

Cost Containment

A Critical Piece of Data Team ROI



Lindsay Murphy

Head of Data



Instructor, Advanced dbt





New episodes every
Wednesday!



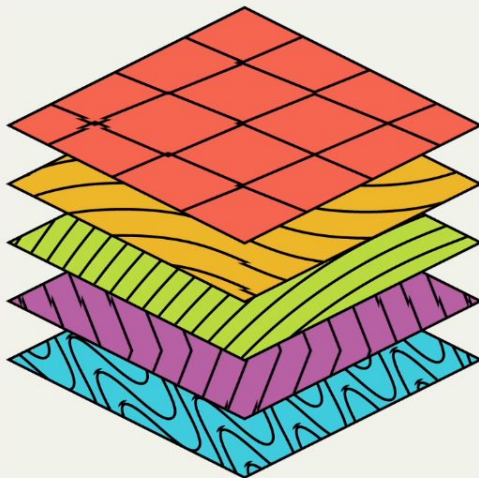
A VIRTUAL DATA CONFERENCE

PRESENTED BY  Secoda

MDS FEST 2.0

2024 LINEUP

OLEG AGAPOV • DYLAN ANDERSON • JAI BALANI
CHRISTOPHER BERRY • JULIE BEYNON
VIPUL BHARAT MARLECHA • CHRISTOPHE BLEFARI
MATTHEW BRANDT • TAYLOR BROWNLOW • KELLY BURDINE
STÉPHANE BURWASH • ERIC CALLAHAN
CYNTHIA CARIDAD • JES CARNEY • TIM CASTILLO
CHRISTOPHER CHIN • BRUCE CODELL • TIMO DECHAU
SHRAVAN DEOLALIKAR • RYAN DOLLEY • CECILIA DONES
TIANKAI FENG • KATIE HINDSON • KAREN HSIEH
MONICA KAY ROYAL • JERRIE KUMALAH
FAITH LIERHEIMER • MEGAN LIEU • ANDREW MADSON
PHOENIX MILLACY JAY • ETAI MIZRAHI
JORDAN MORROW • LINDSAY MURPHY
NATALIE NAKAMINE • PEDRAM NAVID • PARDIS NOORZAD
ZACK OBER • STEFANÍA ÓLAFSDÓTTIR • ABISOLA ONI
MEHDI OUAZZA • JAKE PETERSON • EMILY RIEDERER
BEN ROGOJAN • AUGUSTO ROSA • CHAD SANDERSON
ARCHIE SARRE WOOD • MADISON SCHOTT
EVA SCHREYER • SUMI SINGH • ABHI SIVASAILAM
PÁDRAIC SLATTERY • JEFF SLOAN • MARISA SMITH
BIJAN SOLTANI • GABI STEELE • JACK SWEENEY
MATT WEINGARTEN





- Catalog
- Lineage
- Governance
- Monitoring
- Documentation

The screenshot displays the Secoda web application interface. At the top, there are navigation tabs for Catalog, Lineage, Dictionary, Analysis, Requests, and Monitoring. The left sidebar contains a search bar and a navigation menu with options like Home, Inbox (23), AI Assistant, Teams, All teams, General, Product, Marketing, Data, Sales, and Finance. At the bottom of the sidebar are Integrations, Analytics, Settings, and the user profile for Olivia Morales.

The main content area shows a document titled "Marketing / Documents" with a "Net Promoter Score" report. The report includes a definition of NPS and a bar chart showing scores from Sep 16 to Sep 23. The current score on Sep 22 is 8.4.

Net Promoter Score

Net Promoter Score (NPS) is a measure used to gauge customer loyalty, satisfaction, and enthusiasm with a company that's calculated by asking customers one question: "On a scale from 0 to 10, how likely are you to recommend this product/company to a friend or colleague?"

Net Promoter Score
Sep 22, 2023

Net Promoter Score
8.4

Date	Net Promoter Score
Sep 16	4.5
Sep 17	7.0
Sep 18	8.5
Sep 19	7.0
Sep 20	5.0
Sep 21	6.5
Sep 22	8.4
Sep 23	8.5

Calculation

The NPS score is calculated using the following query (using a Fivetran

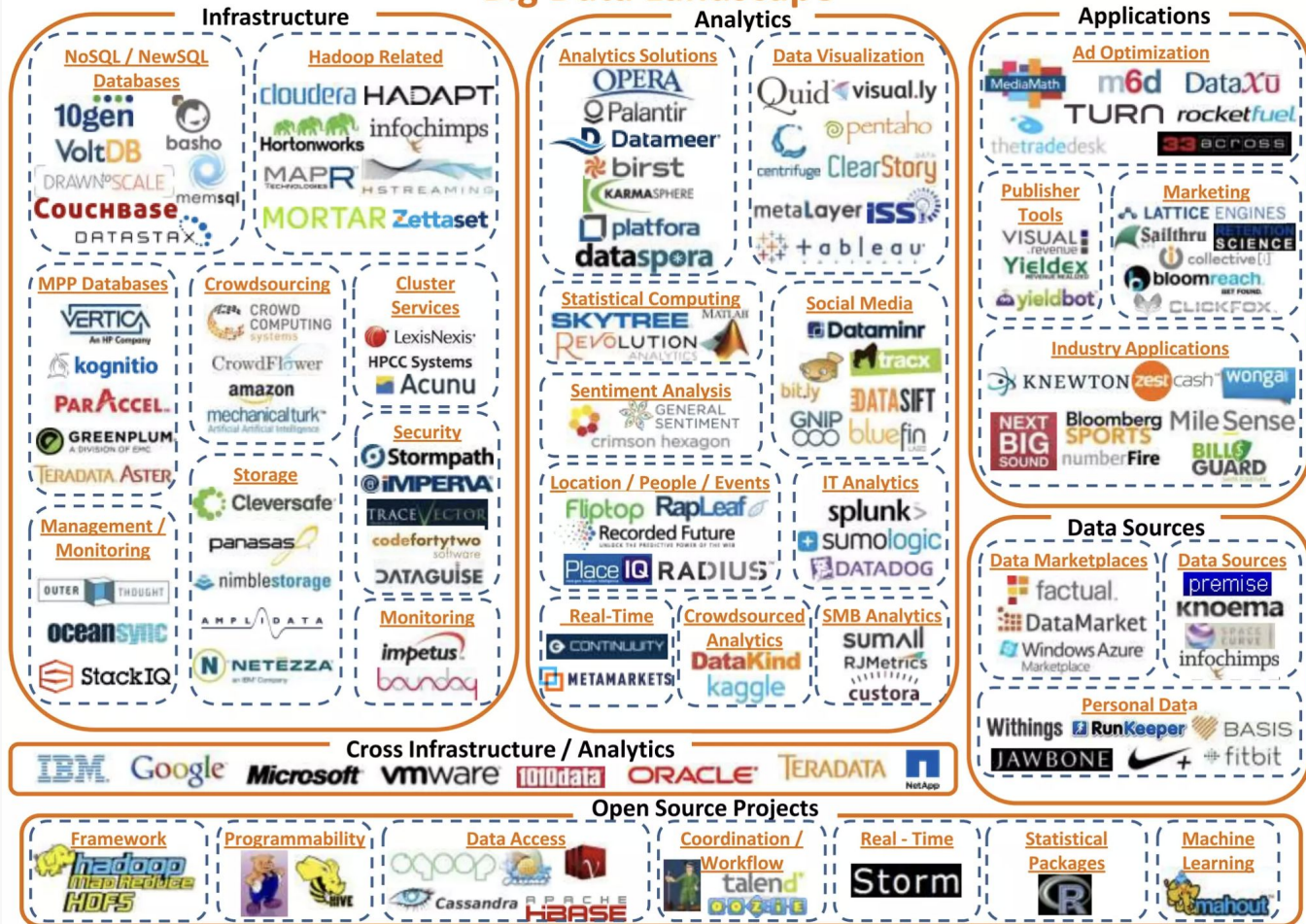
Let's dive in



**A lot has happened in our industry over
the last 10+ years**

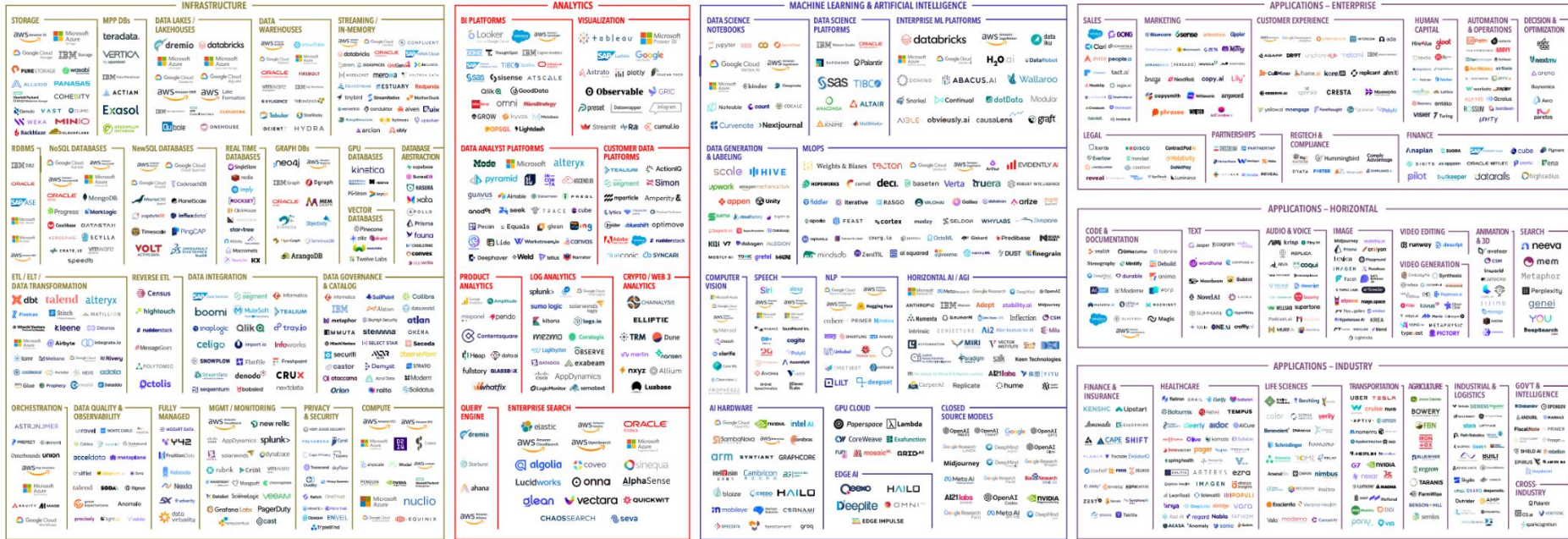
We got some new tools

Big Data Landscape



© Matt Turck (@mattturck) and ShivonZilis (@shivonz)

