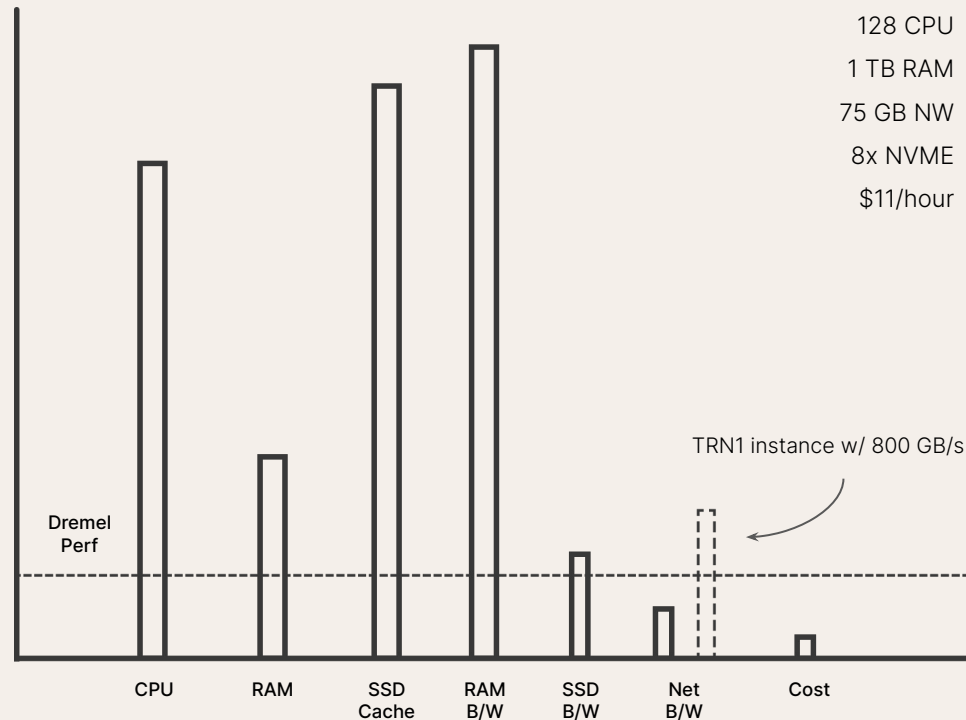Cold (S3): ✘   (✔?)

