



**KubeCon**



**CloudNativeCon**

**Europe 2019**



KubeCon



CloudNativeCon

Europe 2019

# gRPC load balancing and Service Mesh

Vishal Powar (GitHub @vishalpowar)  
Google

# Agenda



KubeCon



CloudNativeCon

Europe 2019

- Load balancing in gRPC
- Centralized balancing
- Load balancing at scale
- Service mesh
- Demo

# Why load balancing?



KubeCon



CloudNativeCon

Europe 2019

- Load balancing is a mechanism to
  - Improve throughput of a service.
  - Improve service availability and reliability.
- Load balancing should help client pick service endpoints
  - Based on client requirements (latency, #connections).
  - Based on endpoint requirements (isolation, #connections).

# Client-side load balancing (gRPC)



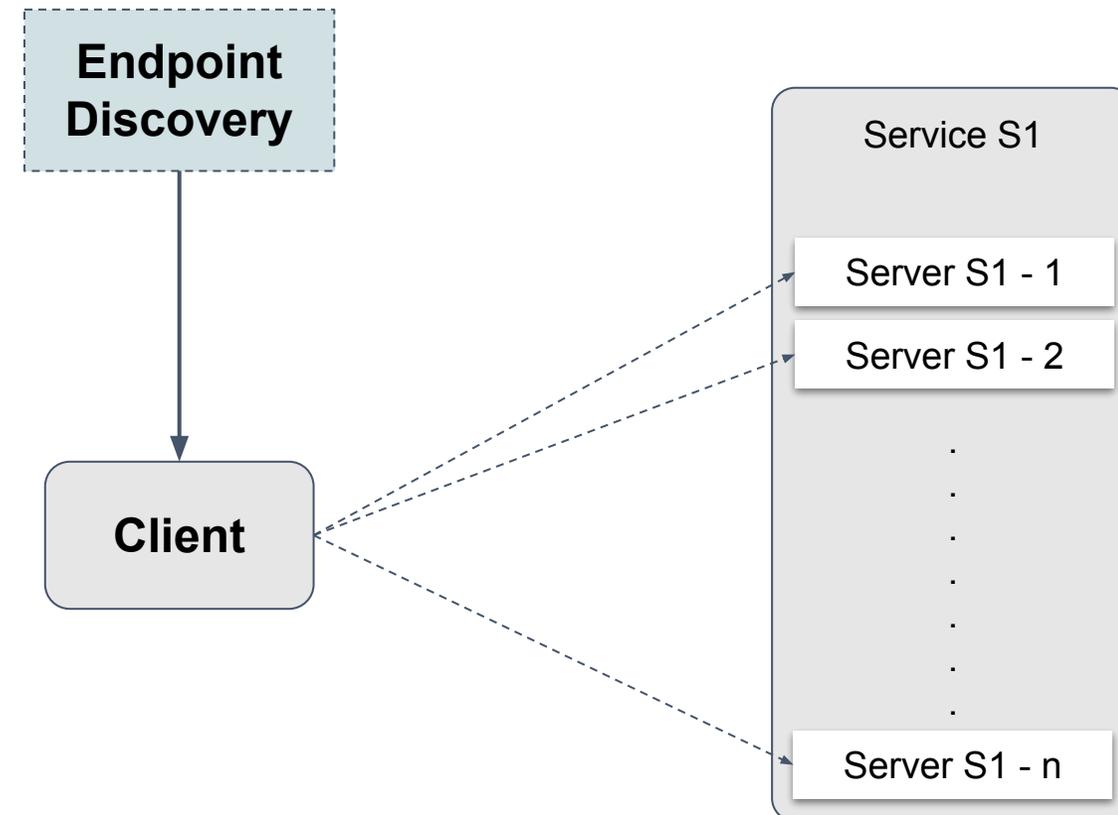
KubeCon



CloudNativeCon

Europe 2019

- Client decides which endpoints to connect and send request
  - e.g
    - Round Robin
    - Pick-First
- Client can also connect to subset of endpoints.



