# Agenda

# Background

- Large language models are getting bigger, a single accelerator node doesn't have the memory to load the entire model, e.g. Llama3-405B, Mistral-123B etc.

- High Bandwidth Memory (HBM) is important for storing intermediate state and enable higher throughput. Thus we need to shard the model on multiple nodes to efficiently serve the LLM.

Single node serving instance

Accelerator node

Model A (70b)

Multi-node serving instance

Model Shard 1

Model Shard 2

Model B (405b)

Accelerator node

Model B shard-1

Tensor parallelism / Pipeline parallelism
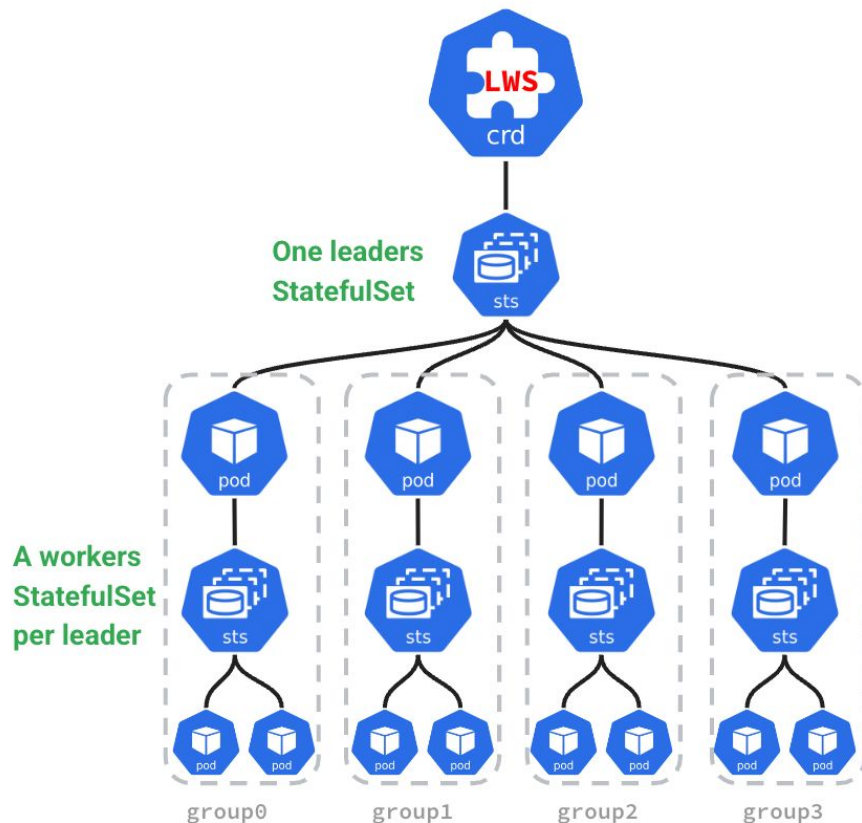
Accelerator node

Model B shard-2

However, there is no automated way to start and replicate a group of pods as a unit.

