

# Honey, I Shrunk the Kids Bill!

Harjot S. Parmar

Featured

## BigQuery's Ridiculous Pricing Model Cost Us \$10,000 in Just 22 Seconds!!!

5 min read · Mar 19, 2025



We had a dev query the entire datalake.

The query ran for like 2-3 weeks and ended up costing about 75K.



**RedditBeaver42** · 1y ago

Customer did training on Cosmos DB and left instances running. Some timer later they got a call from MS because the bill had reached 45000 USD.

Featured

## BigQuery's Ridiculous Pricing Model Cost Us \$10,000 in Just 22 Seconds!!!

5 min read · Mar 19, 2025



We had a dev query the entire datalake.

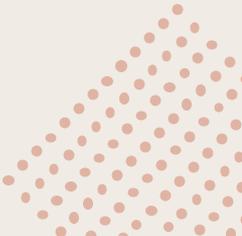
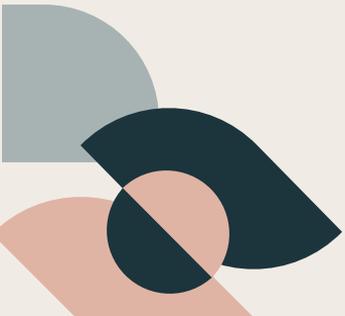
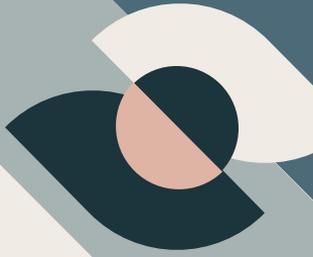
The query ran for like 2-3 weeks and ended up costing about 75K.



**RedditBeaver42** · 1y ago

Customer did training on Cosmos DB and left instances running. Some timer later they got a call from MS because the bill had reached 45000 USD.

**“Cloud comes with convenience...  
and convenience comes at a cost ”**



# HARJOT SINGH PARMAR

LDX3 2025 EDITION

- **Occupation**  
Staff Machine Learning Engineer @ Intuit  
Data Engineering, Generative AI and Agentic
- **Education**  
University of Waterloo, SYDE
- **Accessories**  
Macbook, Hiking Boots, Soccer Ball, Charging Cords
- **Special Interest**  
Data intensive products, architecture optimization,  
Trivia nights

*Trains Models, Trek Mountains, Trust the Process* 

 harjotsparmar





**You are here**



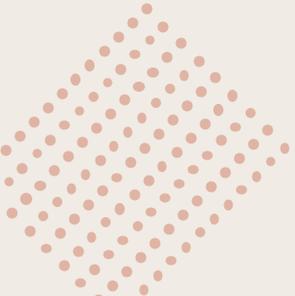
**Start**

**Common  
Patterns  
(Savings 😊)**

**Opex/Audit  
Framework**

**End**

Travel Time: 25m





You are here

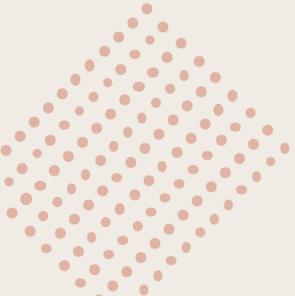


Start

Common  
Patterns  
(Savings 😊)

Opex/Audit  
Framework

End



# Minimal Audit Framework

## Pre-audit

Laying the groundwork.  
Setting up the foundation

---

- Technical Setup
- Cultural Setup

## Audit

Put on your detective hats and start digging.

---

- Weekly OpEx Reviews
- Housekeeping Days
- Bounties and Rewards

## Post-Audit

Sustaining the gains

---

- Diagramming and Optimization
- Education and Record Keeping
- Automate common patterns

