



# Reliability, Cost, Performance

Pick 3

---

Suyog Soti  
2026



# About Me

Why am I here?

University of Colorado Boulder (2020)

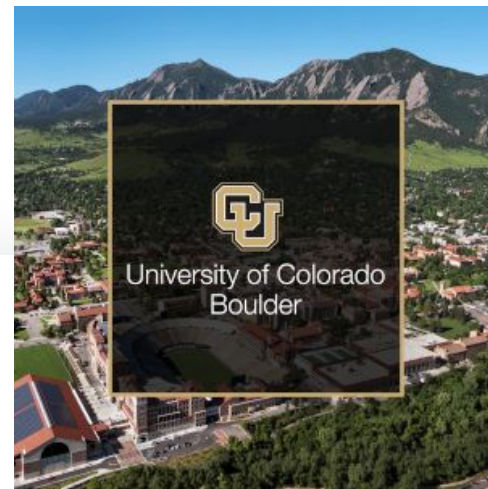
*Applied Math and Computer Science*

Google (2020 – 2023)

*Worked on resource allocation and isolation for multi-tenant services on GCP*

Databricks (2023 – Present)

*Overload Protection, Service Discovery and Load balancing*



# Obligatory: About Databricks :)

The unified data platform



## Store big data

- Delta/Iceberg Formats
- Store it in S3
- Metadata in Unity Catalog



## Query big data

- Spark to query data fast
- Governance via Unity Catalog
- Build business dashboards



## Build on top of the data

- Create agents
- Fine tune models
- Lakewatch
- Lakehouse Apps

# Common Wisdom

Build a low latency, highly scalable service with a 99.999 uptime



# Dimensions

## Autoscale

- 0 downtime deployments
- Scale horizontally with load
- Resource Utilization

### Why care about these

- Reliability
- Cost

## Storage

- Keep consistent and durable state via some remote storage layer

- Performance
- Reliability
- Cost

## Caching

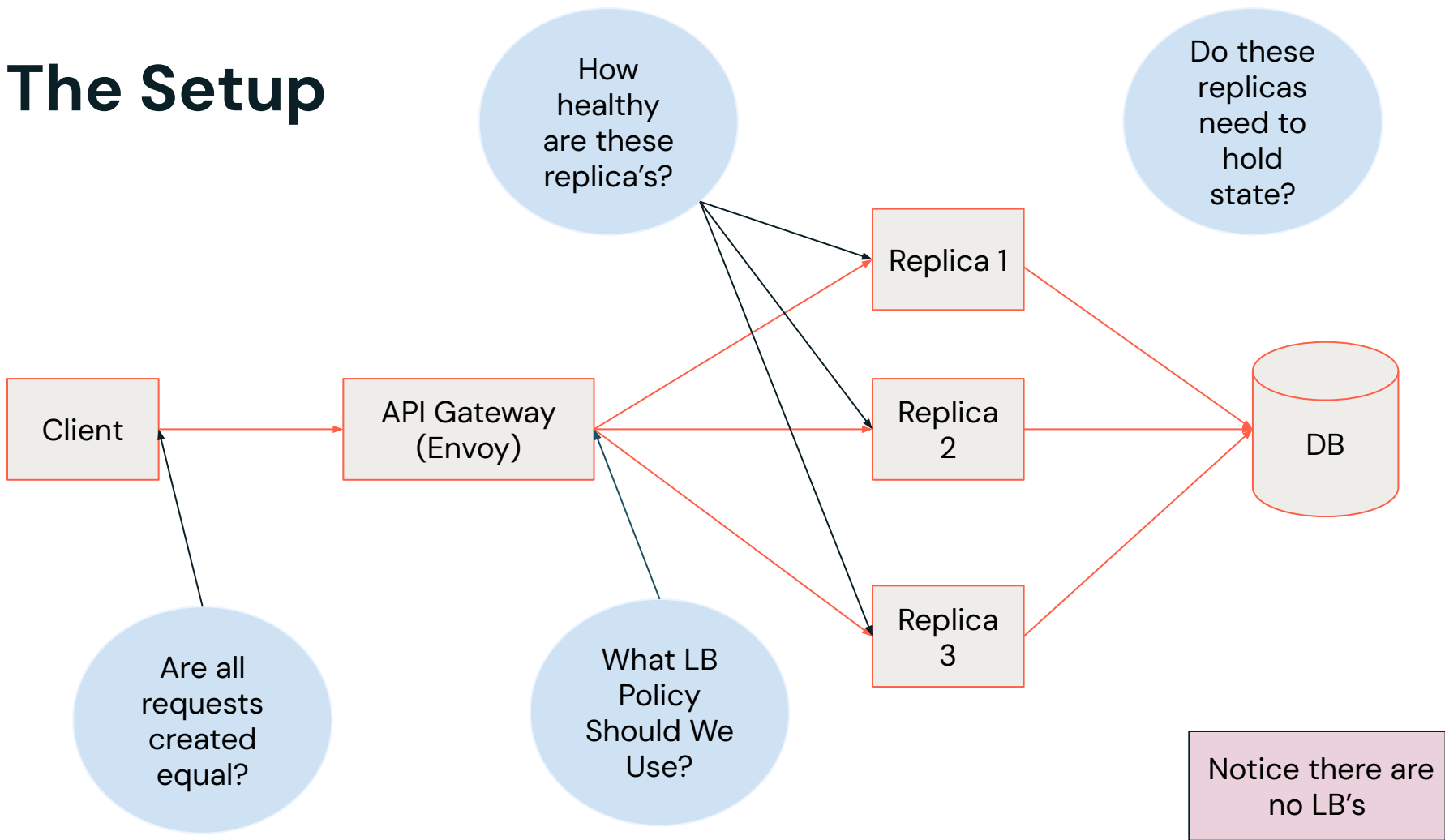
Scalable and eventually consistent reads