**AWS I4i.metal**
128 CPU
1 TB RAM
75 GB NW
8x NVME
$11/hour

Performance vs Benchmark

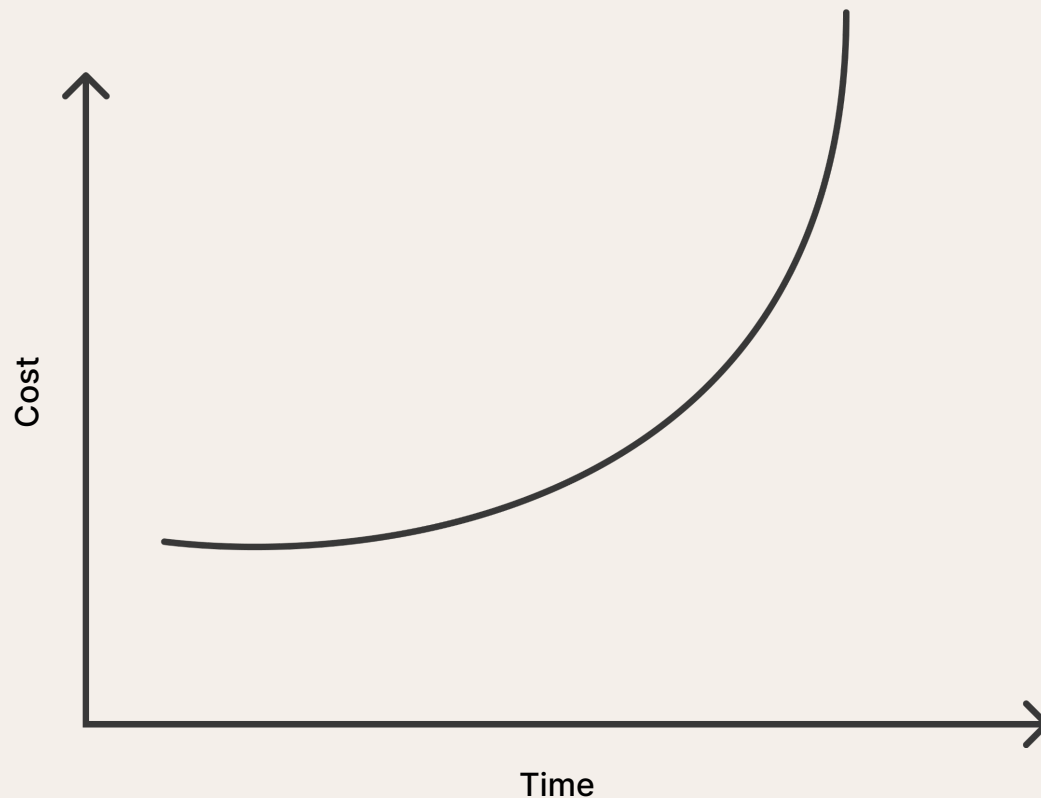Dremel Perf

TRN1 instance w/ 800 GB/s

CPU   RAM   SSD Cache   RAM B/W   SSD B/W   Net B/W   Cost

# THE BIG DATA FRONTIER IS HEADING OFF INTO THE SUNSET

MotherDuck

Percentile of Workloads

Single Node Vintage

2006

2023

# DATA AS A LIABILITY

WHAT IS THE
TRUE COST OF
YOUR DATA?

Cost

Time

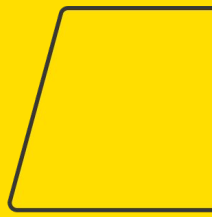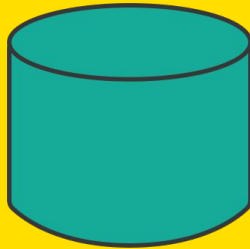# ARE YOU IN THE BIG DATA 1%

(and shouldn't we be building tools for the other 99%)
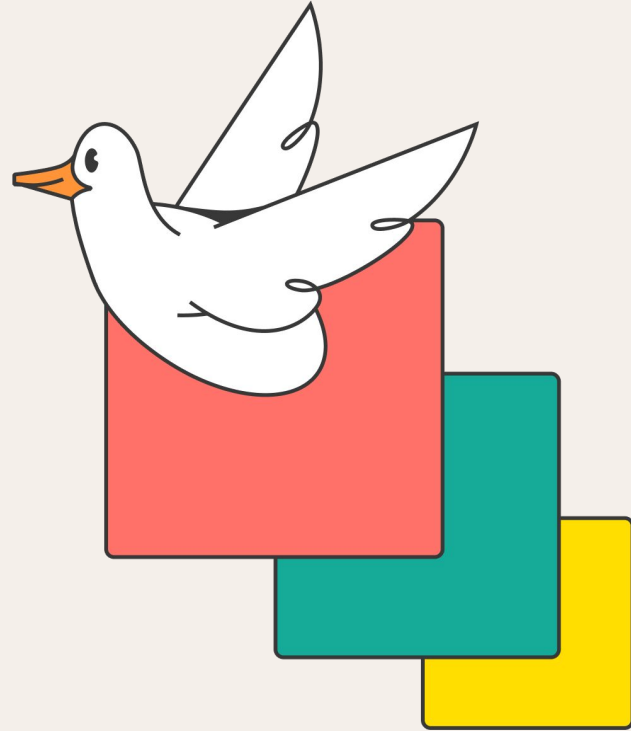
## YOU HAVE
## BIG DATA WHEN:

- You've really got a huge amount of data

- AND you need to a lot of that data at once

- AND what you're accessing won't fit on one machine

- AND you're not hoarding your data
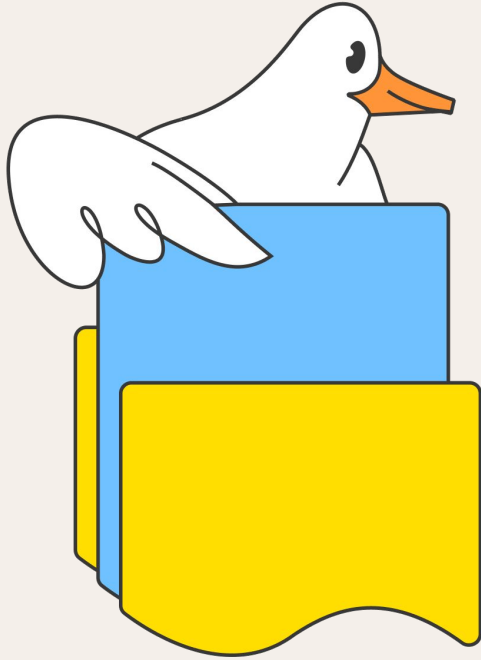
- AND you wouldn't be better off summarizing