# Groundwork

## Pre-audit

Laying the groundwork.  
Setting up the foundation

- Technical Setup
- Cultural Setup



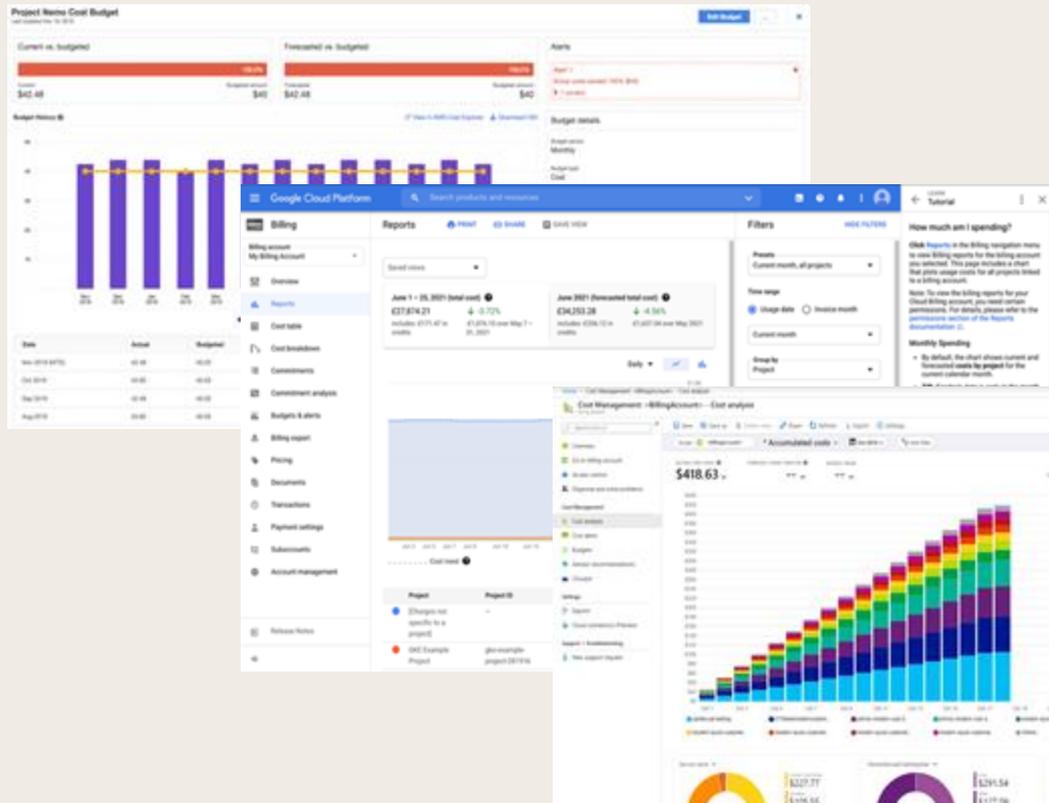
# Cloud Provider Setup

## Pre-audit

Laying the groundwork.  
Setting up the foundation

Technical Setup

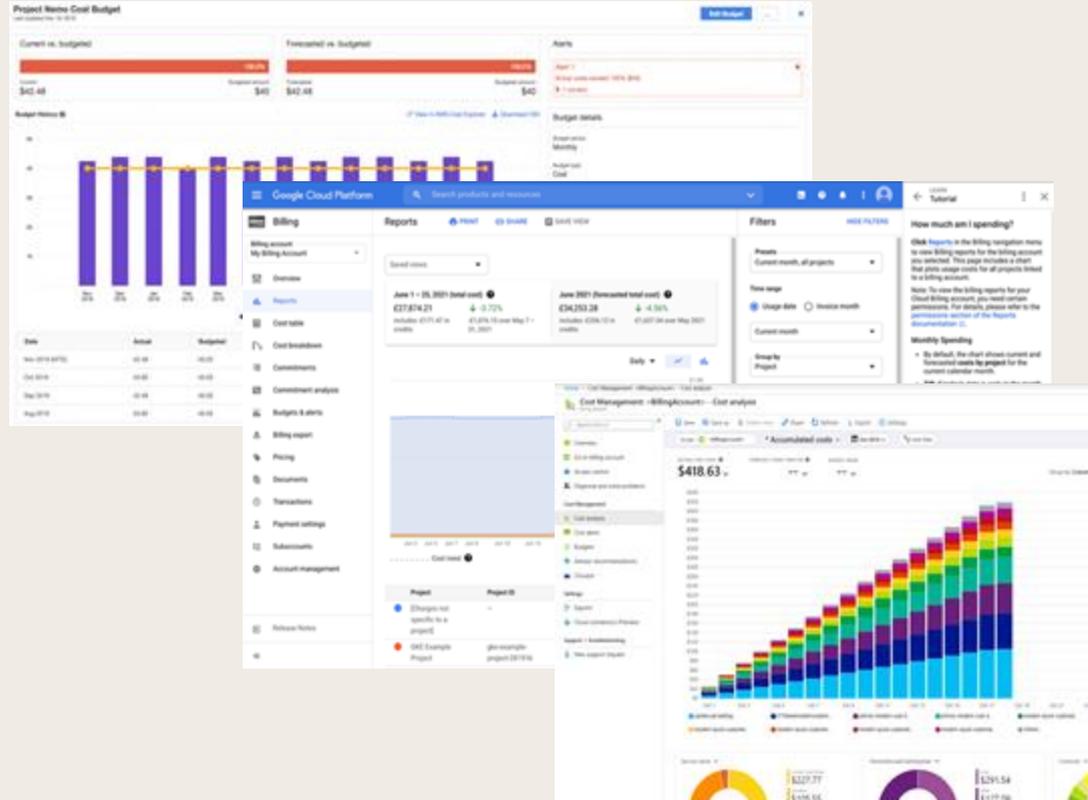
Cultural Setup



# Cloud Provider Setup

## Core Features

- Budgeting
- Forecasting
- Anomaly Detection
- Data Exports & Custom Reporting



# Microservices Setup

## Pre-audit

Laying the groundwork.  
Setting up the foundation

Technical Setup

Cultural Setup

Metrics

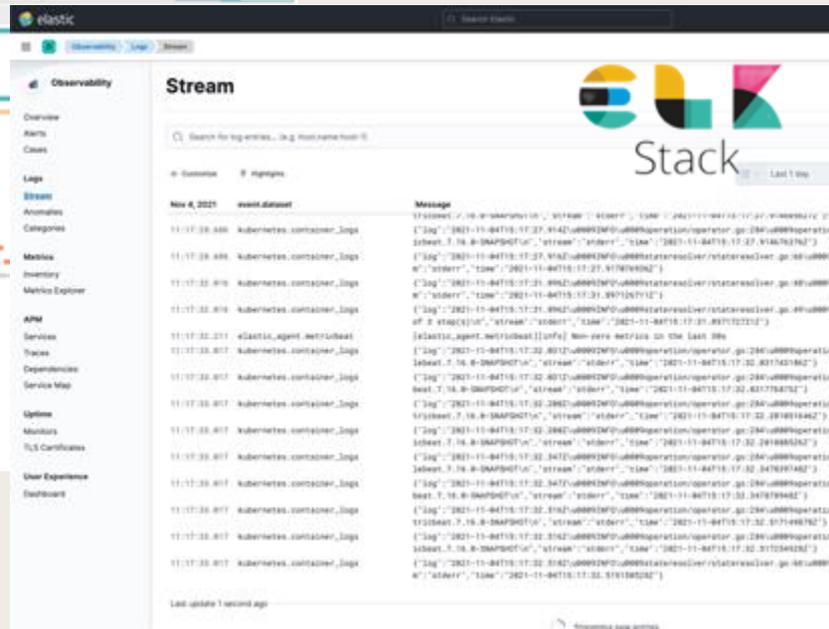
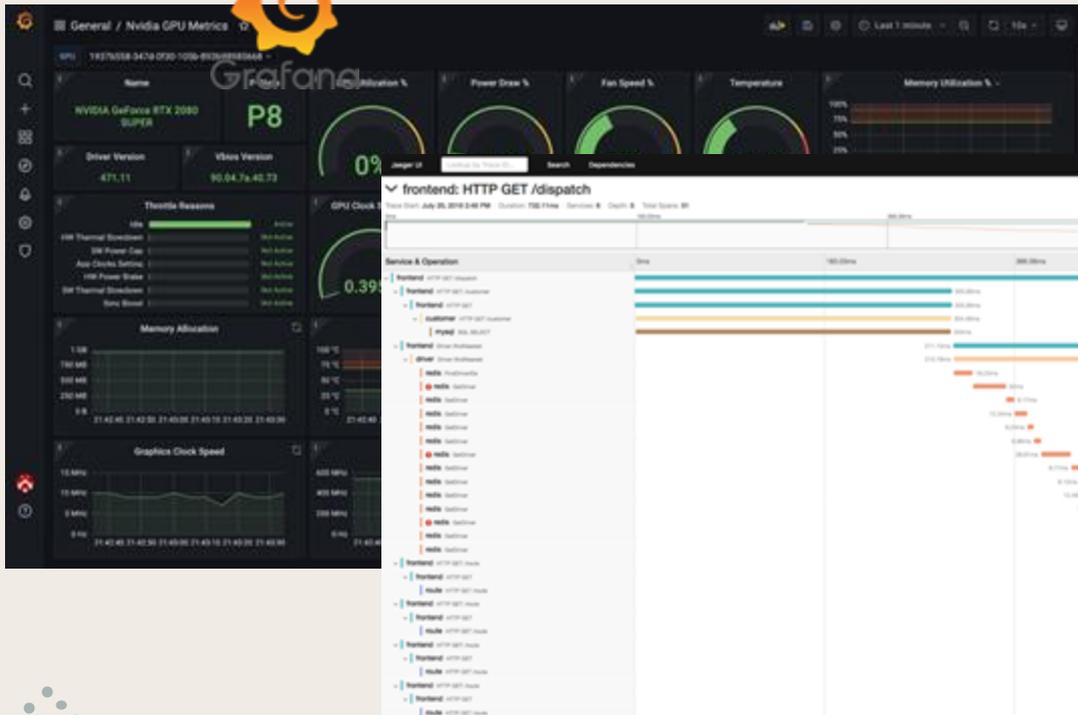


Logs



Tracing

# Microservices Setup



<https://grafana.com/> | <https://www.elastic.co/observability/log-monitoring> | <https://www.jaegertracing.io/docs/1.49/>

