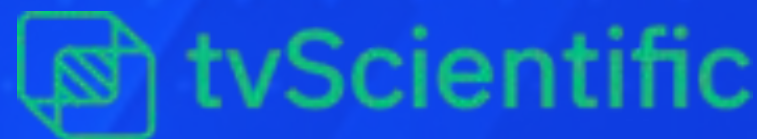# Rill is an operational BI platform

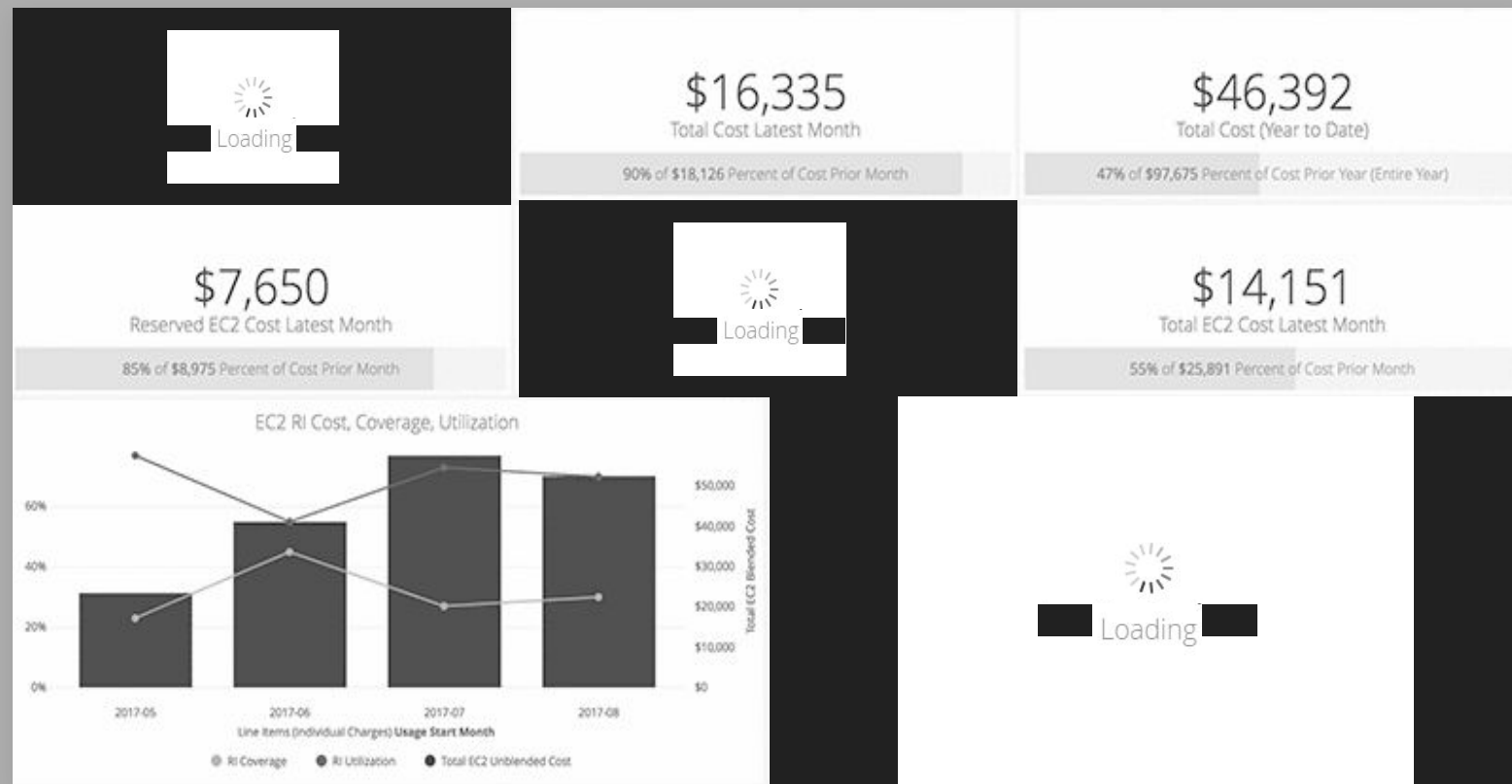**Proven tech**: Spun core IP out of Snap (ex-Metamarkets)