- Fall back to storage layer

- Performance
- Reliability
- Cost

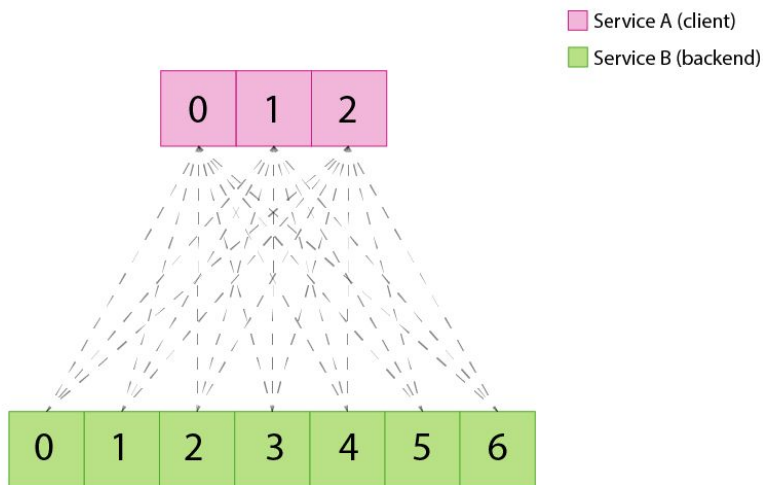


# The Setup



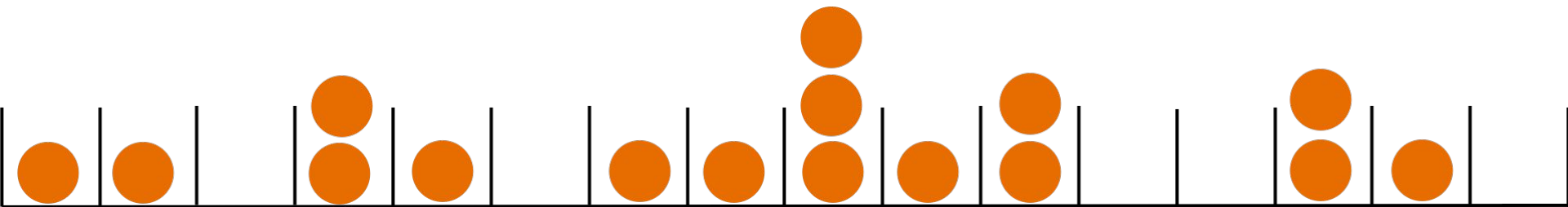
# LB Algorithms – Stateless

- Random Routing
  - Uniform Distribution
- P2C
  - Assign weights, pick 2 random replicas, choose replica with lower weight



# LB Algorithms – Static Stateful

- Recall [LB Lecture](#)
- Hash key matters
- Requests are not homogenous



# Common Wisdom

## Autoscale

- 0 downtime deployments
- Scale horizontally with load
- Resource Utilization

### With stateless services

- 0 downtime deployments
- Scale horizontally with load

## Storage

- Keep consistent and durable state via some remote storage layer

- Databases don't scale well per tenant
- Network IO Overhead
- Serialization Overhead

## Caching

Scalable and eventually consistent reads

- Fall back to storage layer
- Low cache hit rates
- Network IO Overhead
- Serialization Overhead



# Why not static sharding?

## Autoscale

- 0 downtime deployments
- Scale horizontally with load
- Resource Utilization

### With statically sharded services

- 0 downtime deployments
- Scale horizontally with load
- **Hot Spots**



## Storage

- Keep consistent and durable state via some remote storage layer

- Databases scale well per tenant (batch writes/reads)
- Network IO Overhead
- Serialization Overhead