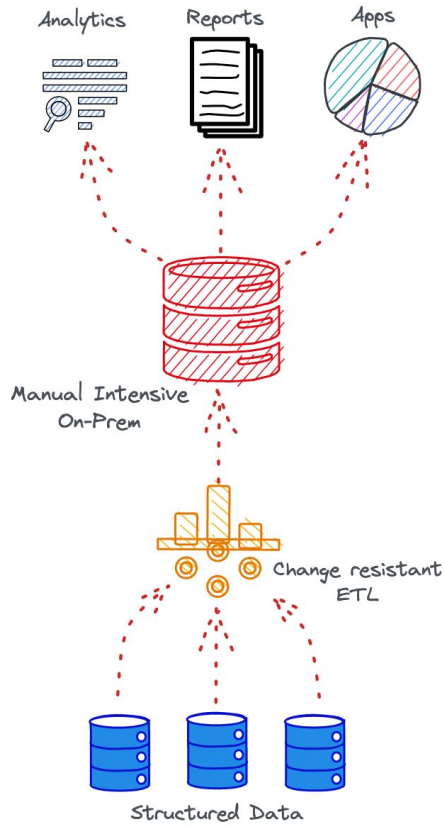
THE 2023 MAD (MACHINE LEARNING, ARTIFICIAL INTELLIGENCE & DATA) LANDSCAPE

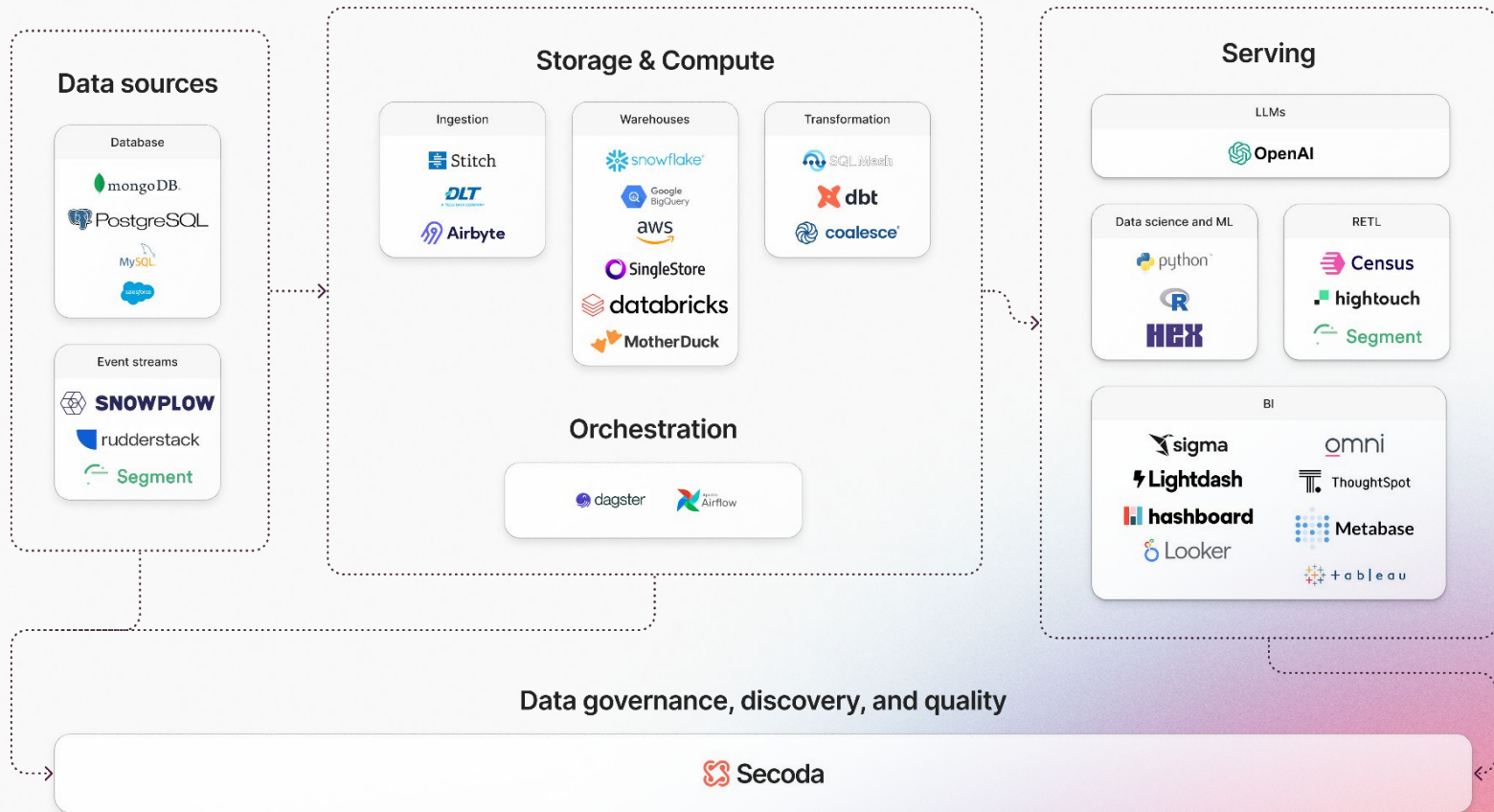


Source: <https://mattturck.com/mad2023/>

**It's gotten a lot easier to buy your entire
data stack**

Traditional Data Stack





Companies have invested a lot in becoming “data-driven”

Data Scientist: *The Sexiest Job of the 21st Century*

Meet the people who can coax treasure out of messy, unstructured data.
by Thomas H. Davenport and D.J. Patil

W

hen Jonathan Goldman arrived for work in June 2006 at LinkedIn, the business networking site, the place still felt like a start-up. The company had just under 8 million accounts, and the number was growing quickly as existing members invited their friends and colleagues to join. But when

Data Engineer: The Sexiest Job of the 21st Century

Why We Need Them Before Data Scientists

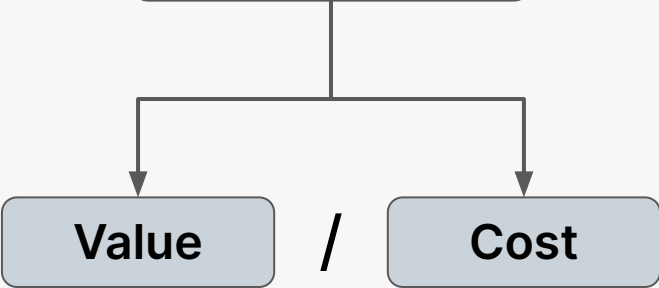
**Move Over Data Scientists,
Analytics Engineers Have the
Sexiest Job**



But some things haven't changed...

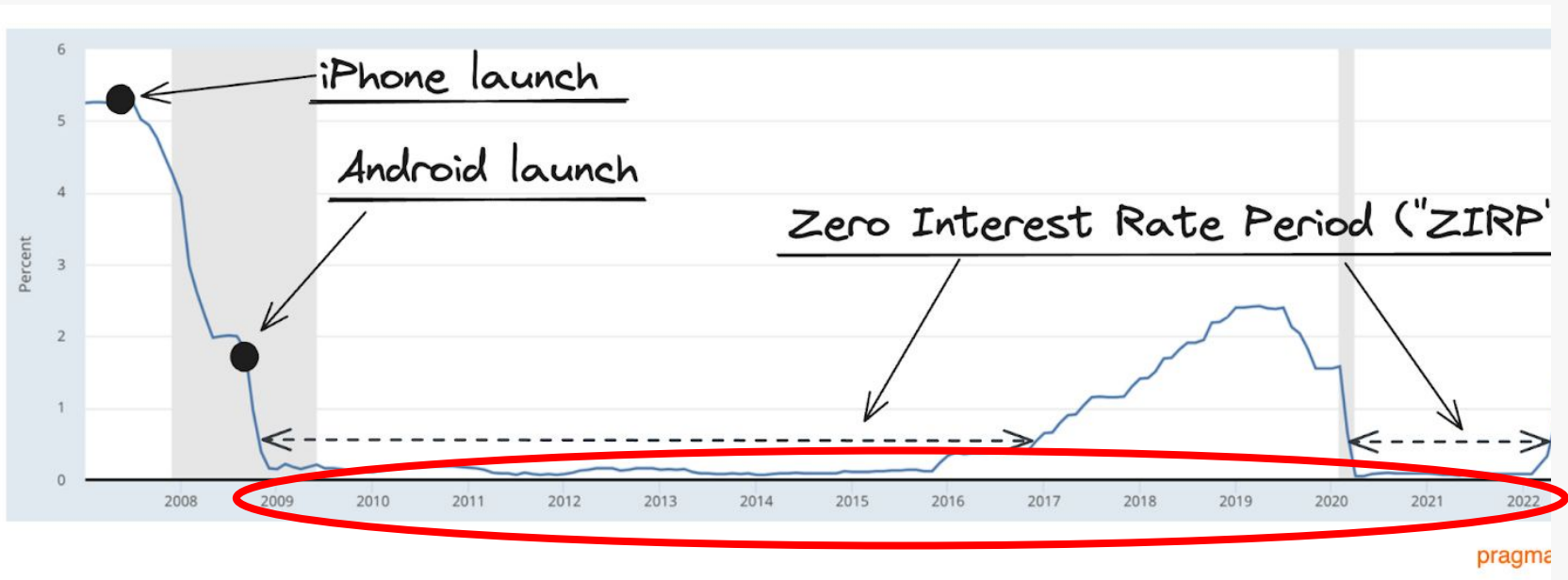
**Data teams still aren't very good at
measuring the value we deliver**