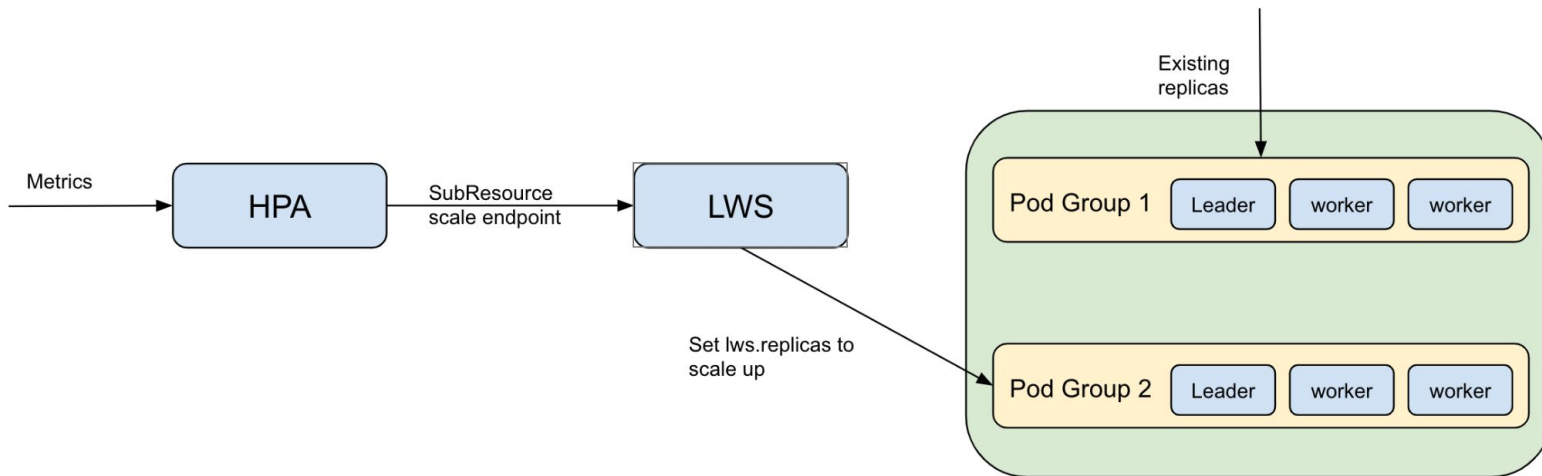Solution: a new API named **LeaderWorkerSet**

# LeaderWorkerSet Overview

- An API for deploying a group of pods as a unit of replication
  - A group is a "Super Pod": a single leader and multiple workers
  - Built on top of StatefulSet and supports most of it semantics

- Features that enables multi-node serving
  - **Pod indices** to assign worker ranks
  - **Dual-template Pods** the leader and the workers have different templates, allows running additional containers on the leader (e.g., the Ray head)
  - **Group startup policy:** create the workers when the leader is ready
  - **Group restart** restarts the pods of a group as a unit when anyone in the group fails
  - **Horizontal-scaling** of groups triggered by HPA
  - **Rolling updates** and rollbacks at the group level
  - **Placement policy** for exclusive 1:1 group to TPU slice placement

# Autoscaling

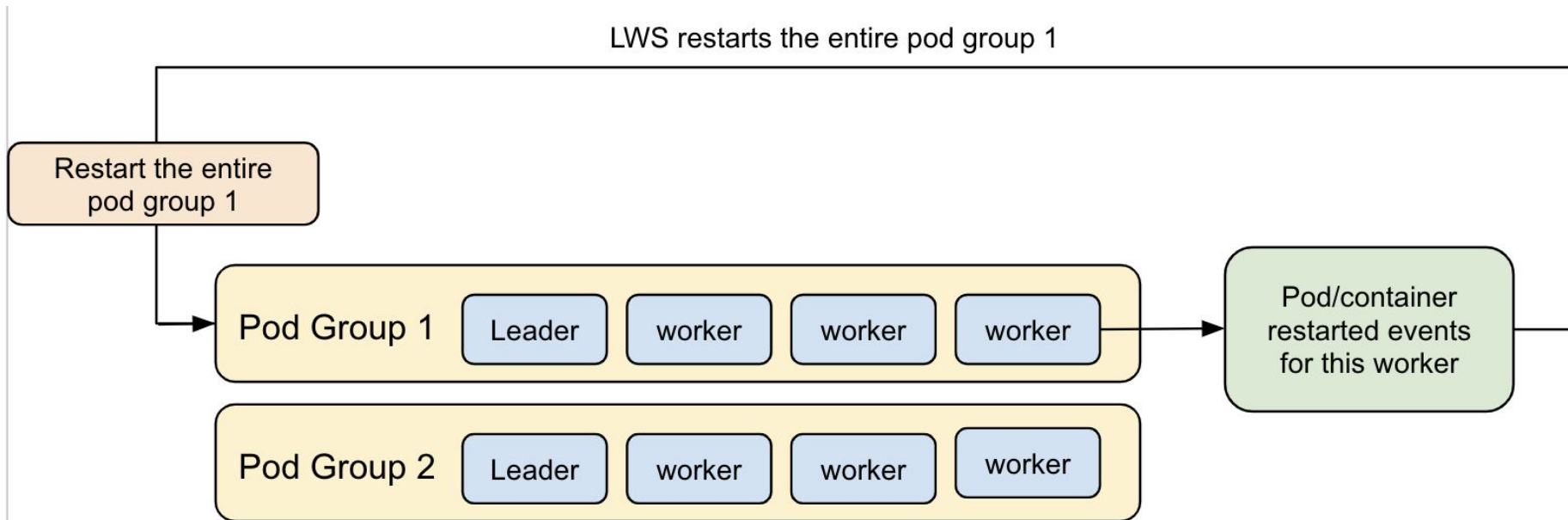LWS exposes a scale subresource to allow HPA to adjust the number of replicas

```
subresources:
  status: {}
  scale:
    specReplicasPath: .spec.replicas
    statusReplicasPath: .status.replicas
    labelSelectorPath: .status.hpaPodSelector
```