**Proven scale**: Managing 100TB+ of data, 1000s of users, millions of queries

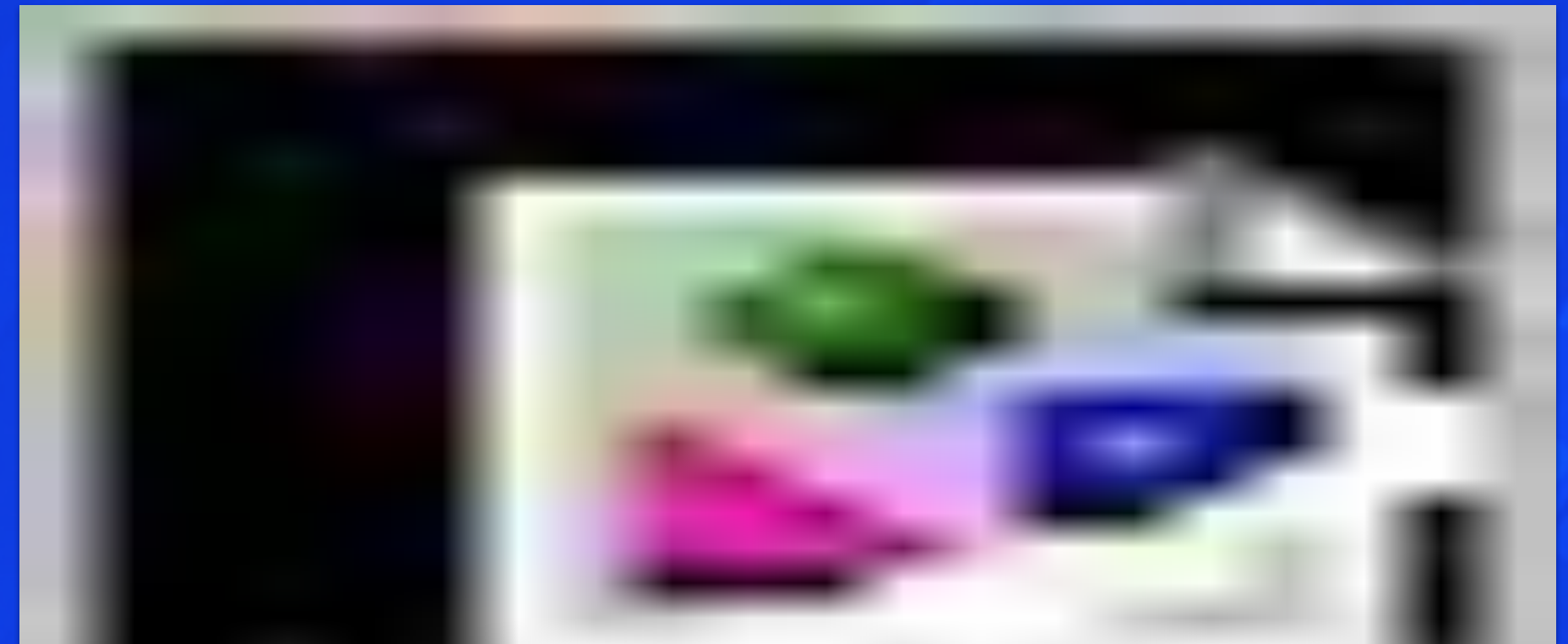**Proven impact**: Trusted by leading media, ecommerce, and tech platforms

# Traditional BI is
## *slow & rigid*



- Canned reports
- Limited exploration
- Works with slow data warehouses

# Rill is
## *fast & flexible*



- Ad hoc, exploratory analysis
- High dimensionality & cardinality
- <u>Requires a high performance DB</u>