Data Team ROI



**Complicated to
measure**

**Pretty easy to
measure**





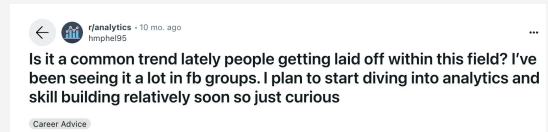
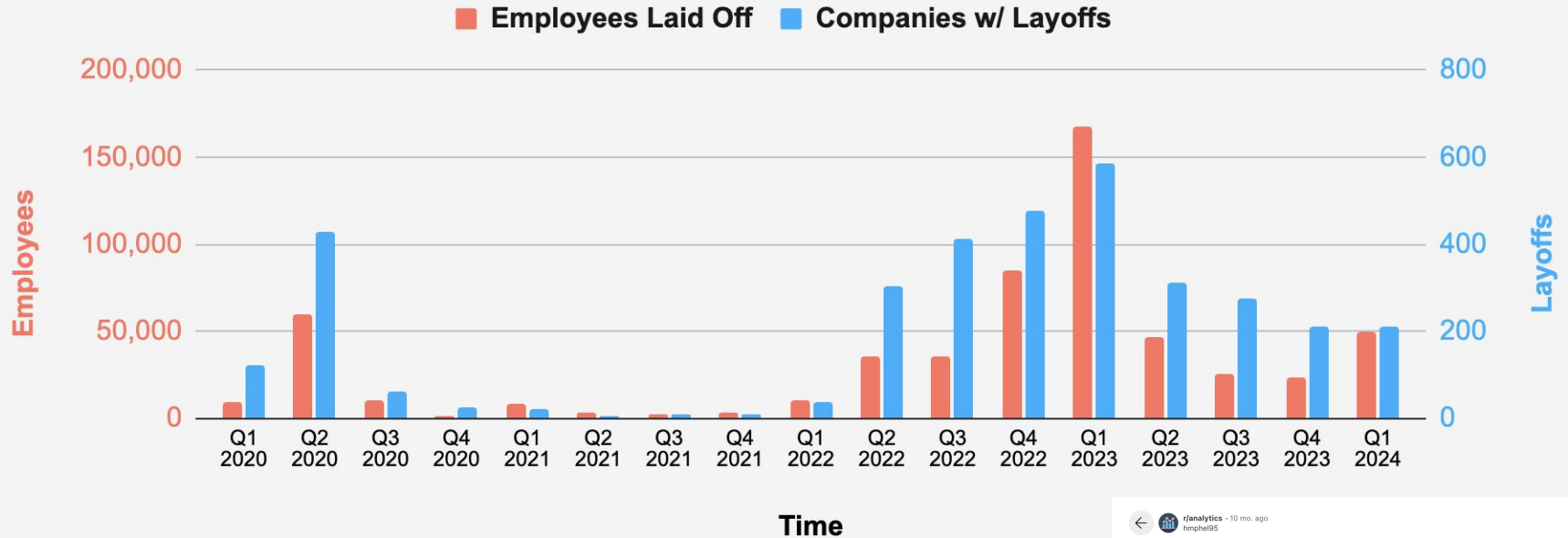


**Data teams get labelled as cost centres,
instead of value drivers**



Tech layoffs since COVID-19

Source: <https://layoffs.fyi>



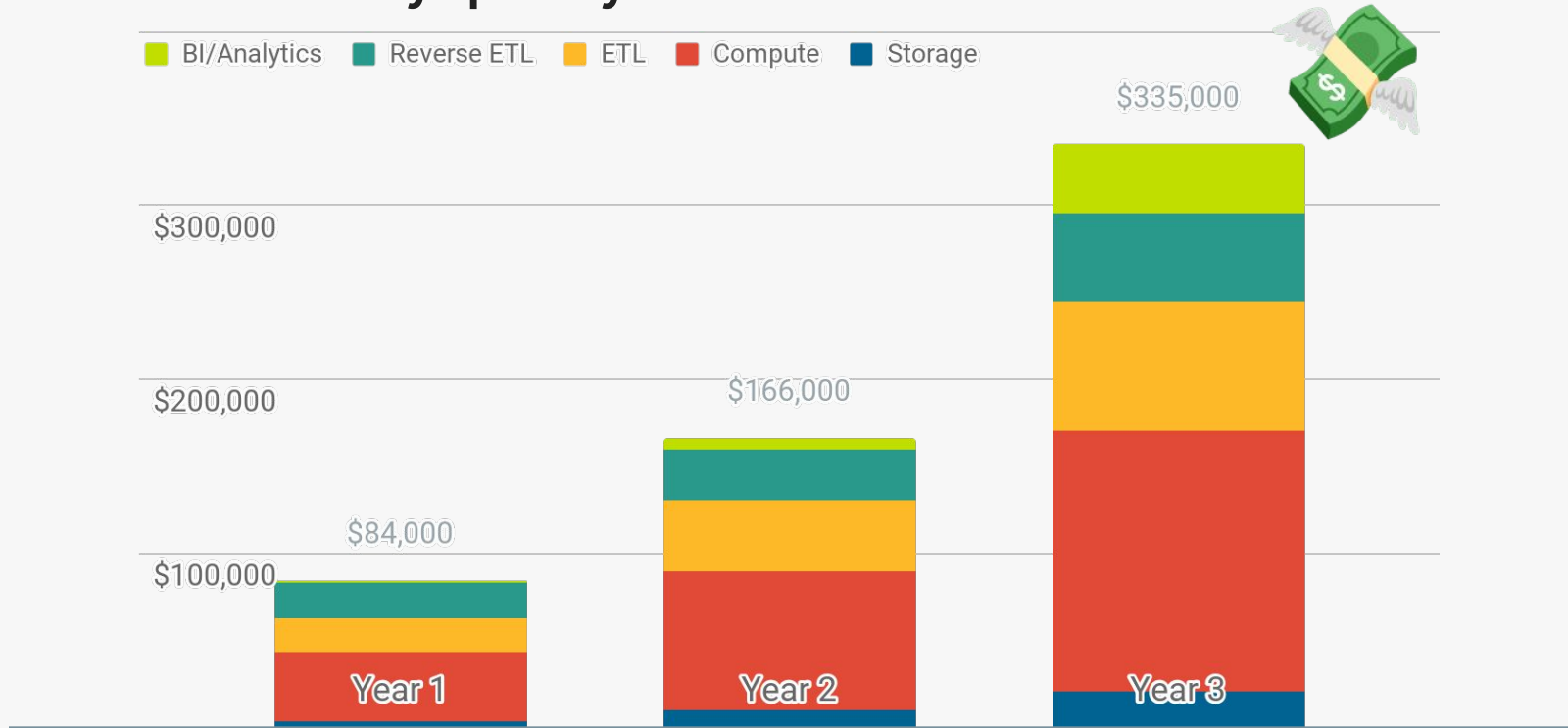


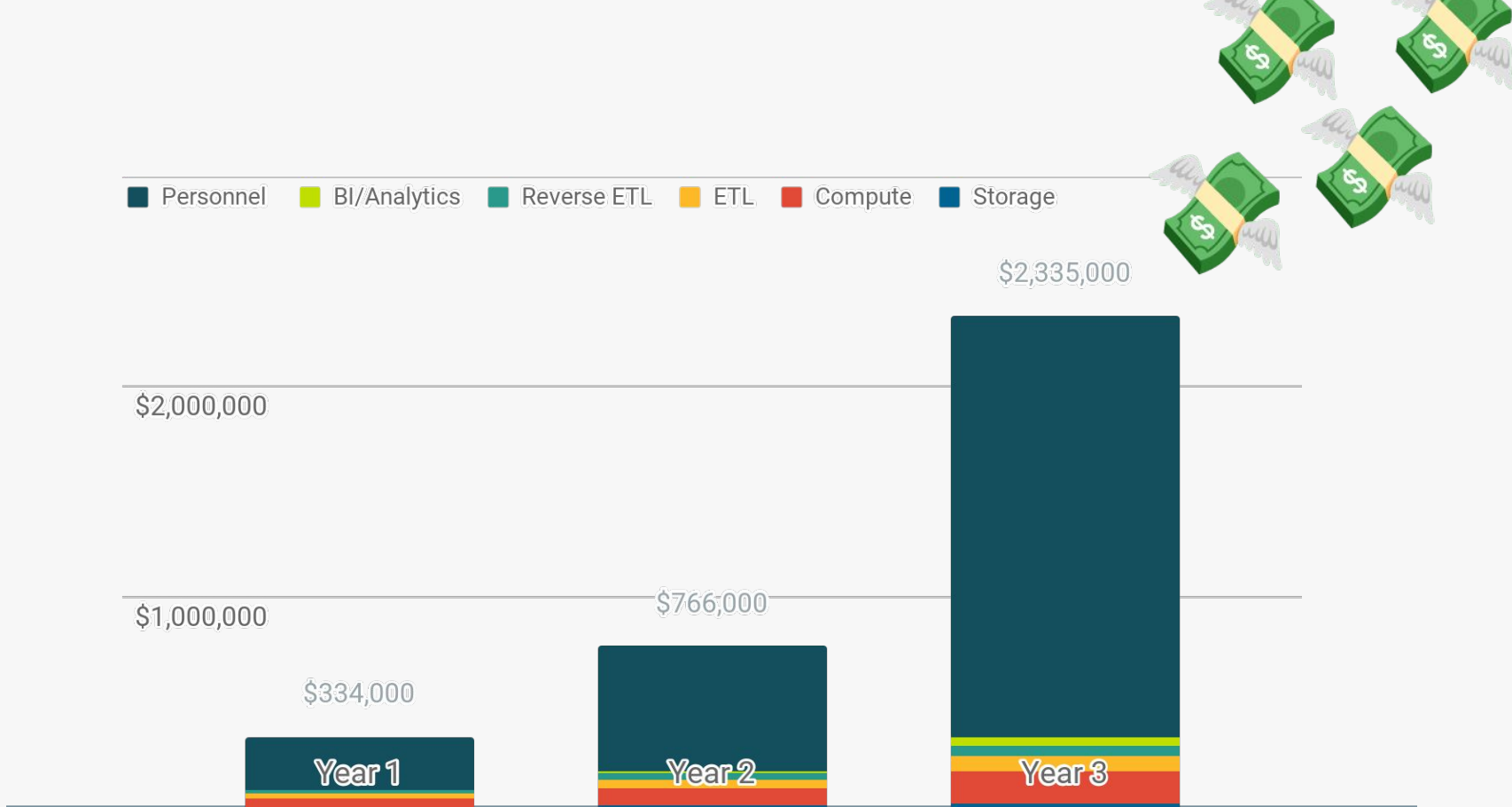
And we've got other problems too...