# Failure Handling

Any pod/container restarts in a single serving replica will require the entire replica to be restarted. (group restart)

This is because the distributed inference workloads need to be re-initiated when any program on different pod in the same group is restarted.
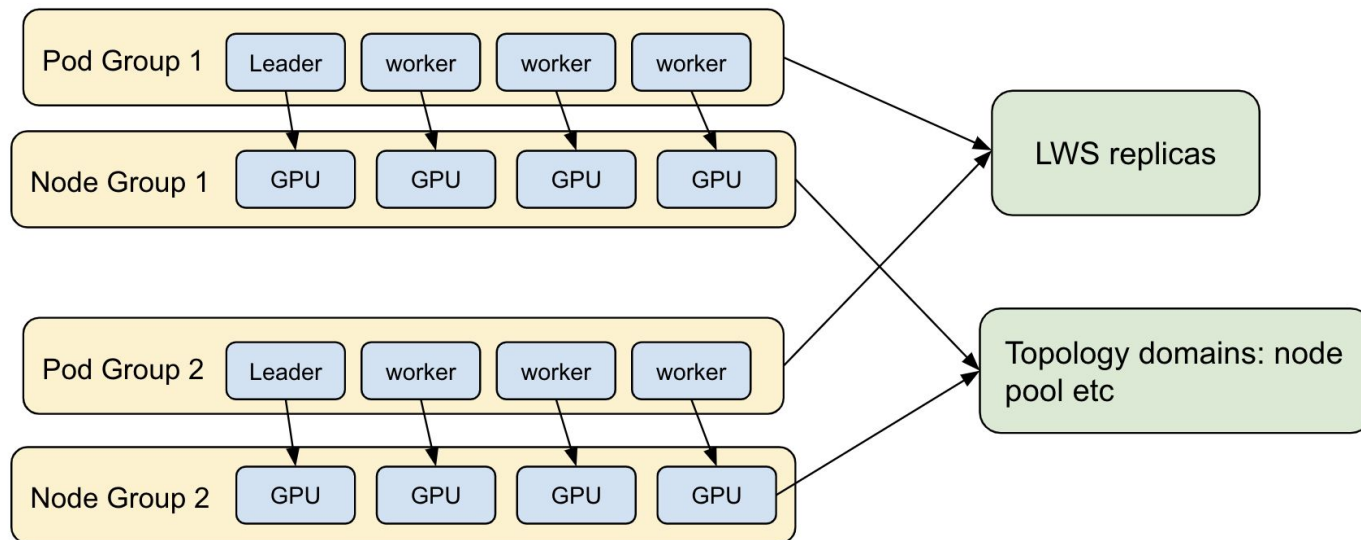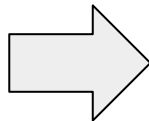
# Topology aware scheduling

One pod group/replica will be exclusively scheduled to one topology domain.