# Data Engineering & Infra

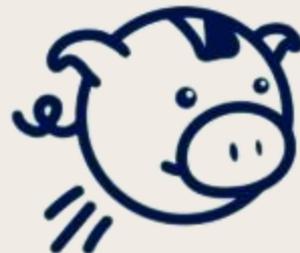
## Pre-audit

Laying the groundwork.  
Setting up the foundation

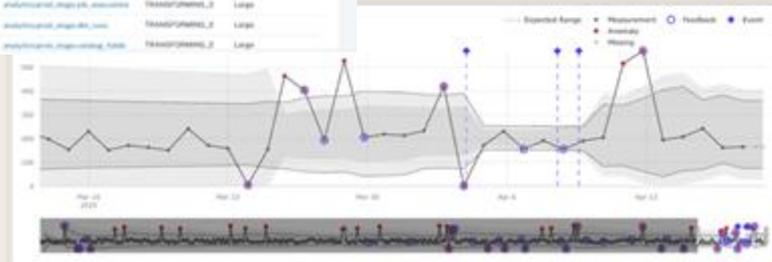
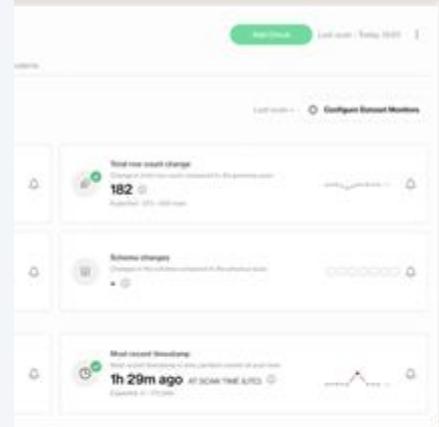
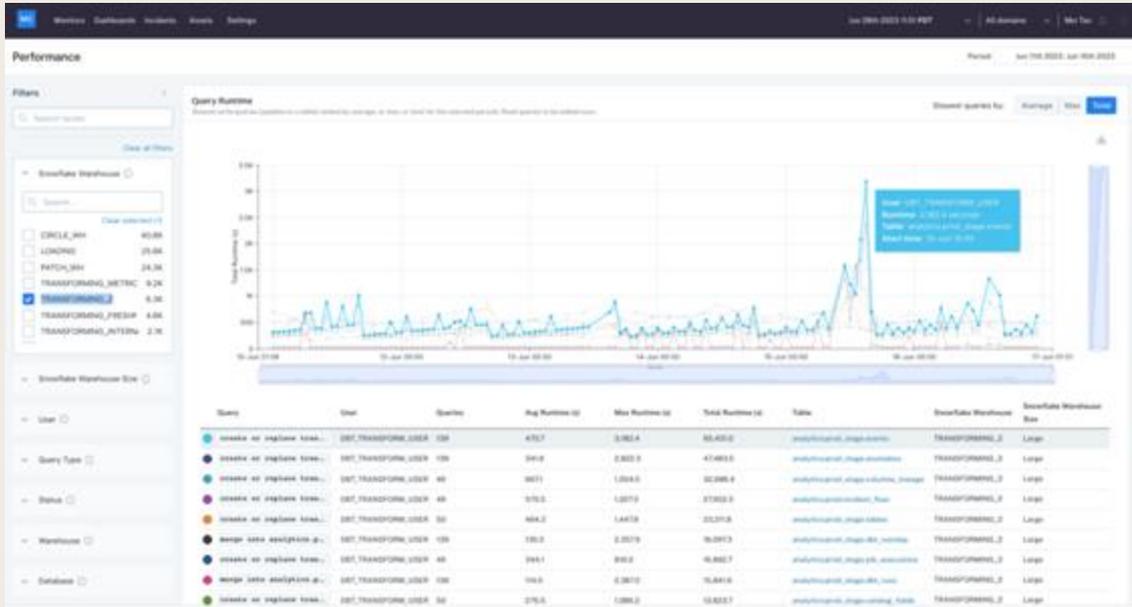
- 
- Technical Setup
  - Cultural Setup

**SODA** 

**MC** **MONTE  
CARLO**



# Data Infra



# Cultural Setup

## Pre-audit

Laying the groundwork.  
Setting up the foundation

- ✓ Technical Setup
- ✓ Cultural Setup



# Enable Engineers

## Pre-audit

Laying the groundwork.  
Setting up the foundation

- ✓ Technical Setup
- ✓ Cultural Setup



# Enable Engineers

## Pre-audit

Laying the groundwork.  
Setting up the foundation

- ✓ Technical Setup
- ✓ Cultural Setup



# Recognition and Incentives

Show me the incentive,  
I'll show you the  
outcome



Charlie Munger

## Pre-audit

Laying the groundwork.  
Setting up the foundation

- Technical Setup
- Cultural Setup



You've received an award!

Emily Wright - The Best Colleague!!!



Congrats Emily!



Job well done!



# Minimal Audit Framework

## Pre-audit

Laying the groundwork.  
Setting up the foundation

---

- Technical Setup
- Cultural Setup

## Audit

Put on your detective hats and start digging.

---

- Weekly OpEx Reviews
- Housekeeping Days
- Bounties and Rewards



# Weekly OpEx Reviews

Regular meetings focused solely on operational expenditures, including cloud costs

## Audit

Put on your detective hats and start digging.

- Weekly OpEx Reviews
- Housekeeping Days
- Bounties and Rewards

# Monthly Housekeeping Days

Routine chores that prevent cost creep

## Audit

Put on your detective hats and start digging.

- Weekly OpEx Reviews
- Housekeeping Days



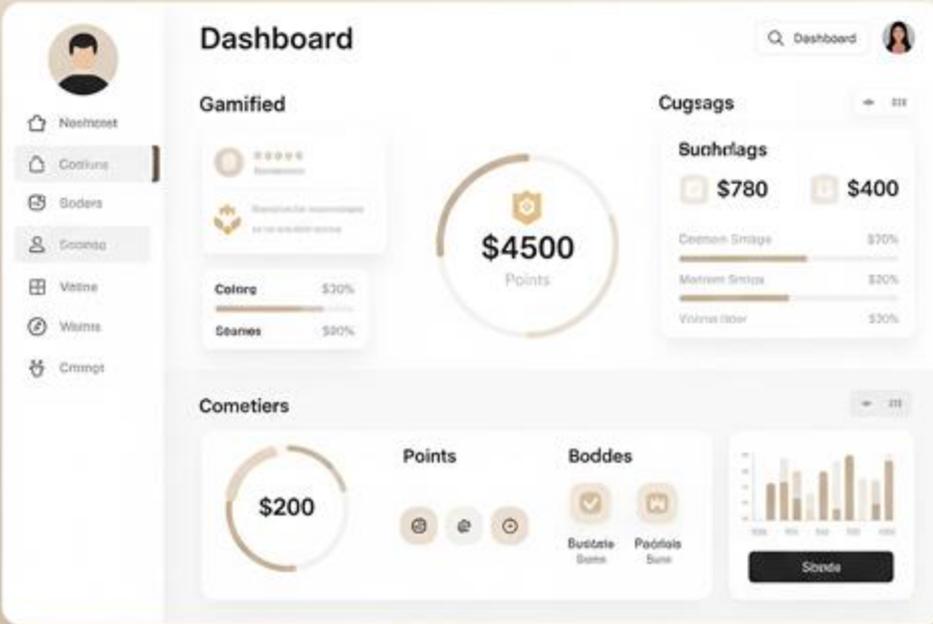
# Bounties And Rewards

Actively incentivize engineers to find and suggest cost-saving opportunities

## Audit

Put on your detective hats and start digging.

- Weekly OpEx Reviews
- Housekeeping Days
- Bounties and Rewards



# Minimal Audit Framework

## Pre-audit

Laying the groundwork.  
Setting up the foundation

---

- Technical Setup
- Cultural Setup

## Audit

Put on your detective hats and start digging.

---

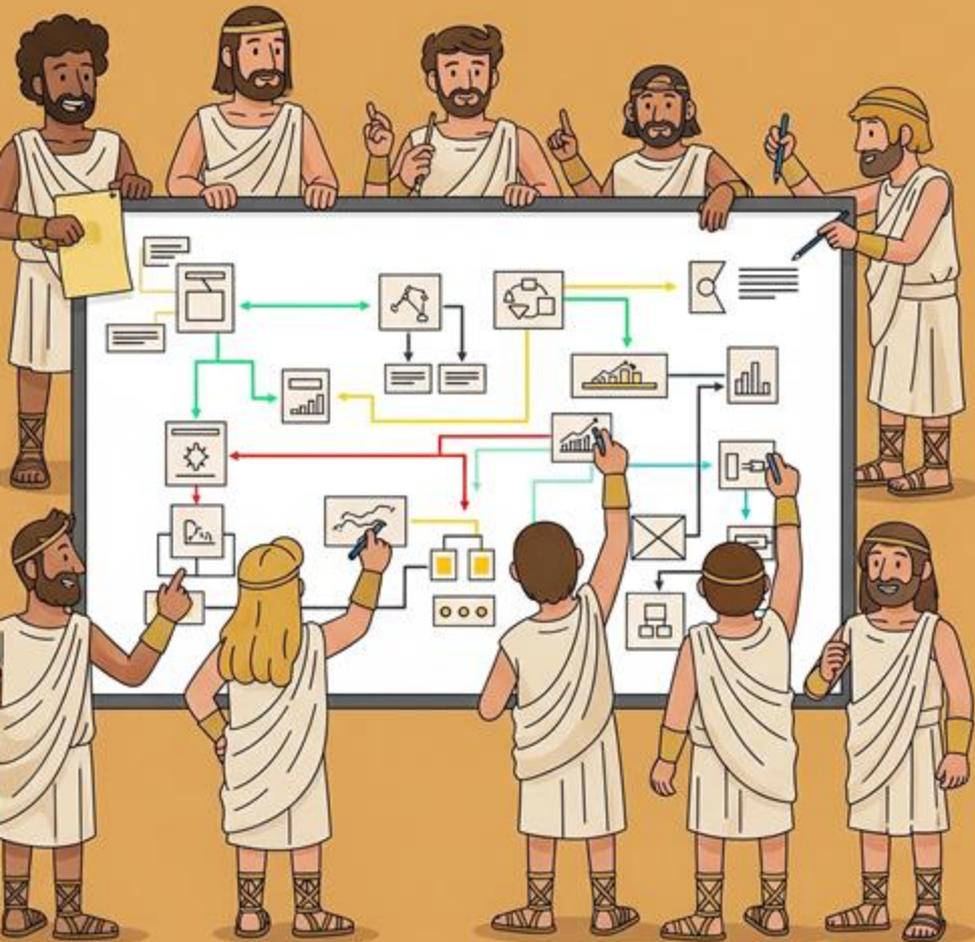
- Weekly OpEx Reviews
- Housekeeping Days
- Bounties and Rewards