**The “modern data stack” is designed to
help you spend \$\$\$**

Initial costs are typically low...

But can scale very quickly...



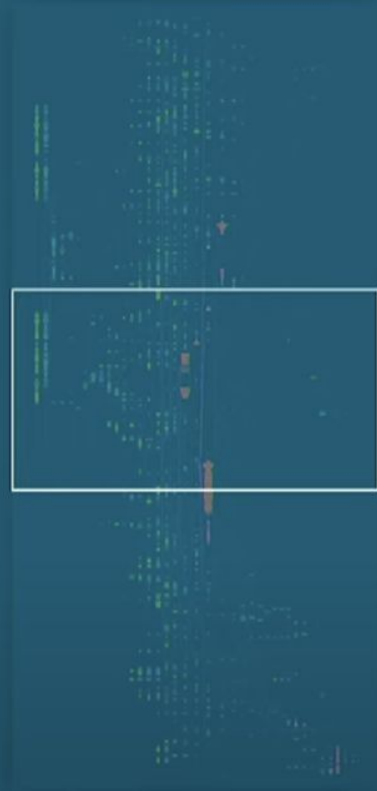


**Tools abstract away complexity in favour
of convenience**



**Addition is easy,
but depreciation can be...**

...a little scary

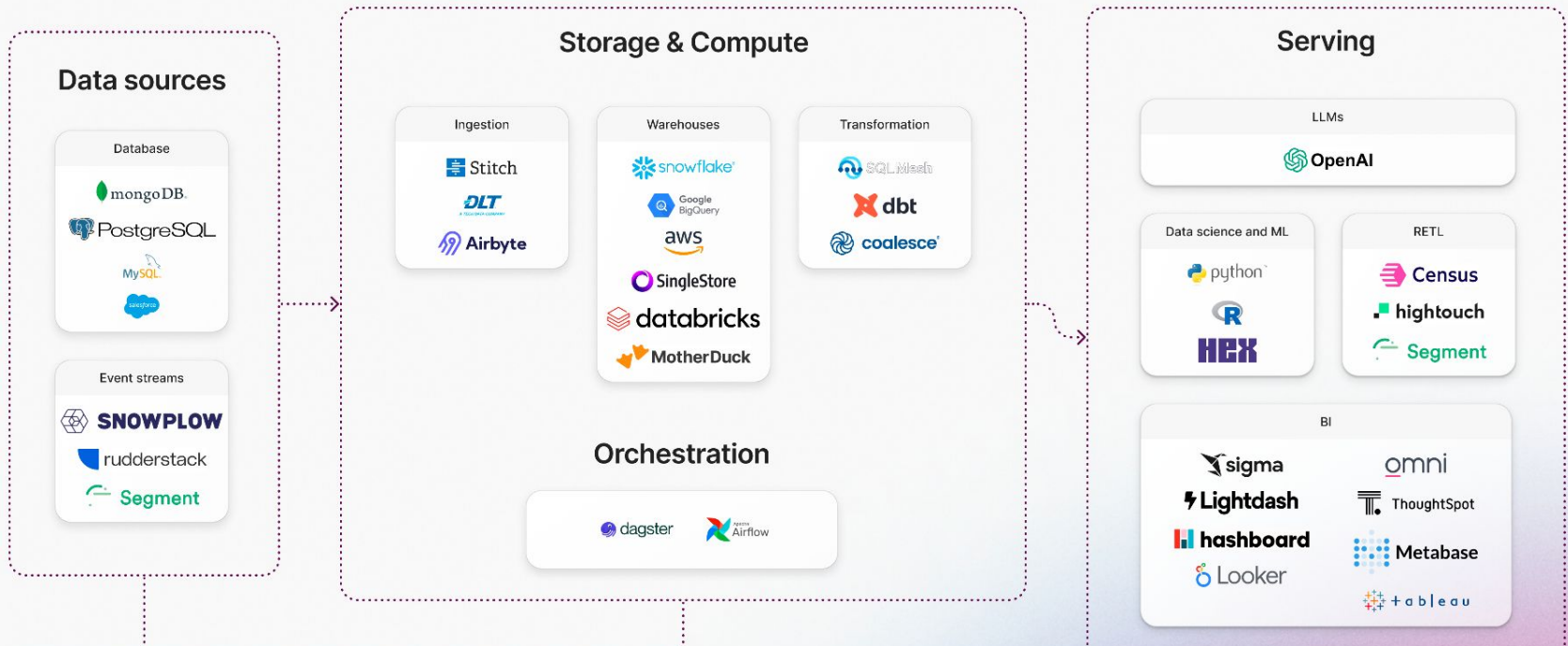


Source: dbt Labs product spotlight & keynote, Coalesce 2023

Interplay of pricing models

Consumption based pricing

Seat-based pricing



Snowball Effect



→ More data volume

→ More feature requests

→ More data models

→ More compute

**Cost and usage data is not easily
exposed or aggregated**



Billing & Usage

Billing Usage Usage Estimator

\$ **spend so far this month**

We've calculated this number based on your MAR consumption this month.

[See usage](#)

Alert me about my monthly spend



We will send you an email when your monthly spend reaches

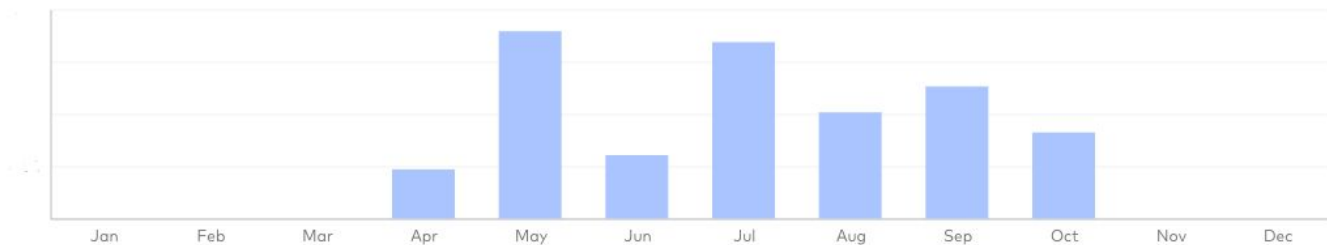
[Edit spend threshold](#) [See change history](#)