- One pod group is a single multi-host serving instance
- A topology domain here is a group of nodes that share the same node label which categorize them to be the same group

# Topology aware scheduling



```
apiVersion: leaderworkerset.x-k8s.io/v1
kind: LeaderWorkerSet
metadata:
  annotations:
    leaderworkerset.sigs.k8s.io/exclusive-topology:
google.com/node-pool
  name: multi-node-inference
spec:

...
...
```

```
kind: LeaderWorkerSet
  ...
    affinity:
      podAffinity:
        requiredDuringSchedulingIgnoredDuringExecution:
        - labelSelector:
            matchExpressions:
            - key: leaderworkerset.sigs.k8s.io/group-key
              operator: In
              values:
              - multi-node-inference-0   # pod group name
          topologyKey: google.com/node-pool   # the grouping domain
      podAntiAffinity:
        requiredDuringSchedulingIgnoredDuringExecution:
        - labelSelector:
            matchExpressions: # exclude the pods of this job
            - key: leaderworkerset.sigs.k8s.io/group-key
              operator: NotIn
              values:
              - multi-node-inference-1
            - key: leaderworkerset.sigs.k8s.io/group-key
              operator: Exist
          topologyKey: tpu-slice-id
```
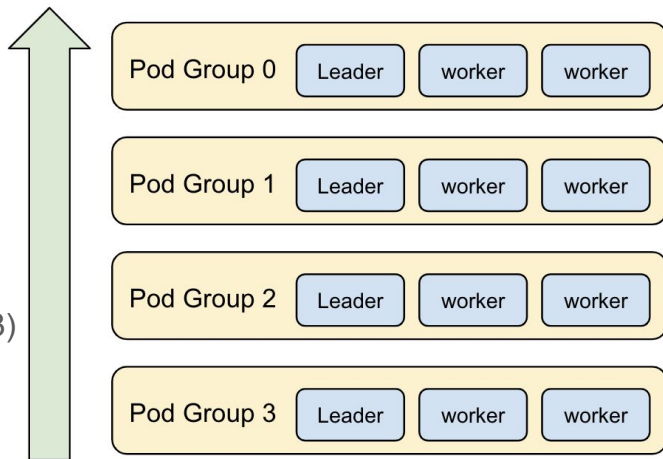
**Collocate the pods of a group**