## Post-Audit

Sustaining the gains

---

- Diagramming and Optimization
- Education and Record Keeping
- Automate common patterns



# Post-Audit: Diagramming and Optimization

## Post-Audit

Sustaining the gains

- Diagramming and Optimization
- Education and Record Keeping
- Automate common patterns

# Post-Audit: Education and Record Keeping

## Post-Audit

Sustaining the gains

- Diagramming and Optimization
- Education and Record Keeping
- Automate common patterns

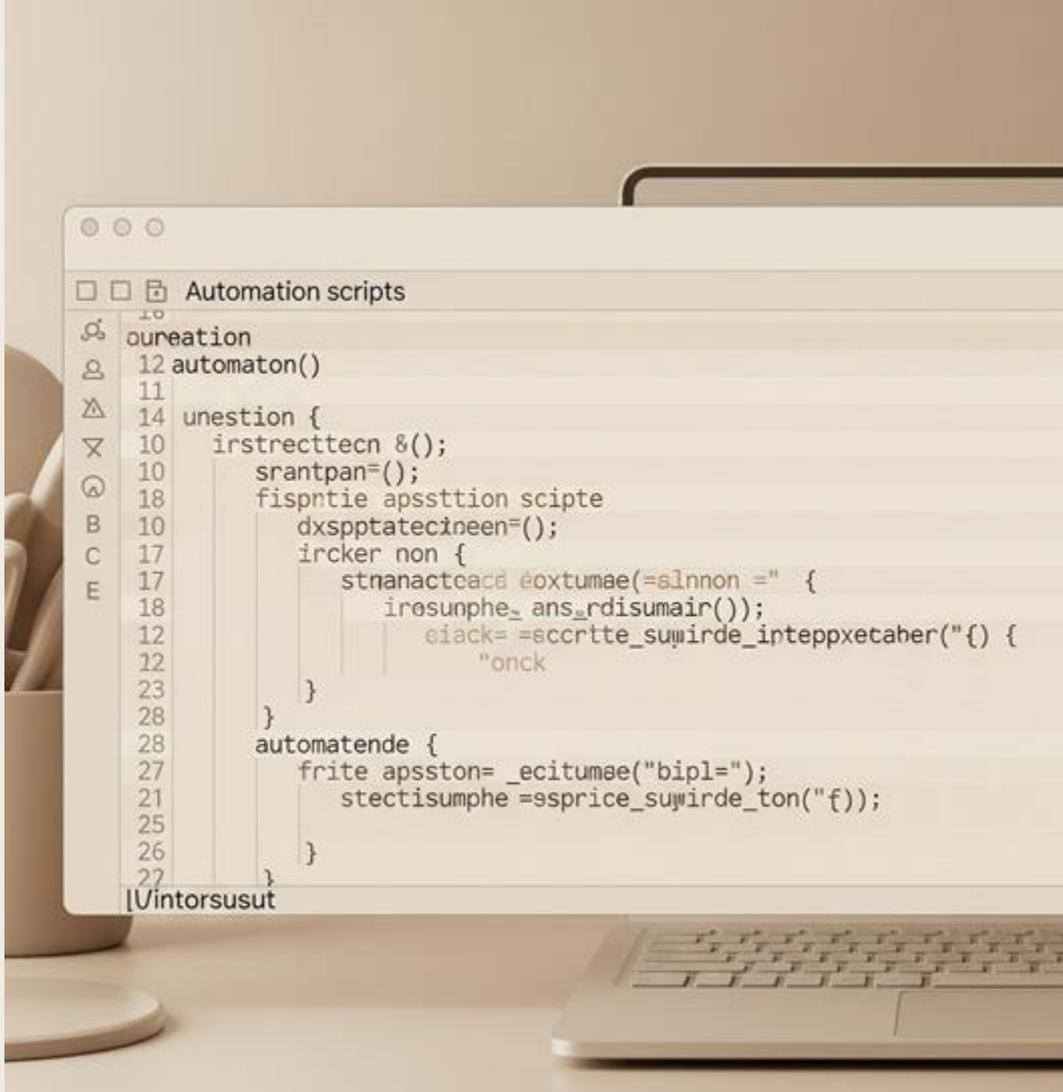


# Post-Audit: Automate Common Patterns

## Post-Audit

Sustaining the gains

- ✓ Diagramming and Optimization
- ✓ Education and Record Keeping
- ✓ Automate common patterns



```
Automation scripts
10
11
12 automaton()
13
14 unestion {
15     irstrecttecn &();
16     srantpan=();
17     fisptie apssttion scipte
18     dxspptatecineen=();
19     ircker non {
20         stnanactead eoxtumae(=alnnon =" {
21             irosunphe_ ans_rdisumair());
22             eiack= =eccrite_suyirde_inteppxetaber("{} {
23                 "onck
24         }
25     }
26 }
27
28 automatende {
29     frite apsston= _ecitumae("bipl=");
30     stectisumphe =sprice_suyirde_ton("f));
31 }
32
33 |Uintorsusut
```

# Post-Audit: Automate Common Patterns

Eg:

- Automated shutdown
- Script for the cleanup
- lifecycle policy



```
Automation scripts
10
11
12 automaton()
13
14 unestion {
15     firstrecttecn &();
16     srantpan=();
17     fisptie apssttion scipte
18     dxsptatecineen=();
19     ircker non {
20         stnanactead eoxtumae(=alnnon =" {
21             irosunphe_ ans_rdisumair());
22             eiack= =eccrite_suyirde_inteppxetaber("{} {
23                 "onck
24         }
25     }
26 }
27
28 automatende {
29     frite apsston= _ecitumae("bipl=");
30     stectisumphe =esprice_suyirde_ton("f));
31 }
32 }
```