Monthly spend ⓘ

■ Spend

2023

All destinations





Billing & Usage

Billing **Usage** Usage Estimator

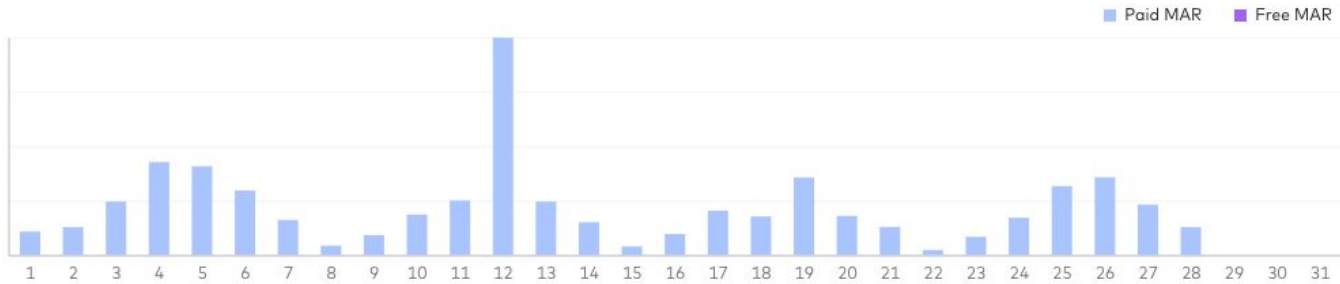
Connectors Transformations

Daily MAR ⓘ

All connectors ▼

All destinations ▼

Paid MAR × ▼





Time period

Last 30 Days

Source

All Sources

Destination

All Destinations

Total credits usage

Billed



Usage per connection

CONNECTION

SOURCE

DESTINATION

SCHEDULE

USAGE

US Read Replica → Snowflake



US Read Replica

CERTIFIED



Snowflake

CERTIFIED

6 hours



APAC Read Replica → Snowflake



APAC Read Replica

CERTIFIED



Snowflake

CERTIFIED

6 hours



EU Read Replica → Snowflake



EU Read Replica

CERTIFIED



Snowflake

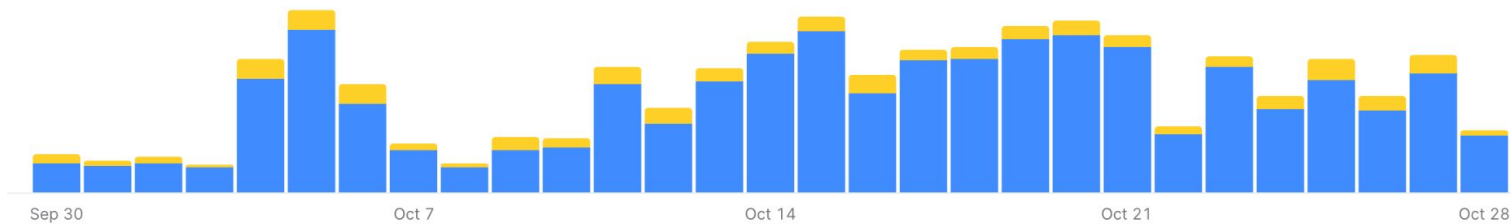
CERTIFIED

6 hours





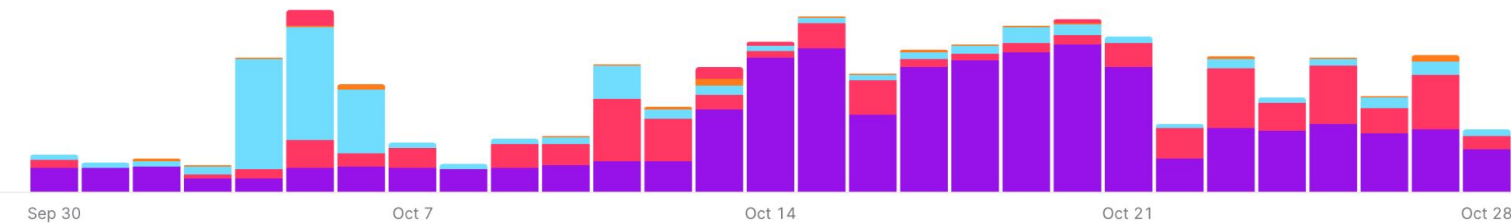
■ COMPUTE ■ CLOUD



■ CLOUD_SERVICES_ONLY

■ LOADING ■ MONITORING ■ REPORTING

■ TRANSFORMING





History ⓘ



Find a Field

All Fields

In Use

- ▶ Custom Fields + Add
- ▶ Dashboard
- ▶ Dashboard Creator
- ▶ Group
- ▶ History
- ▶ Look
- ▶ Merge Query
- ▶ Merge Query Source Query
- ▶ Model Set
- ▶ Node
- ▶ Permission Set
- ▶ Query
- ▶ Result Maker
- ▶ Source Query

Guided Analysis

A new experimental approach to analyzing your data



Which users are most active in your instance?

Understand who's getting the most value out of your...



What content is taxing your instance?

Identify heavily used content



User Audit

Understand user activity by type and role



Instance usage over time

Compare different activity metrics over time

Quick Start

Explore from a prebuilt analysis in History



Hourly Source Activity

What query sources are most active over the last...



Daily User Activity

Which users are most active over the last week?



Dashboard Creators

Which user's dashboards are most popular?



Historical Dashboard Usage

How many times has a particular dashboard been...



Query Panel

What are the most recently run queries?



Historical Look Usage

How many times has a particular look been used...



Average Runtime By Model

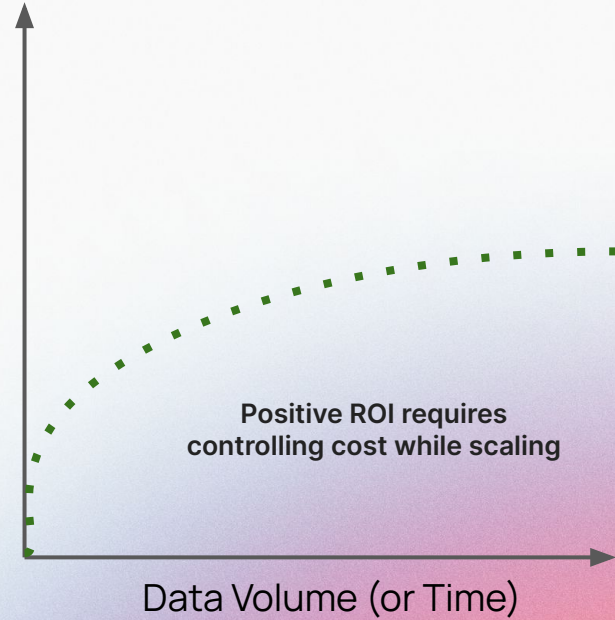
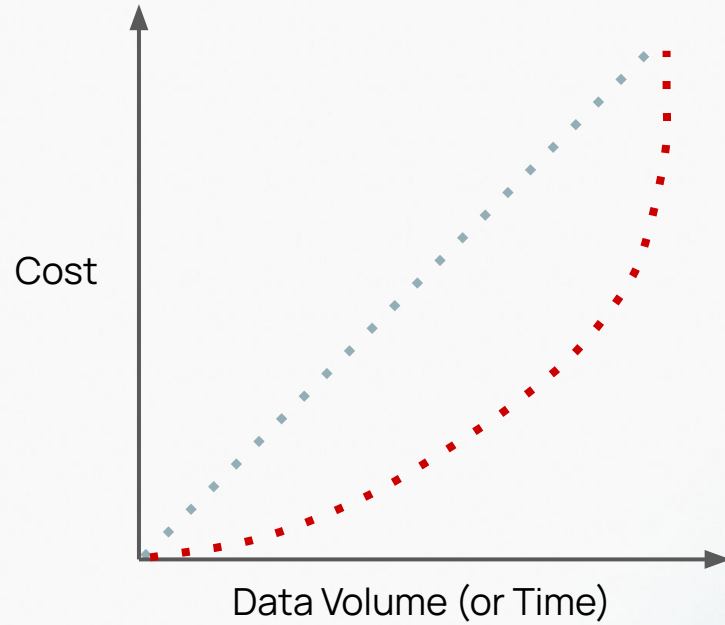
What are the average query runtimes of individual...



Recent User History

What is the recent query activity of a particular user?