# Why is Rill's BI tool different?

10X Faster Dashboards

- DuckDB & ClickHouse-powered, so queries return instantly

BI-as-Code

- Develop locally, deploy globally with benefits of Git workflows

Metrics-first Philosophy

- Users declare metrics with SQL expressions, Rill *auto-generates* dashboards

# 3-in-1 Architecture: ETL + OLAP + BI

# Our criteria for choosing OLAP engines

Speed

- Profiling 10s of GBs of datasets with sub-second performance

Simplicity

- Lightweight, embeddable, low maintenance overhead

Scalability

- Serve 100s of concurrent queries and scales to 100s of GBs

Open source

- Permissive license with a vibrant developer community

# Live Demo with the BlueSky Firehose



**INSTALL RILL**

```
$ curl https://rill.sh | sh
```

Visit www.rilldata.com to copy command

# Thank you!

✉ mike@rilldata.com

[www.rilldata.com](www.rilldata.com)

```
$ curl https://rill.sh | sh
```

## Follow us @RillData

# Appendix

# Speed at Scale Requires Coordinated Optimization Across the Stack

- Application layer
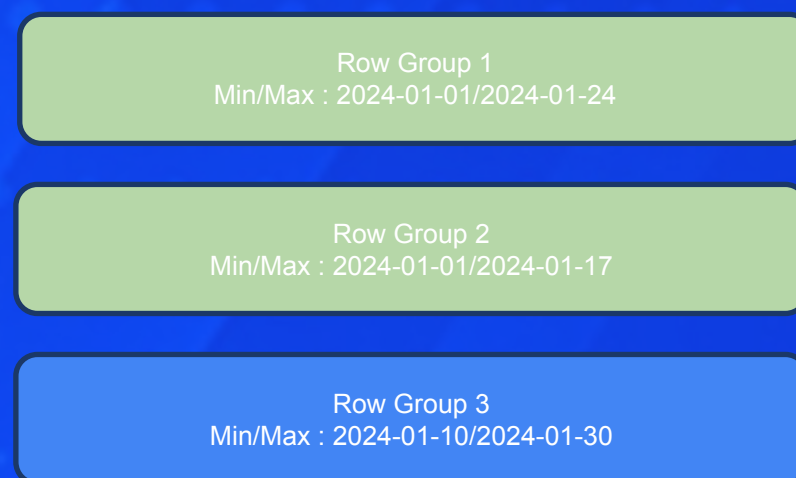- Data model
- Data engine (DuckDB)

# Data reordering speeds up min-max indexes

DuckDB's storage format stores the data in row groups, i.e., horizontal partitions of the data.

Sorting by time -> efficient use of min-max indexes

```
SELECT product_id, sum(total_sales) FROM sales
WHERE time BETWEEN '2024-01-01' AND '2024-01-07'
ORDER BY 1 DESC LIMIT 10
```

Without ordering

Row Group 1
Min/Max : 2024-01-01/2024-01-24

Row Group 2
Min/Max : 2024-01-01/2024-01-17

Row Group 3
Min/Max : 2024-01-10/2024-01-30

With ordering by time

Row Group 1
Min/Max : 2024-01-01/2024-01-10

Row Group 2
Min/Max : 2024-01-10/2024-01-20

Row Group 3
Min/Max : 2024-01-20/2024-01-30

Tip: For datasets already partitioned by time, preserving insertion order during ingestion is faster and leads to natural partitioning in row groups as well

# Enums are faster than strings

Converting strings to enum types leads to lower RAM usage, improving speed at scale

```
1   CREATE TYPE campaign_enum AS ENUM (SELECT DISTINCT campaign_name FROM events);
2   ALTER TABLE events ALTER COLUMN campaign_name SET TYPE campaign_enum;
```

Trade-offs

- Higher data ingestion time

- Incremental ingestion becomes harder as you need to rewrite column types for new values in enum

More details - https://duckdb.org/2021/11/26/duck-enum.html

# Query cancellation sheds unnecessary load

Exploratory dashboards can queue 100s of queries per view
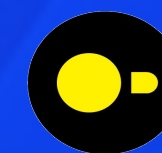
When the user changes a view, pending queries are no longer needed

Drill down

Cancel in-flight
API requests

Interrupt & cancel
query

# Priority queues keep applications snappy

Some application queries ought to be served faster than others:

- Interactive dashboard queries are the highest priority

- Scheduled reports and machine generated queries are lower priority

How

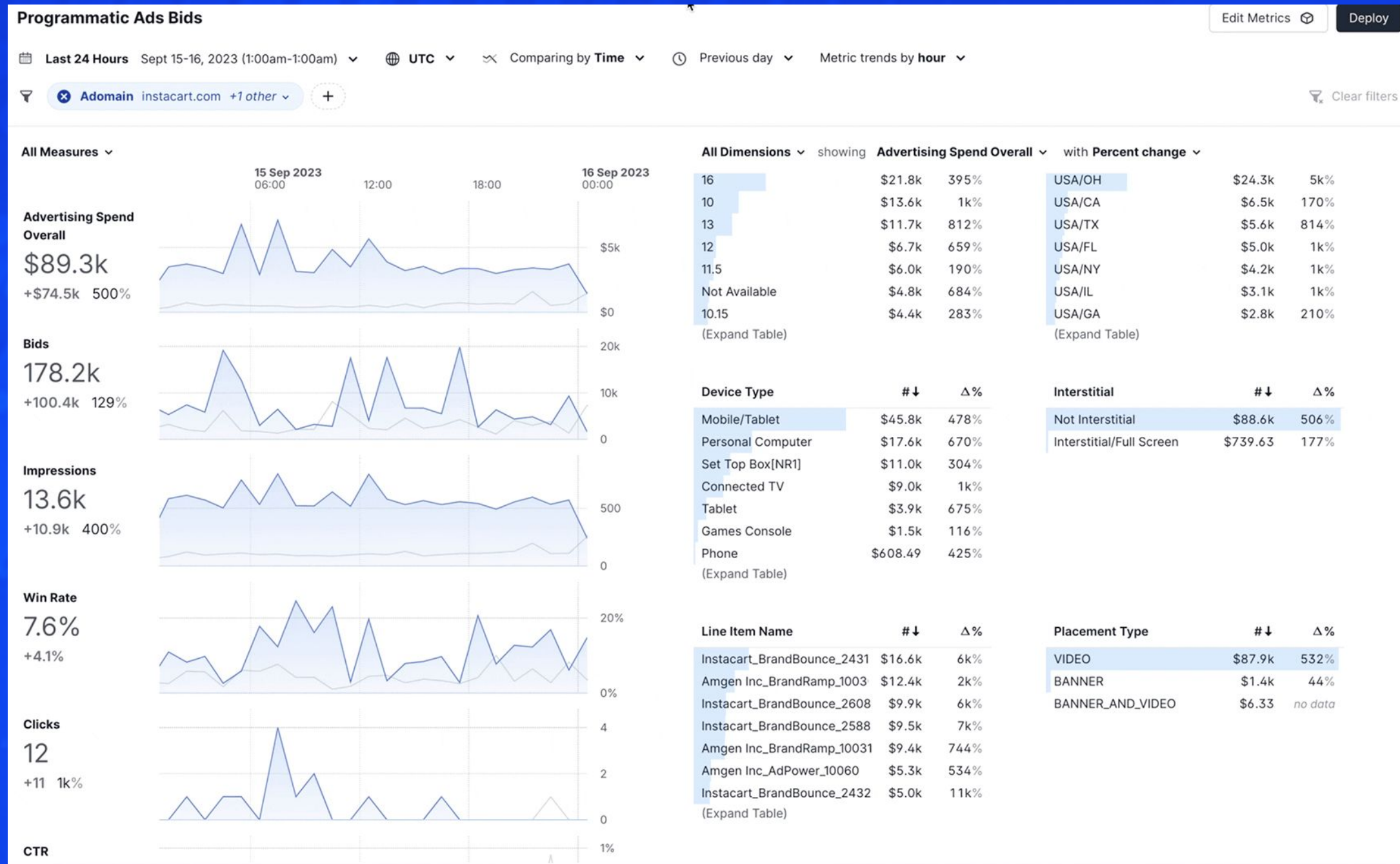- Create a priority queue for database queries

Results

 -  Improved dashboard interactivity

# WYSIWYF (What you see is what you fetch)
## Delayed evaluation of queries

# Data modeling: reduce data, retain insights

```
1   -- @materialize: true        -> Materialize model output
2   SELECT
3   DATE_TRUNC('DAY', event_datetime) AS event_date,    -> Aggregate to reduce data by 10X
4   store_id,
5   product_id,
6   SUM(sales_amount) as total_sales_amount,
7   SUM(quantity_sold) as total_quantity_sold
8   FROM sales
9   WHERE  event_date > current_date() - INTERVAL '1' YEAR    -> Prune per business needs
10  GROUP BY ALL
11  ORDER BY event_date ASC        -> Order by time for better use of min-max indexes
12  |
```