# Minimal Audit Framework

## Pre-audit

Laying the groundwork.  
Setting up the foundation

- ✓ Technical Setup
- ✓ Cultural Setup

## Audit

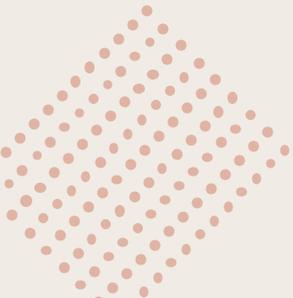
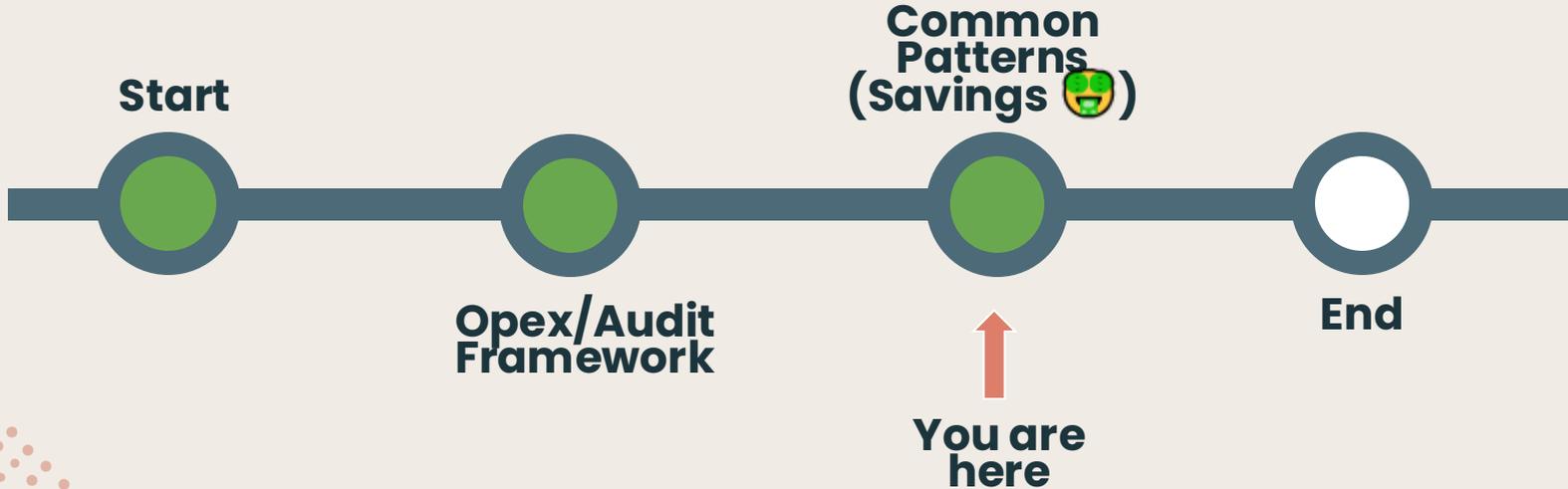
Put on your detective hats and start digging.

- ✓ Weekly OpEx Reviews
- ✓ Housekeeping Days
- ✓ Bounties and Rewards

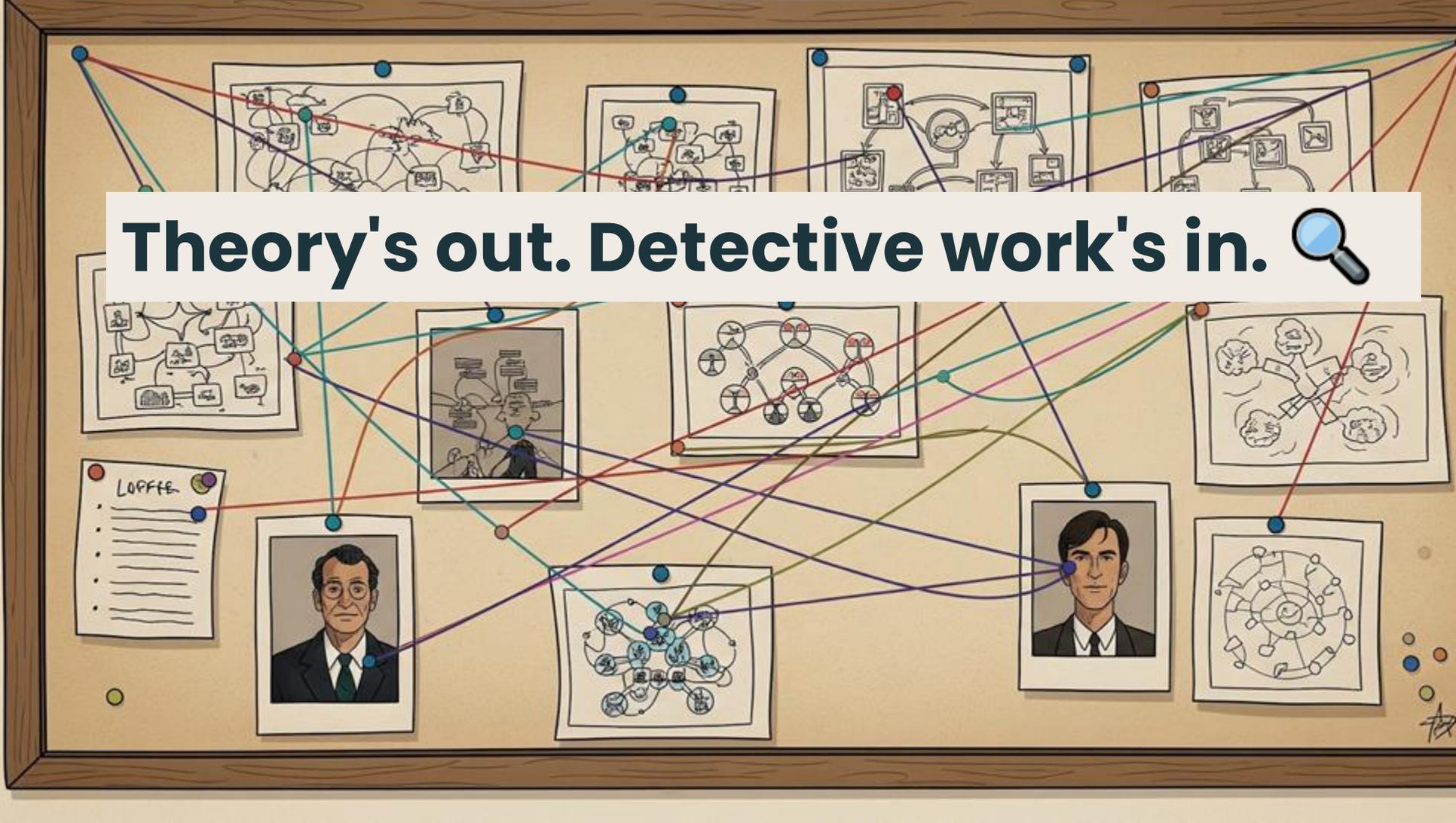
## Post-Audit

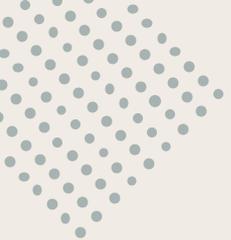
Sustaining the gains

- ✓ Diagramming and Optimization
- ✓ Education and Record Keeping
- ✓ Automate common patterns