Relationship between cost, volume or time



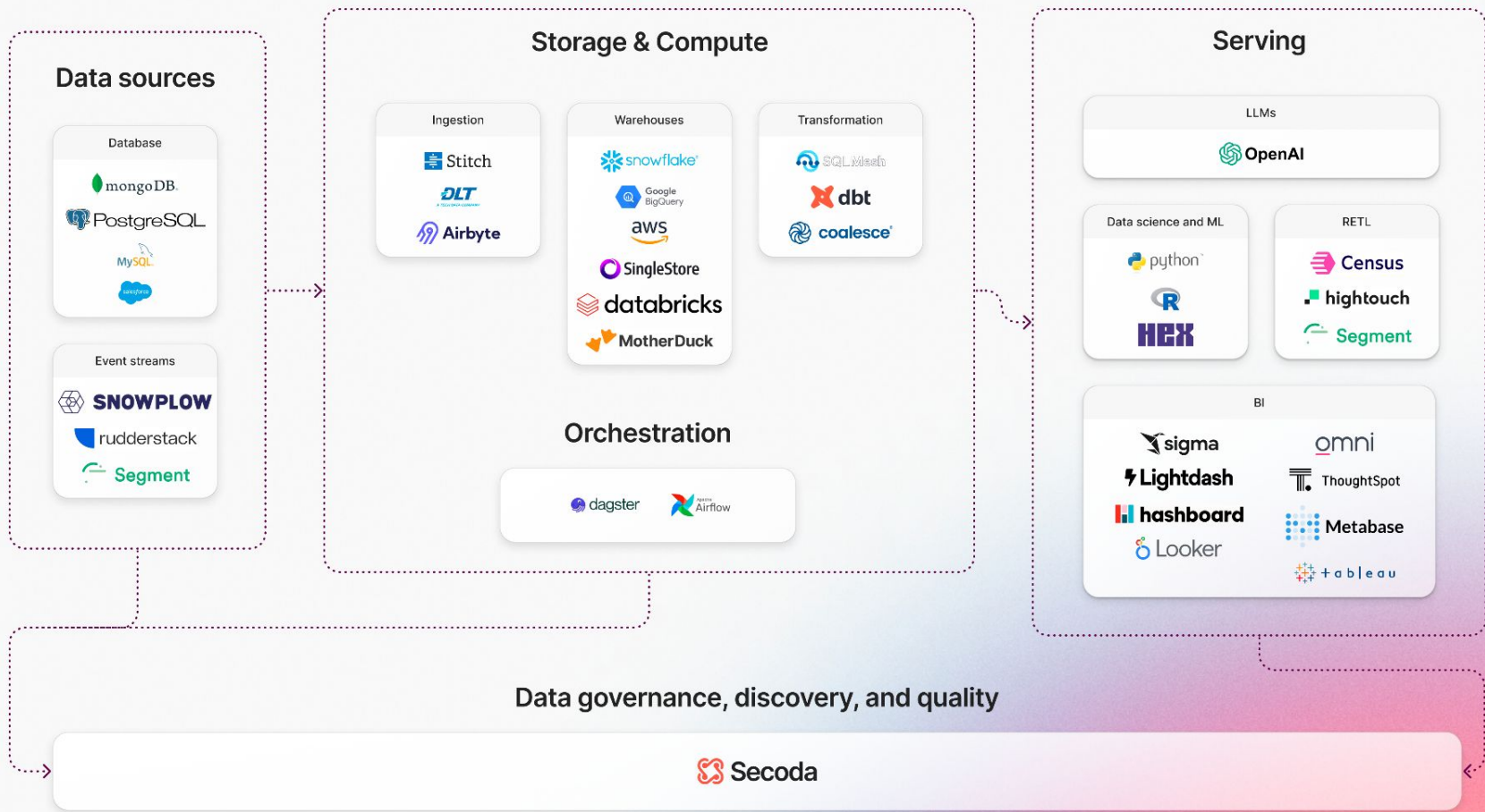


Recommendations



Measure

1. Map out your tech stack



2. Categorize pricing models and identify sources of cost related data

Tool	Type	Pricing Model	Pricing Model Category	Billing Cadence	Data Source
Fivetran	ETL	usage based	variable	monthly	Fivetran logs
Airbyte	ETL	open source	variable	monthly	AWS billing API
Snowflake	Warehouse	usage based	variable	monthly	Snowflake information schema
GitHub Actions	Orchestration	plan based	static	monthly	N/A
Lightdash	BI	open source	variable	monthly	AWS billing API
Retool	Data Applications	plan based	static (currently credit-based)	monthly	Float card
Segment	CDP	plan based	static (up to usage tier)	monthly	Float card
Mixpanel	Product analytics (quant)	annual license based	static	annually	Contract
Zapier	Automation	plan based	static (up to usage tier)	monthly	Float card

Static

Annual Contracts

Monthly Plans

Salaries

Variable

ELT

Warehouse

Transformation

Orchestration

Headcount

4. Source and model cost data

Sourcing variable cost data

Warehouses:

- <https://github.com/get-select/dbt-snowflake-monitoring>
- <https://github.com/kayrnt/dbt-bigquery-monitoring>
- <https://github.com/dbt-labs/redshift/>

ETL:

- https://github.com/fivetran/dbt_fivetran_log

Lineage Graph



dbt_snowflake_monitoring

Created by get-select

get-select/dbt-snowflake-monitoring 186



resources

packages

tags

--select

--exclude

All selected



dbt_snowflake_monito...

All selected



...

...

Update Graph



Lineage Graph



dbt_snowflake_monitoring

Created by get-select

get-select/dbt-snowflake-monitoring 186



SELECT



resources

All selected



packages

2 selected



tags

All selected



--select

+dbt_queries

--exclude

...

Update Graph



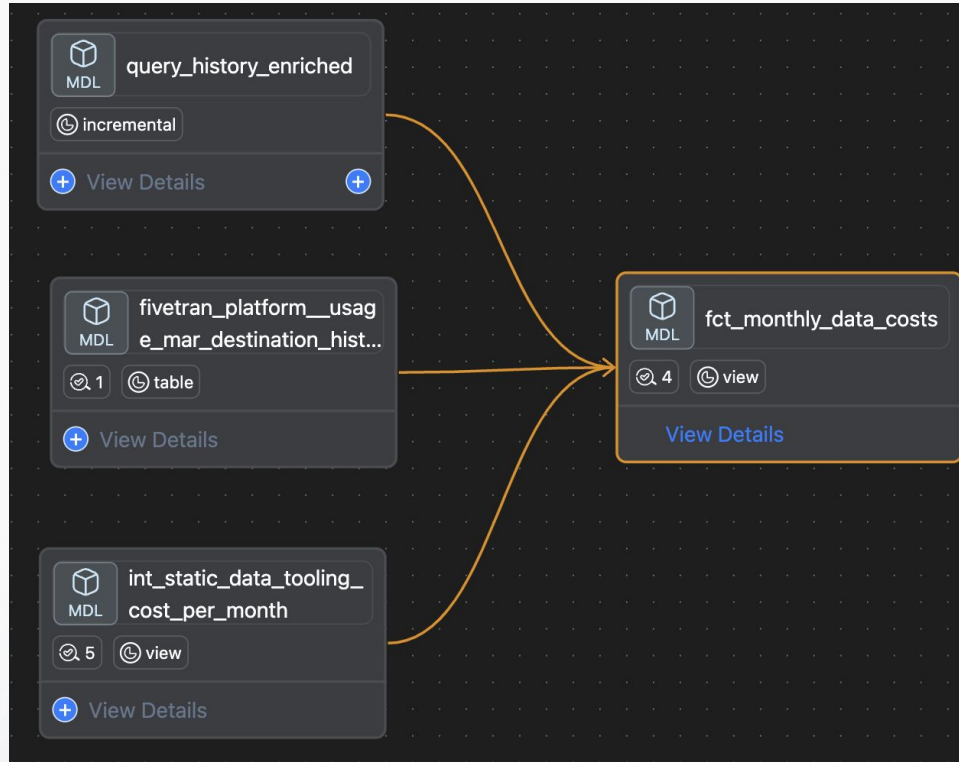
Building static costs from contracts

```
seeds > 📄 static_data_costs.csv
You, 1 second ago | 1 author (You)
1  tool,pricing_model,billing_cadence,billing_start_date,billing_end_date,cost,notes
2  segment,plan_based,monthly,02/19/2024,01/01/2030,###
3  mixpanel,plan_based,monthly,08/01/2023,07/31/2024,###
4  zapier,seat_based,monthly,01/01/2023,01/01/2025,###,seat based
5  fullstory,plan_based,annually,01/01/2023,01/01/2025,###,plan based
```

Building static costs from contracts

```
models > intermediate > ❌ lint_static_data_tooling_cost_per_month.sql
You, 20 hours ago | 1 author (You)
1 WITH
2
3 date_spine AS (
4   SELECT
5     calendar_date AS month_start_date
6   FROM
7     {{ ref ('dim_dates') }}
8   WHERE
9     day_of_month = 1
10    AND month_start_date <= CURRENT_DATE
11    AND month_start_date >= '2023-01-01'
12 ),
13
14 static_cost_tools AS (
15   SELECT
16     tool,
17     pricing_model,
18     billing_cadence,
19     DATE_TRUNC('month', billing_start_date) AS billing_month_start_date,
20     DATE_TRUNC('month', billing_end_date) AS billing_month_end_date,
21     cost
22   FROM {{ ref ('static_data_costs') }}
23 ),
24
25 final AS {{ You, 20 hours ago * add cost model and metrics
26   SELECT
27     date_spine.month_start_date,
28     static_cost_tools.tool AS spend_source,
29
30     CASE
31       WHEN static_cost_tools.billing_cadence = 'monthly' THEN static_cost_tools.cost
32       WHEN static_cost_tools.billing_cadence = 'annually' THEN static_cost_tools.cost / 12.0
33       ELSE NULL
34     END AS monthly_spend
35   FROM
36     date_spine
37   LEFT JOIN
38     static_cost_tools ON date_spine.month_start_date BETWEEN static_cost_tools.billing_month_start_date AND static_cost_tools.billi
39
40
```

Create a cost model



4. Build Operating Cost Assets



Data Team Operating Cost Dashboard



Add filter End time in the last 180 completed days Warehouse name is any value

Date Zoom

Dashboard Details

This dashboard contains spend information from Snowflake and specifically dbt queries.