MotherDuck

WHAT DOES THE WORLD LOOK
LIKE IF SIZE IS NO LONGER
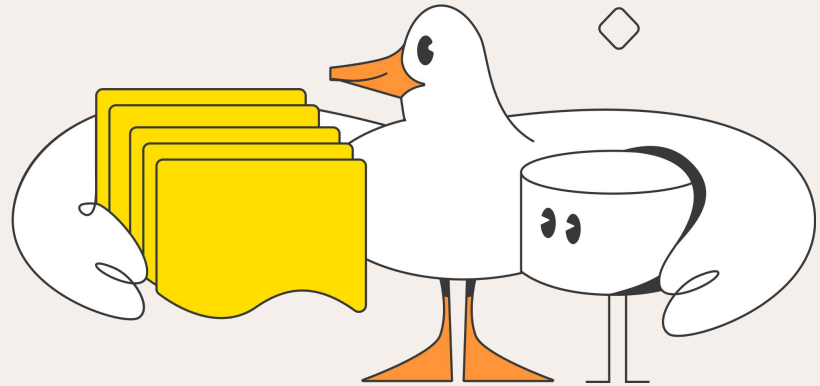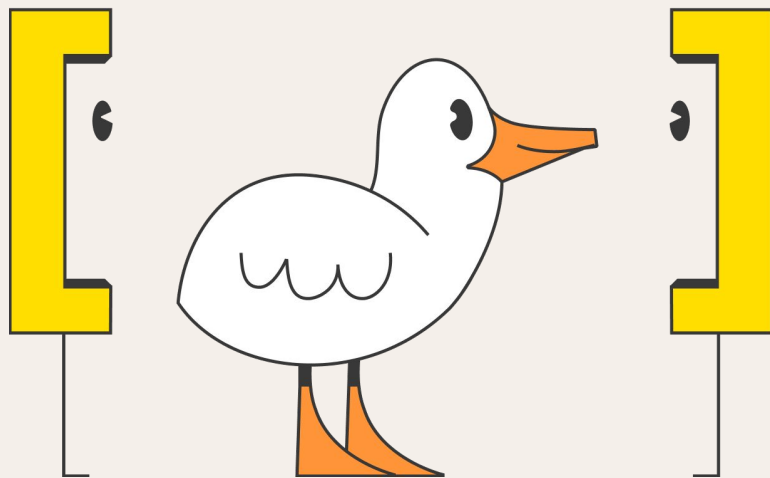THE PRIMARY DRIVER OF DATA
ARCHITECTURES?

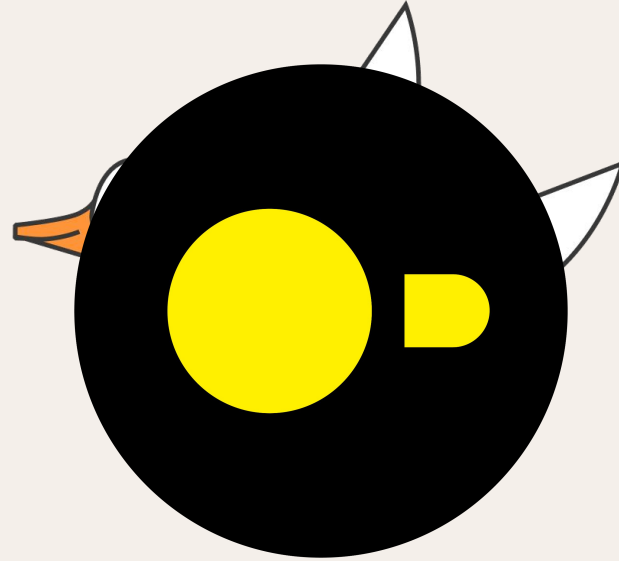DON'T BE AFRAID TO SCALE UP

MotherDuck

CLEAN UP
YOUR DATA

MotherDuck

BRING DATA
TO USERS

MotherDuck

LIFE IS
BETTER WITH
A DUCK

THE MODERN DUCK STACK HAPPY MEAL

MotherDuck

Ingest
+
DuckDB

Analyze
+
DuckDB

Visualize
+
DuckDB

# THANK YOU!

Don't be a data hoarder

**MotherDuck**