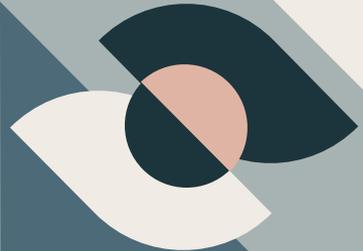
**Theory's out. Detective work's in.**





# Optimizing for Hidden Cost

Identifying Common patterns for wastage

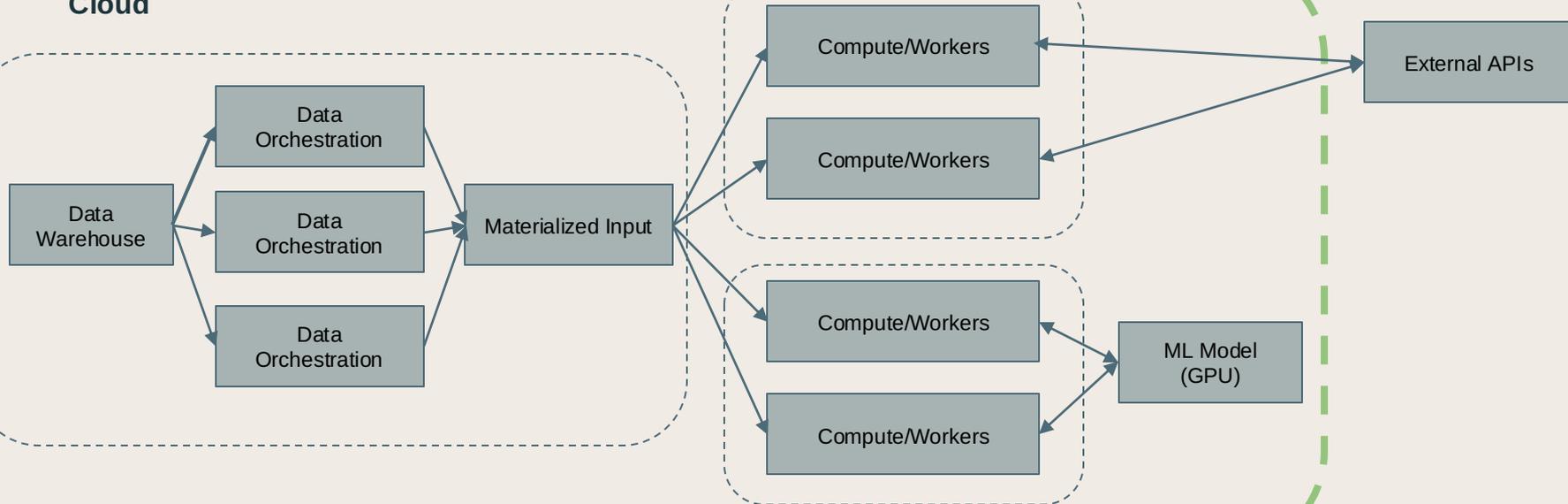
An illustration of an iceberg floating in blue water. The top part of the iceberg is above the water line, and the bottom part is submerged. The text 'OBVIOUS COSTS' is written in blue, bold, uppercase letters on the visible part of the iceberg, and 'HIDDEN COSTS' is written in blue, bold, uppercase letters on the submerged part.

**OBVIOUS  
COSTS**

**HIDDEN  
COSTS**

# ML Pipeline

Cloud





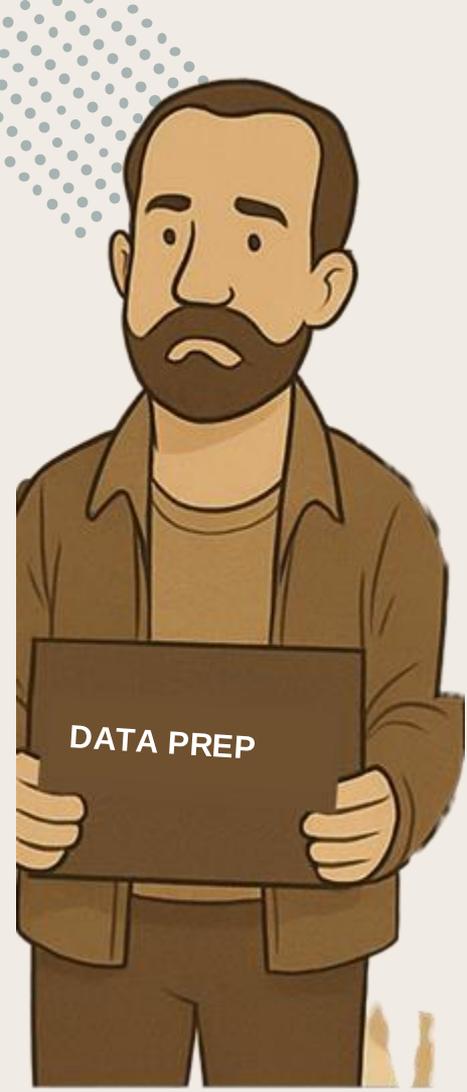
DATA PREP

COMPUTE

GPU

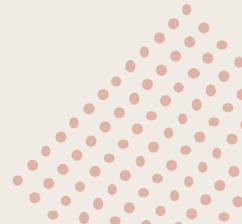
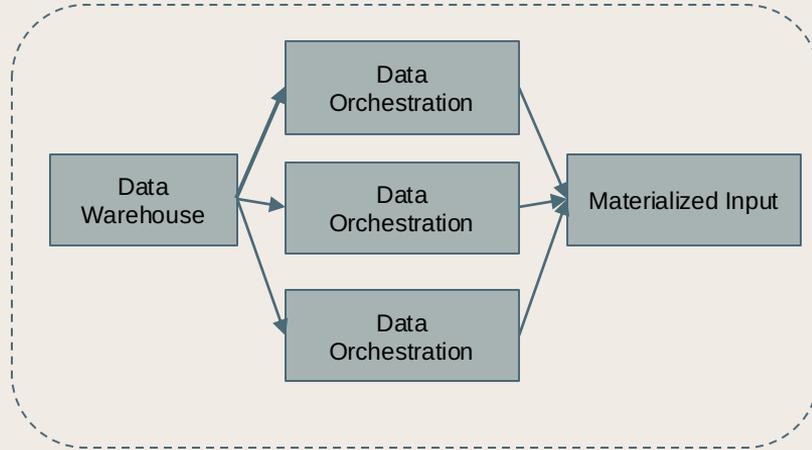
External APIs

LOGS  
NETWORKING



# Suspect #1

## Data Preparation





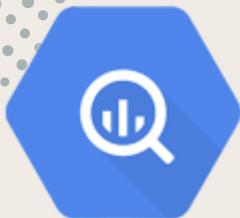
DATA LAYER

## Excessive data scans

**\$4,000 per run**

- Select \* from \_\_\_\_\_

The screenshot shows a SQL query editor interface. At the top, there is a search icon, the text "Untitled query", a blue "Run" button, a "Save" button with a dropdown arrow, and a "Download" button. Below the toolbar, a SQL query is entered: `1 SELECT * FROM `project.dataset.giant_table` LIMIT 1000`. The query is highlighted in a light blue background. Below the query, a green checkmark icon is followed by the text: "This query will process 818.17 TB when run."



# Bigquery

## Clustered and Partitioned Tables

**Orders table**  
Not Clustered; Not partitioned

Order_Date	Country	Status
2022-08-02	US	Shipped
2022-08-04	JP	Shipped
2022-08-05	UK	Canceled
2022-08-06	KE	Shipped
2022-08-02	KE	Canceled
2022-08-05	US	Processing
2022-08-04	JP	Processing
2022-08-04	KE	Shipped
2022-08-06	UK	Canceled
2022-08-02	UK	Processing
2022-08-05	JP	Canceled
2022-08-06	UK	Processing
2022-08-05	US	Shipped
2022-08-06	JP	Processing
2022-08-02	KE	Shipped
2022-08-04	US	Shipped

**Orders table**  
Clustered by Country; Not partitioned

Order_Date	Country	Status
2022-08-04	JP	Shipped
2022-08-04	JP	Processing
2022-08-05	JP	Canceled
2022-08-06	JP	Processing
2022-08-06	KE	Shipped
2022-08-02	KE	Canceled
2022-08-04	KE	Shipped
2022-08-02	KE	Shipped
2022-08-05	UK	Processing
2022-08-06	UK	Canceled
2022-08-02	UK	Canceled
2022-08-06	UK	Processing
2022-08-02	US	Shipped
2022-08-05	US	Processing
2022-08-05	US	Shipped
2022-08-04	US	Shipped

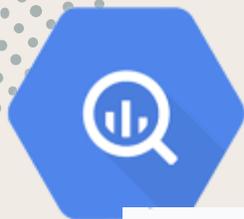
**Orders table**  
Clustered by Country; Partitioned by Order\_Date (Daily)

	Order_Date	Country	Status
Partition: 2022-08-02	2022-08-02	KE	Shipped
	2022-08-02	KE	Canceled
Clusters: Country	2022-08-02	UK	Processing
	2022-08-02	US	Shipped

	Order_Date	Country	Status
Partition: 2022-08-04	2022-08-04	JP	Shipped
	2022-08-04	JP	Processing
Cluster: Country	2022-08-04	KE	Shipped
	2022-08-04	US	Shipped

	Order_Date	Country	Status
Partition: 2022-08-05	2022-08-05	JP	Canceled
	2022-08-05	UK	Canceled
Cluster: Country	2022-08-05	US	Shipped
	2022-08-05	US	Processing

	Order_Date	Country	Status
Partition: 2022-08-06	2022-08-06	JP	Processing
	2022-08-06	KE	Shipped
Cluster: Country	2022-08-06	UK	Canceled
	2022-08-06	UK	Processing



# Bigquery

```
1 CREATE TABLE
2 | project.dataset.giant_table_partitioned (transaction_id INT64, transaction_date DATE)
3 PARTITION BY
4 | transaction_date
5 AS (
6 | SELECT
7 | transaction_id, transaction_date
8 | FROM
9 | project.dataset.giant_table
10
11 );
```



Untitled query



Run



Save ▾



Download

```
1 select * from `project.dataset.giant_partitioned_table` limit 1000
```



This query will process 818.17 TB when run.