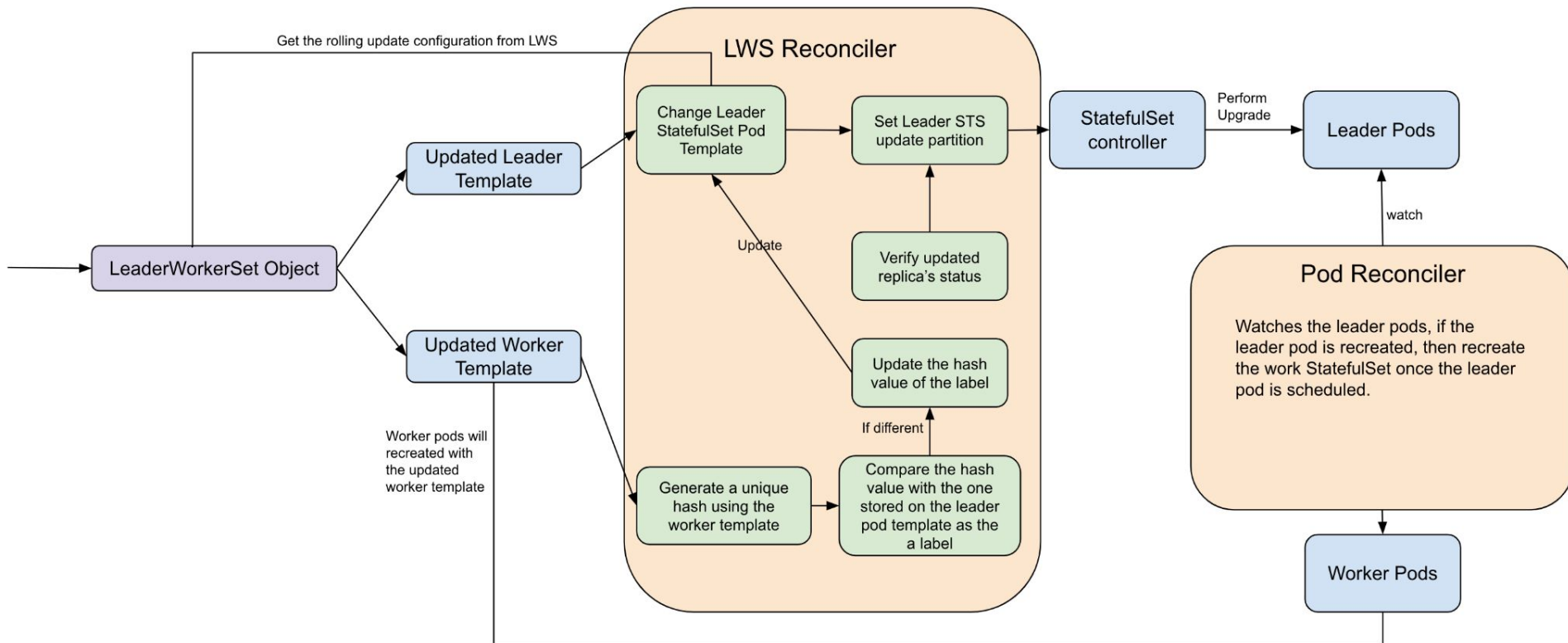**Repels all other pods**

# Rolling update

LWS supports rolling update strategy with
MaxUnavailable and MaxSurge settings

- Each replica is unit for updating
- **Triggering of the upgrade**: a template hash will be generated for the LeaderWorkerTemplate and some API fields, this hash will be stored as a label on the leader statefulset
- **The upgrade process**: LWS relies on the leader statefulset to manage the process of the upgrade. The upgrade will stop when the upgrade replicas are unhealthy and we use the leader Statefulset's Partition field to control this
- **The restart of the single replica**: The leader statefulset will update the leader pod through recreate and the worker statefulset coupled with the leader pod will be recreated as well through owner reference

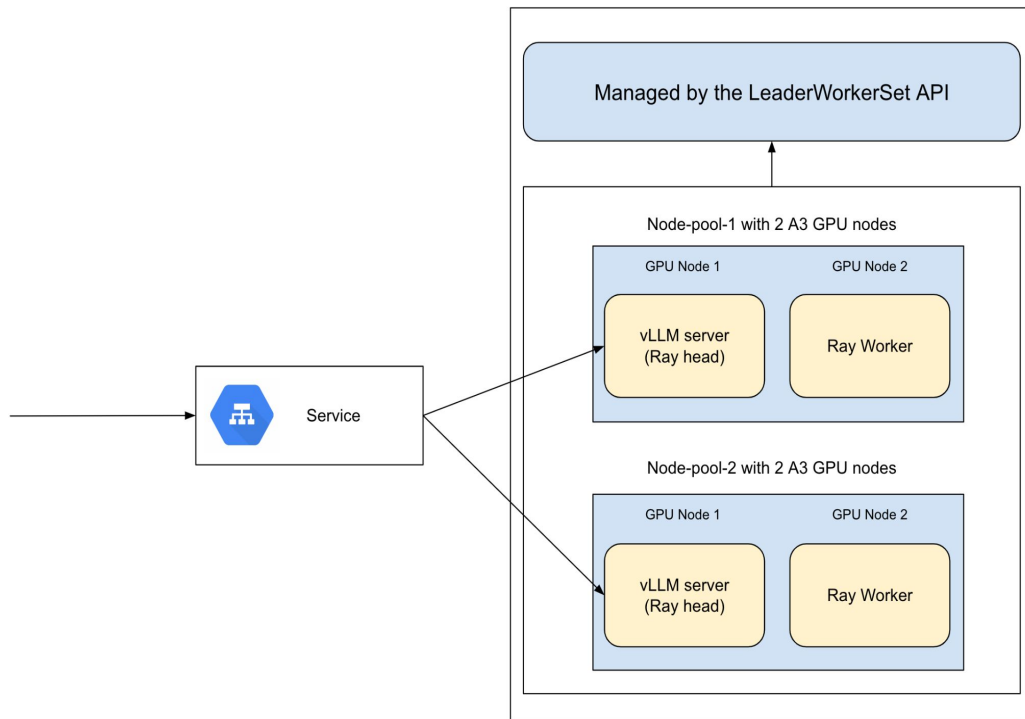The rolling upgrade progresses from the highest group index (group 3) to the lowest (group 0)