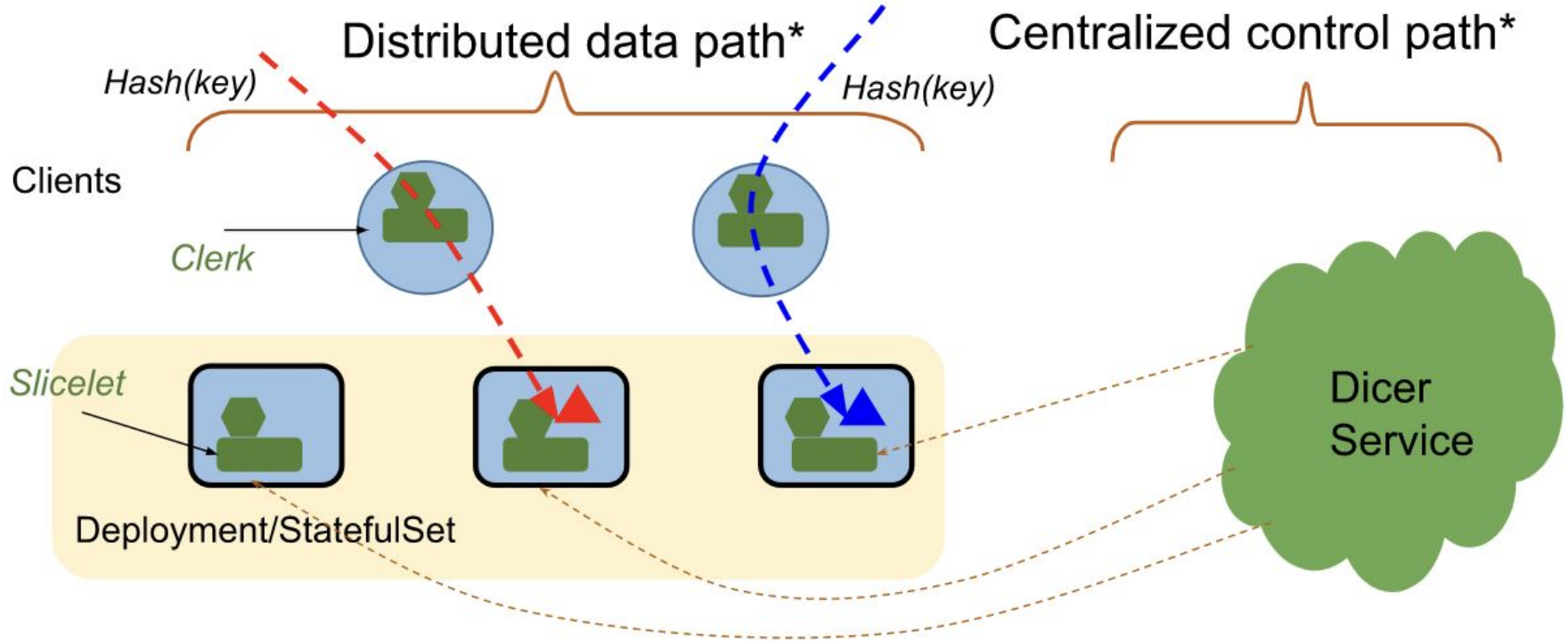
## Caching

Scalable and eventually consistent reads

- Fall back to storage layer
- High cache hit rates
- Low Network IO Overhead
- Low Serialization Overhead



# The Way - Dynamic Autosharder



# Dynamic Autosharder

## Autoscale

- 0 downtime deployments
- Scale horizontally with load
- Resource Utilization

### With dynamically sharded services

- 0 downtime deployments
- Scale horizontally with load

## Storage

- Keep consistent and durable state via some remote storage layer

- Databases scale well per tenant (batch writes/reads)
- Network IO Overhead
- Serialization Overhead