# Bigquery

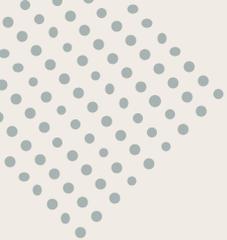
```
1 CREATE TABLE
2 | `project.dataset.giant_partitioned_table_with_required_filter`
3 (transaction_id INT64, transaction_date DATE)
4 PARTITION BY
5 | transaction_date
6 OPTIONS (
7 | require_partition_filter = TRUE) as
8 | select * from `project.dataset.giant_table`
9
```



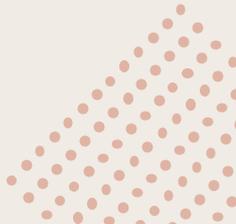
```
1 select transaction_id, transaction_date from
2 | `project.dataset.giant_partitioned_table_with_required_filter`
3 | where transaction_date = "2025-06-16"
```

✓ This query will process 64.23 MB when run.



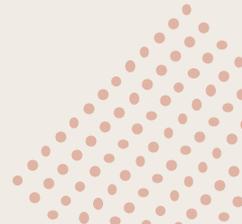
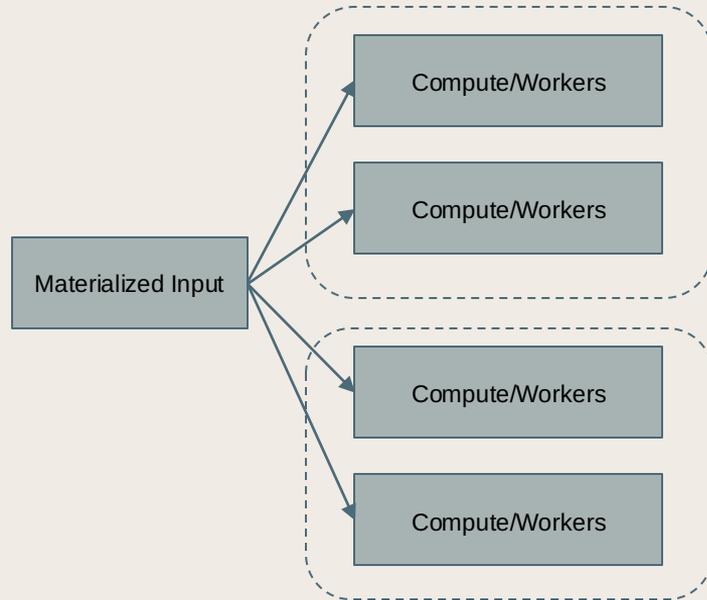


# Remedies

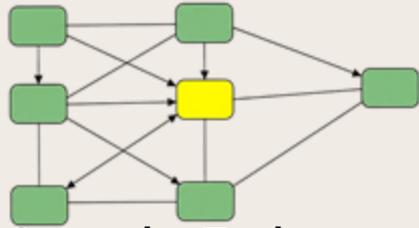
- **Reduce data before joins:** Apply `WHERE` clauses as early as possible in queries to minimize data shuffling during joins.
  - **Choose appropriate pricing model:** Select on-demand pricing for unpredictable workloads and capacity pricing for consistent, high-volume workloads.
  - **Avoid frequent table overwrites:** Use incremental data loads instead of overwrites to prevent hidden storage costs from time travel retention.
  - **Use query quotas:** Set custom daily query quotas at the project or user level to limit the amount of data processed.
    - Utilize `Maximum bytes billed` per query: Define a maximum number of bytes a query can process to avoid unexpected expenses.
- 

# Suspect #2

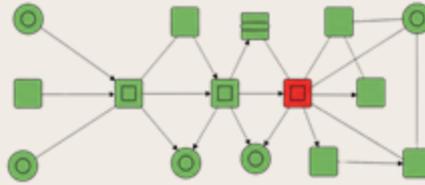
## Distributed Compute



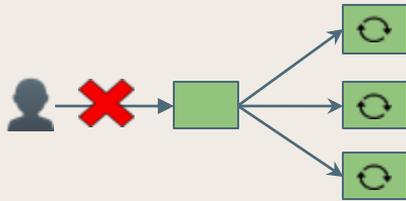
# Compute Wastage



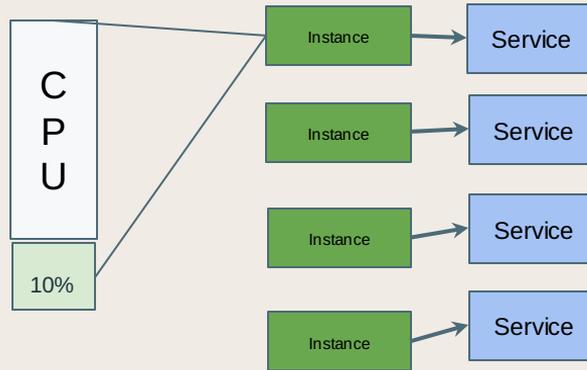
**Straggler Tasks**



**Partial/Failed Executions**



**Orphan Computations**



**Mixed I/O and CPU bound executions**

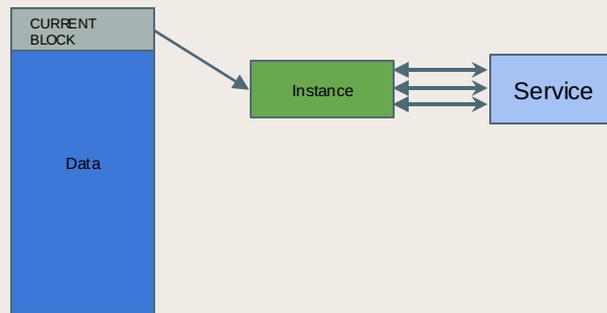
- Underutilized compute
- Excessive Idle time
- Repeated tasks
- Overprovisioning
- Dangling temporary artifacts

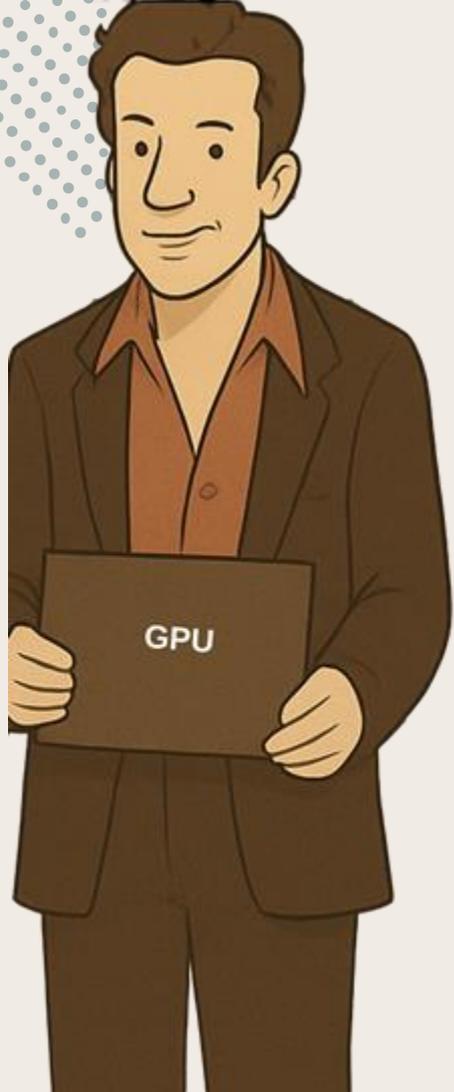
## Remedies

- Implement custom checkpoints behaviour in batch pipelines
- Delta tracking queries as intermediate steps
  