Total Warehouse Spend

\$

Snowflake Annualized Total Spend

\$

+3% vs. Previous Day

Total Spend on dbt

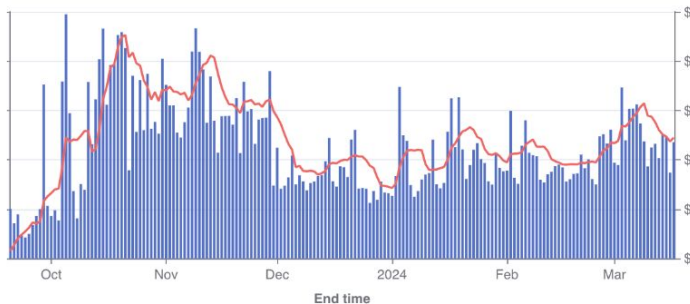
\$

Count of dbt Queries

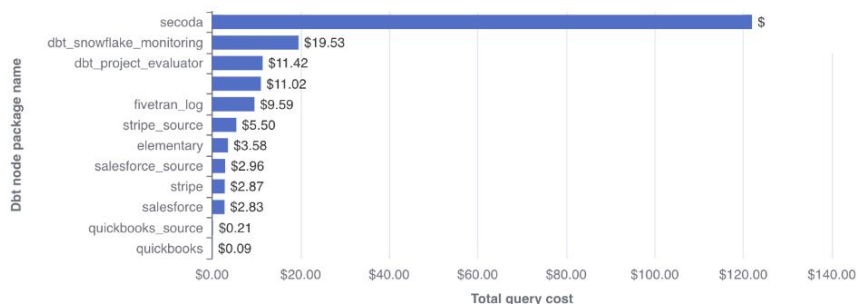
1,698,871

Snowflake: Daily and Annualized Query Cost

Total query cost Annualized Query Cost



Total Query Cost by dbt Package



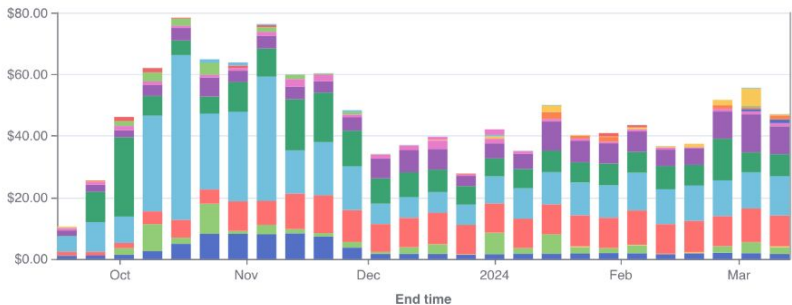


Data Team Operating Cost Dashboard



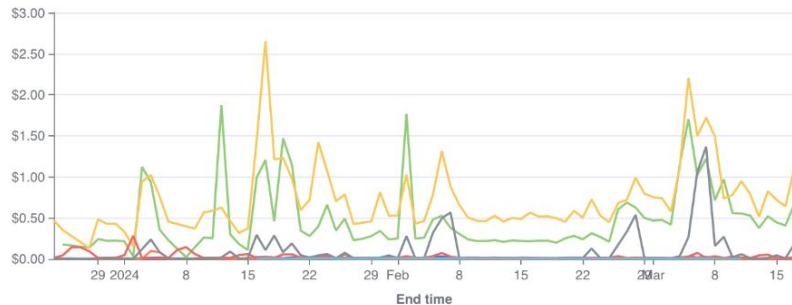
Daily Query Cost by User Name

SEGMENT_USER MODE_USER DBT_LINDSAY ZAPIER_USER FIVETRAN_USER AIRBY 1/5

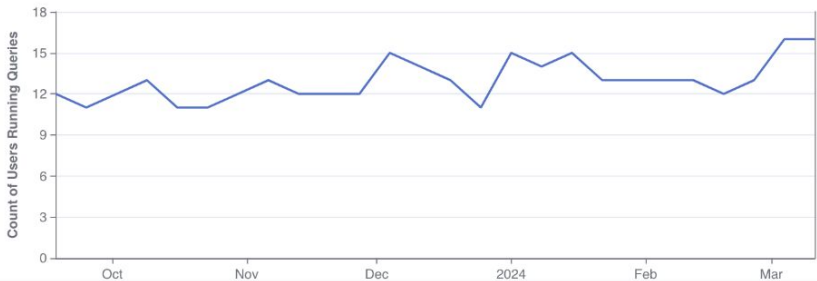


Daily dbt Cost by Resource Type

snapshot seed model test source ∅ operation



Unique Users Running Queries



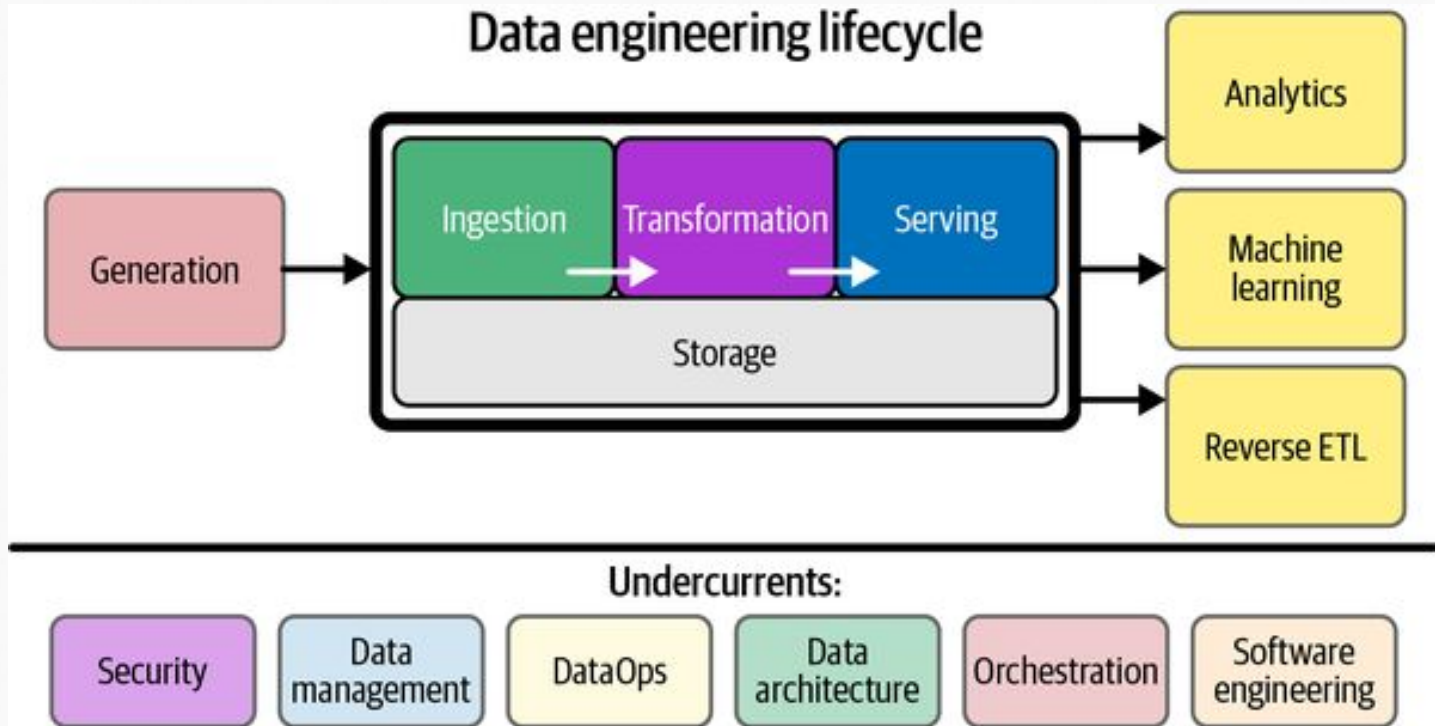
Most Expensive dbt Resources Last 30 Days



#	dbt Query Monitoring Dbt node name	dbt Query Monitoring Dbt node resource type	dbt Query Monitoring Total query cost	Annualized Query Cost
1	∅	∅		
2	stg_mixpanel_app_events	model		
3	hourly_spend	model		
4	query_history_enriched	model		
5	fct_marketing_conversions	model		
6	stg_secoda_admin_analyticmetric	model		

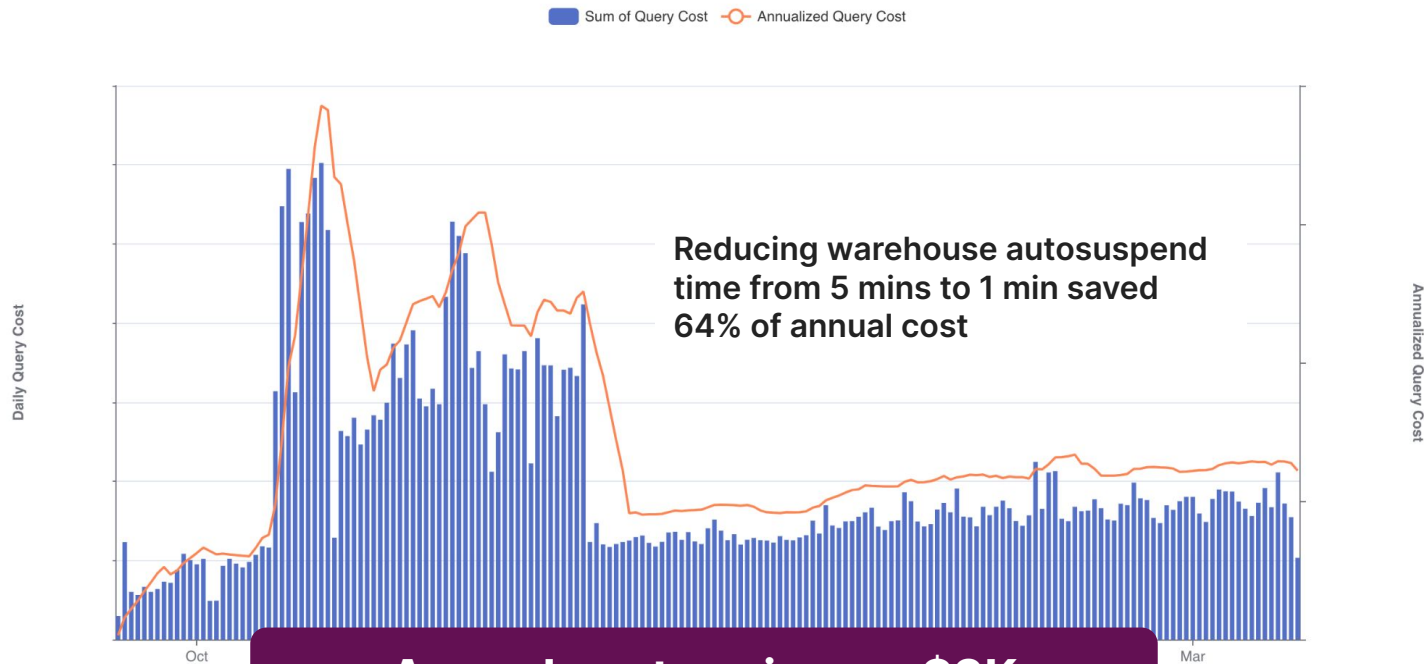
Optimize

Controlling Costs by Lifecycle Stage



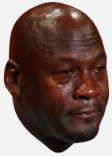
Snowflake Autosuspend Configuration

Daily Warehouse Cost and Annual Projection



Annual cost savings: ~\$3K

Open Source and Free Tiers: ETL



We got burned by
Fivetran one too
many times, switch
to Airbyte OS

Moved to Fivetran
Free Tier for SaaS
applications

Switching power
and familiarity with
multiple ETL tools

Annual cost savings: ~\$7.5K

Consider Total Cost of Ownership



Free

Start sending customer data to your favorite marketing and analytics tools.

Create a free account

\$0/month

What you can do:

- ✓ Includes 1,000 visitors/mo
- ✓ 500,000/mo Reverse ETL Records
- ✓ 2 sources
- ✓ 450+ Integrations
- ✓ 1 data warehouse destination

Team

Unify your customer data and collect every touch point with full access to Connections.