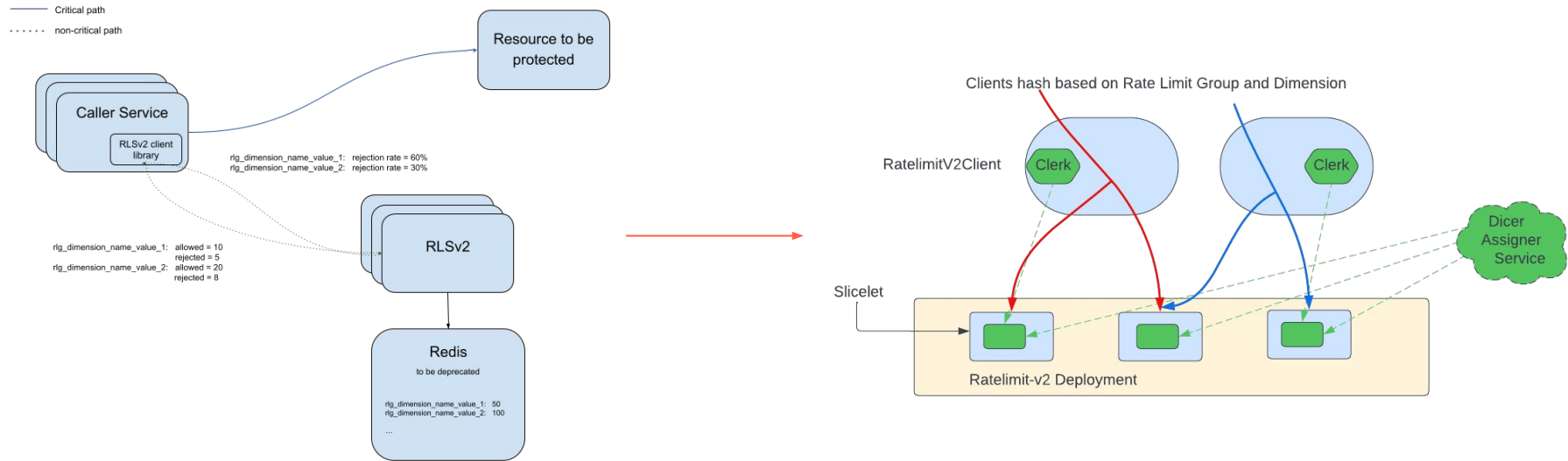
## Caching

Scalable and eventually consistent reads

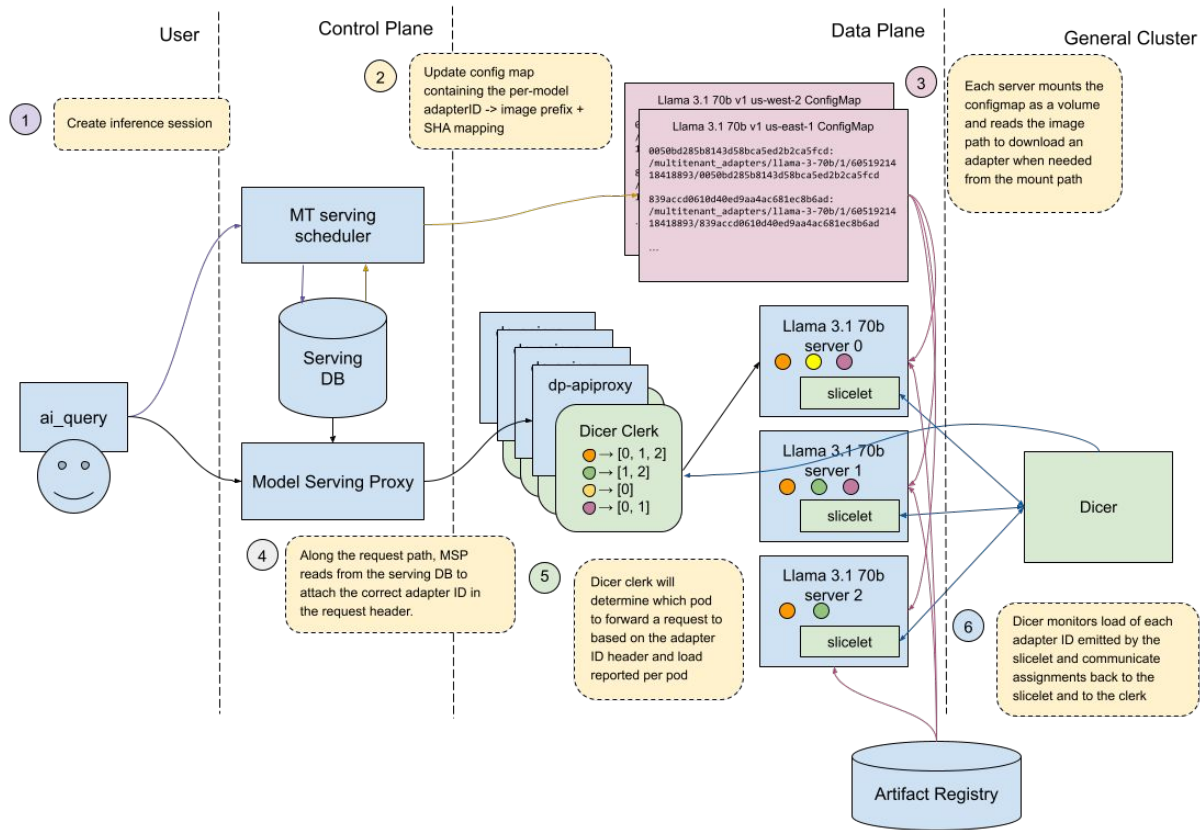
- Fall back to storage layer
- High cache hit rates
- Low Network IO Overhead
- Low Serialization Overhead



# Use Case: Ratelimit and in memory stores



# Use Case: Model Serving - Batch Inference



4x more PEFT adapters loaded

Save 10 mil/year



# The new problems

## Clients in control

- Proxies
- Pod to Pod Forwarding

## Overload Protection

- Interactions with load shedder
- Load and health metrics not perfect
  - Clients often do reachability checks



# Case Study: Unity Catalog