- Implement custom data loaders for paged reading from tables
- Investing in building async API executors on the single nodes



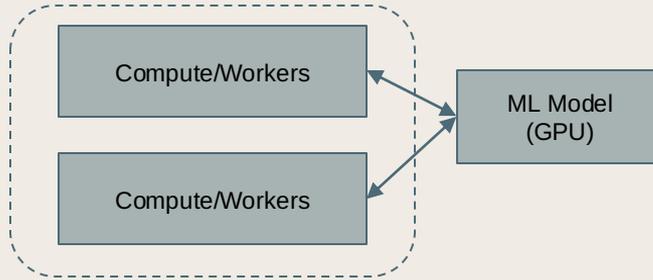
Complete input - partial result = 



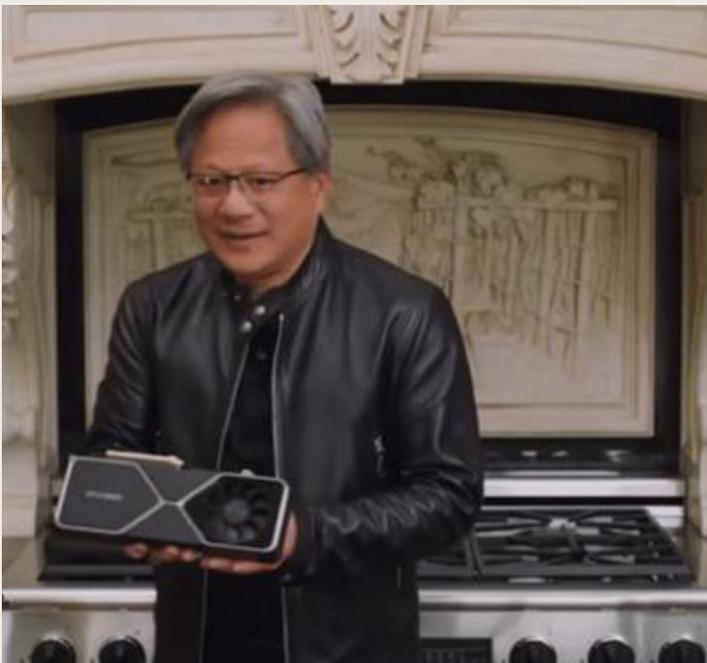


# Suspect #3

## Graphical Processing Unit (GPU)



# GPU Wastage



Beautiful beast, very powerful for ML inferences, but a few understand them well...and that leads to underutilized instances.

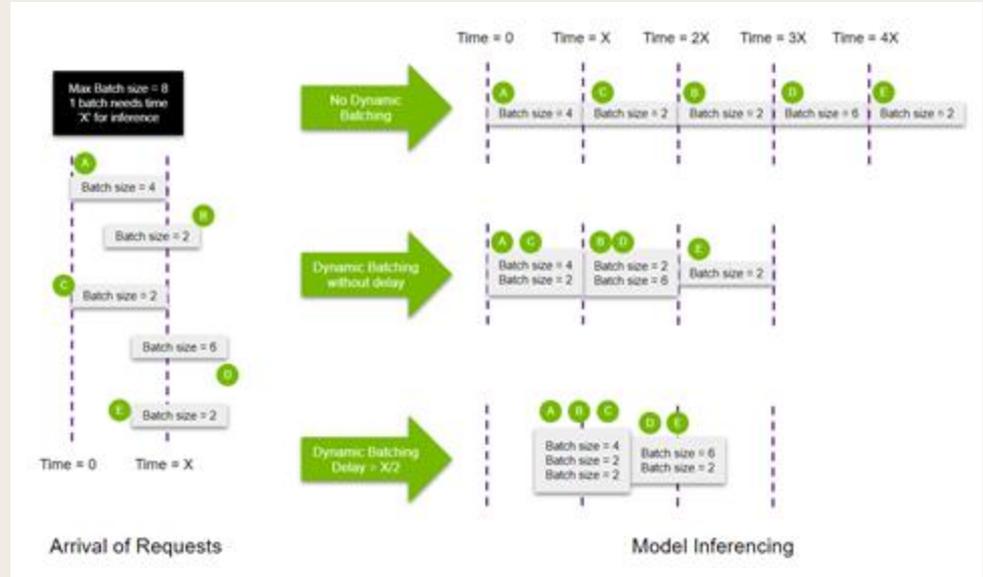


## Misconfigured executions

- VRAM not utilized
- Incompatible Driver and Toolkit versions on host and image
- FastAPI service wrappers for inference

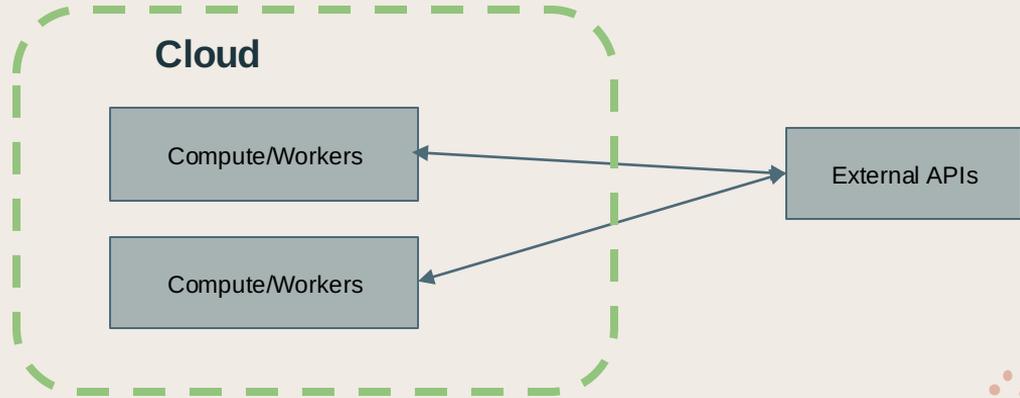
# Remedies

- Automate the shutdowns
- Monitor for GPU memory utilization with custom dashboards
- Optimize GPU inference with proprietary inference servers with features:
  - Multi instance deployment
  - server side batching for inference



# Suspect #4-5

## External API Logs Egress



# API and Networking kerfuffle

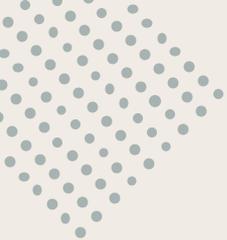
- Chatty Services generating lot of logs, incorrect log levels
- External API calls/Data Processing and Sharing over public internet
- Cross-Region and Cross-Availability Zone Transfers and Networking

## No such thing as too much logging?



engineerL · 3y ago

A team at my company accidentally blew ~100k on Azure Log Analytics during the span of a few days. They set the logging verbosity to a hitherto untested level and threw in some extra replicas as well. When they announced their mistake on Slack, I learned that yes, there is such a thing as too much logging.



# Remedies

- Understand logging service pricing model
- Ensure proper log levels and verbosity, proper Structured Logs
- Implement log segregation and offload to cold storage if needed
- Avoid cross region service calls
- Use VPC peering for cross cloud service communication

JULY 6, 2020

**Cloud Logging Optimization - How We Saved Over \$140k  
in Logging Costs**