Try for free

Starts at \$120/month

All of Free plus:

- ✓ Includes 10,000 visitors/mo ⓘ
- ✓ 1,000,000 Reverse ETL Records
- ✓ Unlimited Sources
- ✓ Public API Access

Business

Solve complex business problems with standardized data and deploy the industry's #1 CDP for market share per IDC, 2022.

Get a demo

Custom Pricing

All of Team plus:

- ✓ Custom Volume
- ✓ Single View of the Customer
- ✓ Data Governance
- ✓ Advanced Roles & Permissions
- ✓ Personalized Customer Experiences
- ✓ HIPAA-eligibility (BAA Required)
- ✓ Regional Segment (EU or US) ⓘ

Consider Total Cost of Ownership

Community Edition



SNOWFLOW

HOSTED AND MANAGED BY YOU

A do-it-yourself solution for early-stage prototypes

- ✓ Manage your own infrastructure, including scaling, upgrades, failovers and costs
- ✓ Self-hosted QA and Dev pipelines
- ✓ Non-production workloads

**Human Behaviour:
Processes, Workflows, and
Accountability**

Processes

- “Use only what you need”
- Business need driven SLAs
- Data exhaust vs. data creation
- Asset deprecation process

Workflows

- Query optimization training
- Training for data producers and consumers
- Build cost optimization into your workflows
- Set cost safeguards where possible (e.g. resource monitors)
- Build cost monitoring dashboard and set up alerting
- Code reviews and feedback loops

Filters

All of the following dimension conditions match:

Dbt node name



is

stg_mixpanel_app_events



Dbt target schema



includes

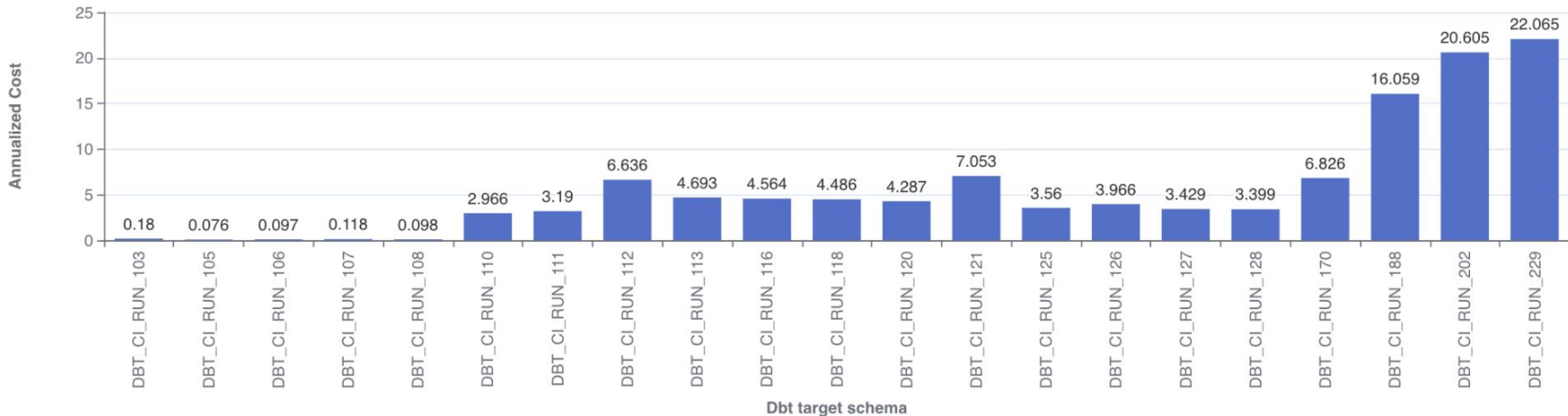
dbt_ci_run_



+ Add filter

Chart

Close configure



Accountability

- Cost of feature delivery
- Share accountability
- Normalize cost discussions
- Educate on cost reduction

User name	Total query cost
AIRBYTE_USER	\$:
FIVETRAN_USER	\$
GITHUB_ACTIONS_DBT	\$
SEGMENT_USER	
DBT_LINDSAY	
RETOOL_USER	
LINDSAY	
DBT_CAMERON	
DBT_LIGHTDASH_USER	
ANDREW	

The background features a gradient of purple and pink colors, with several large, overlapping, curved shapes that create a sense of depth and movement. The word "Scale" is centered in the middle of the image.

Scale

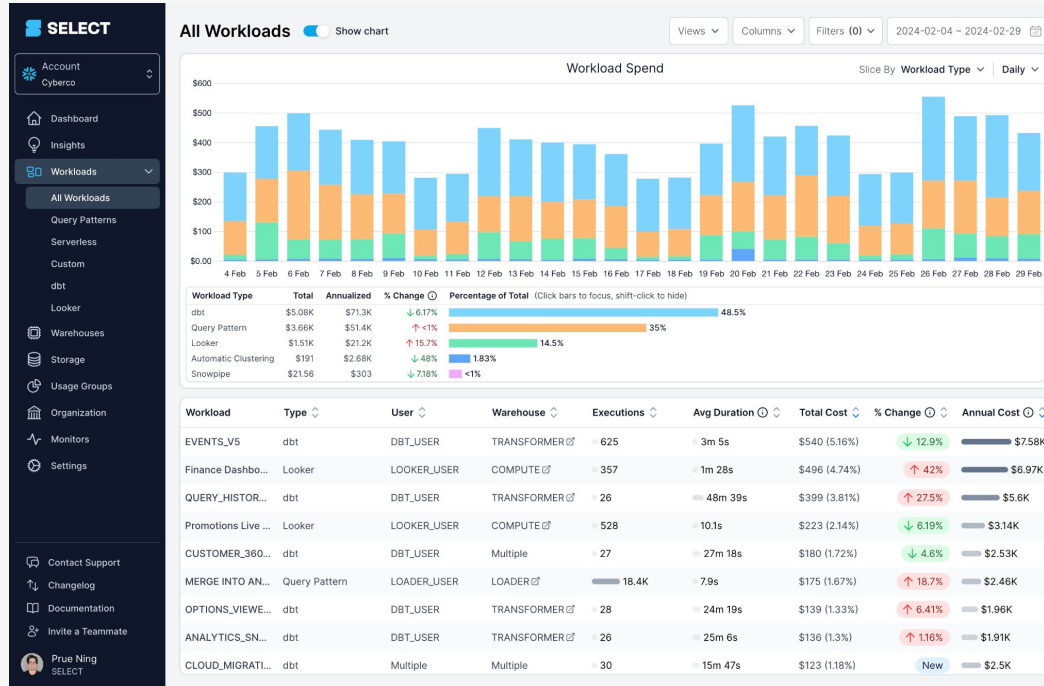
Forecasting

- Create targets and projections
- Use annualized metrics
- Track monthly and quarterly trends
- Alerting

Watch Outs

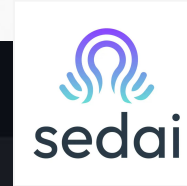
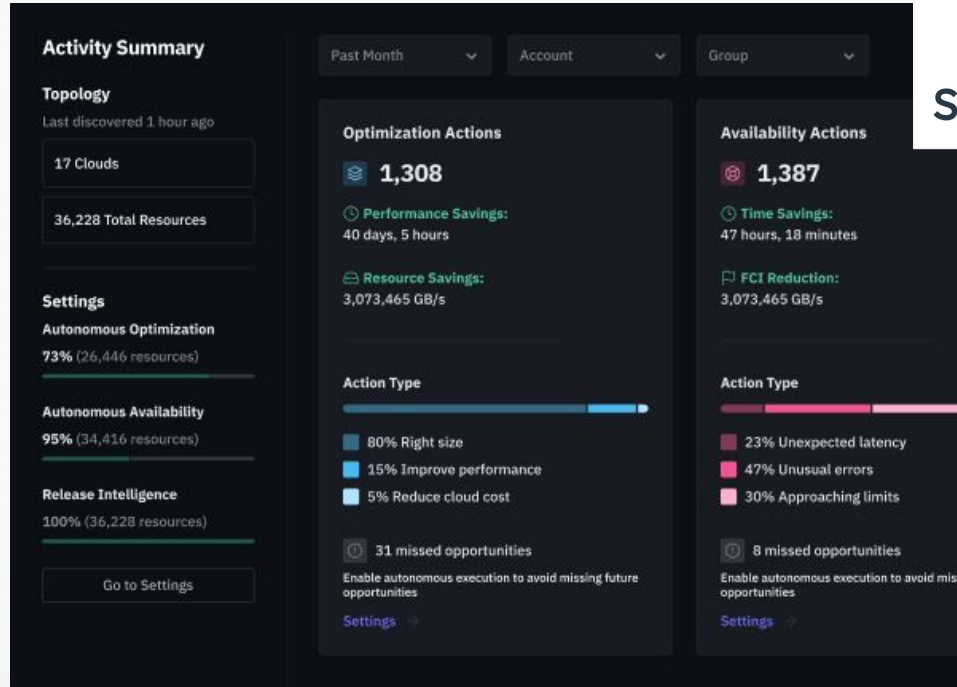
- Be skeptical and avoid linear scaling
- Ensure tool owners understand TCO and scaling costs
- Negotiate contract terms to reduce cost
- Avoid tool lock-in and weigh migration efforts

Cost Reduction Tools



*Disclaimer: I'm not sponsored by SELECT

Autonomous AI



*Disclaimer: I'm not sponsored by Sedai

Takeaways

Takeaways

- 1 Invest in cost measurement and monitoring
- 2 Build cost containment feedback loops
- 3 Drive accountability for cost containment



If you'd like to chat, you can
find me on LinkedIn 🎉