| Pod Group 0 | Leader | worker | worker |
|---|---|---|---|

| Pod Group 1 | Leader | worker | worker |
|---|---|---|---|

| Pod Group 2 | Leader | worker | worker |
|---|---|---|---|

| Pod Group 3 | Leader | worker | worker |
|---|---|---|---|

# Rolling update

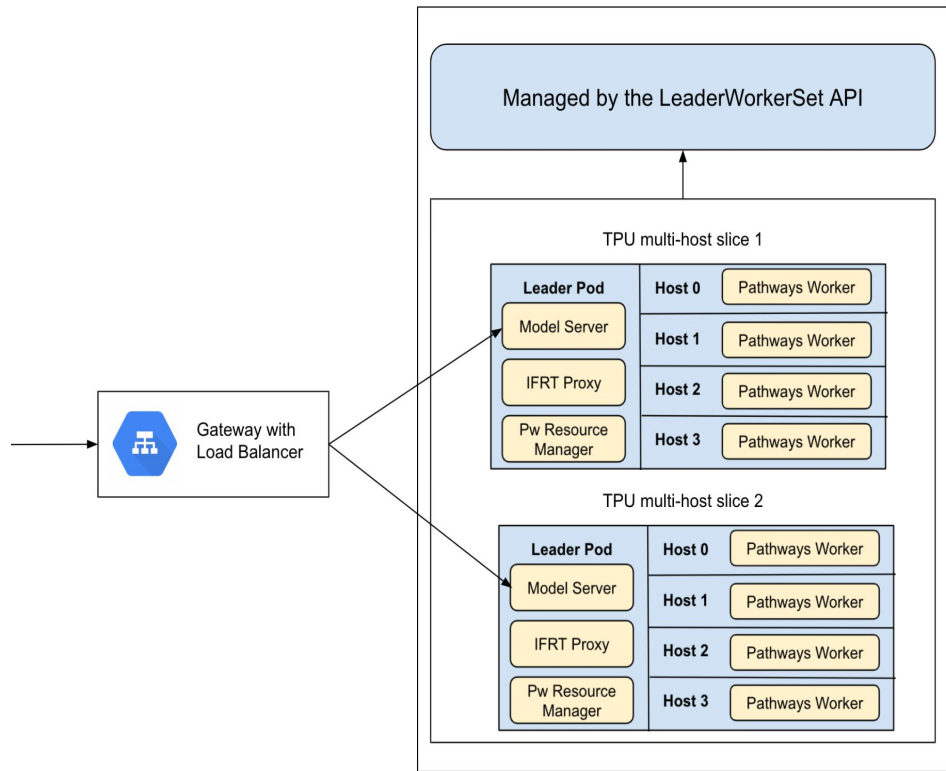# Multi-node inference on GPU

Infrastructure & runtime setup:

- Two A3 nodes with 8 H100 GPU chips using compact placement (co-locating to minimize cross node communication latency)
- vLLM + Ray as the ML runtime
- Llama3-405B model with tensor parallelism = 8 (# GPUs per node) and pipeline parallelism = 2 (# of nodes)
- all_gather happens within the node over nvLink; cross node data transfer over datacenter network links (or RDMA if available)

# Multi-node inference on TPU

Infrastructure & runtime setup:

- One TPU slice of v5e-16 (5 pods per replica running on 4 node)
- JetStream (TPU inference server) + Pathways (distributed compute substrate)
- Llama3-405B model with tensor parallelism
- Cross hosts communication through InterChip interconnect (ICI)

# Multi-node inference on TPU

# Disaggregated Serving

The processing of a LLM inference can be split into two phases: **Prefill & Decode**.

Prefill is more **compute bound** while decode is more **memory bound**. The two phases will interfere with each other running on one server which causes resource inefficiency and unpredictable latency.

To address this problem, disaggregated serving splits prefill & decode stages to run separate servers. This brings the following benefits:
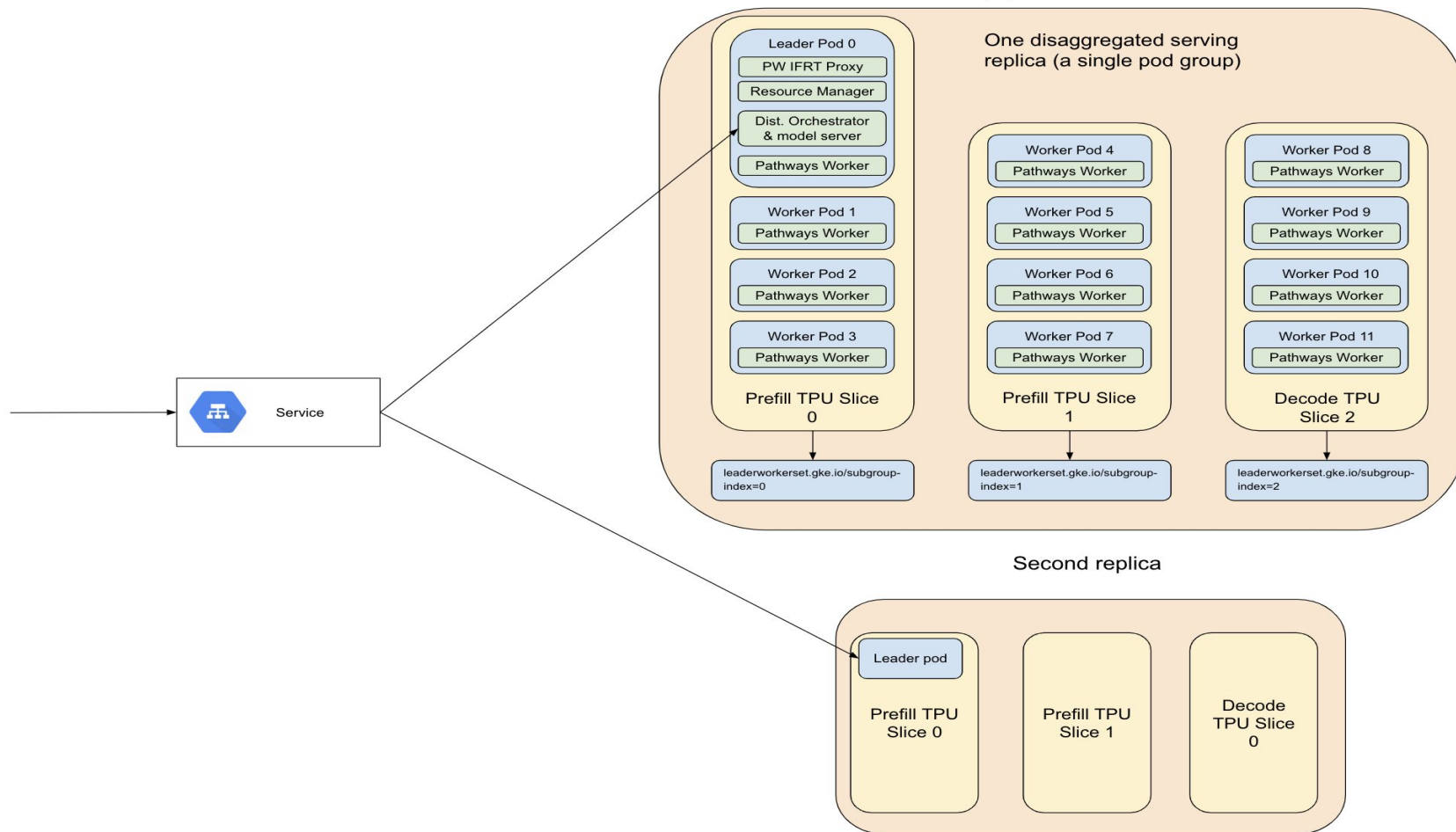
- Reduce prefill/decode interference. Prefill usually work at batch size 1 and decode phase has much larger batch size, whenever there is a request coming in, the decode loop will be interrupted to prefill the new request
- Reduce time to first token latency.
- Enable heterogeneous VM types for prefill/decode that better matches the compute characteristics of each phase
- Better overall resource efficiency and lower cost

LeaderWorkerSet helps orchestrate this type of multi-node serving paradigm

# Getting Involved

**LeaderWorkerSet**
- https://github.com/kubernetes-sigs/lws
- A Serving WG sponsored project


**Serving WG:**
- Biweekly meetings,  Wednesdays at 9:00 PT (Pacific Time)
- slack.k8s.io #wg-serving
- wg-serving@k8s.io
- git.k8s.io/community/wg-serving

# Implementation details