# Client-side load balancing (gRPC)



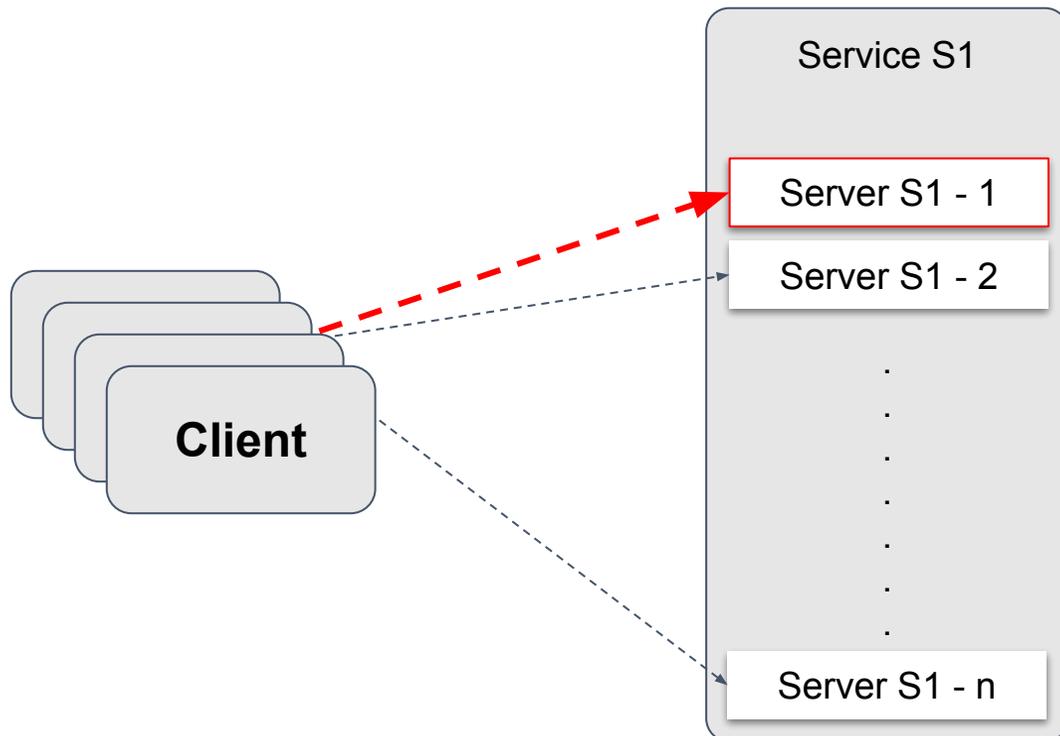
KubeCon



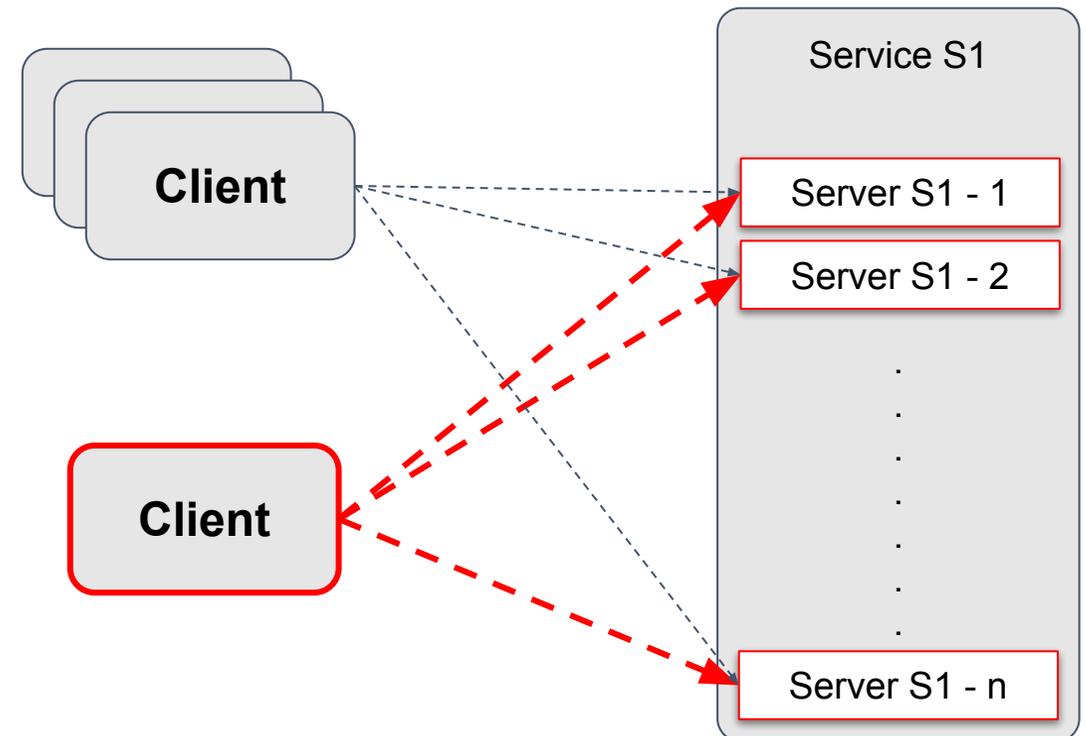
CloudNativeCon

Europe 2019

- Pick first - Server Overload



- Round robin - No isolation



# Centralized load balancing



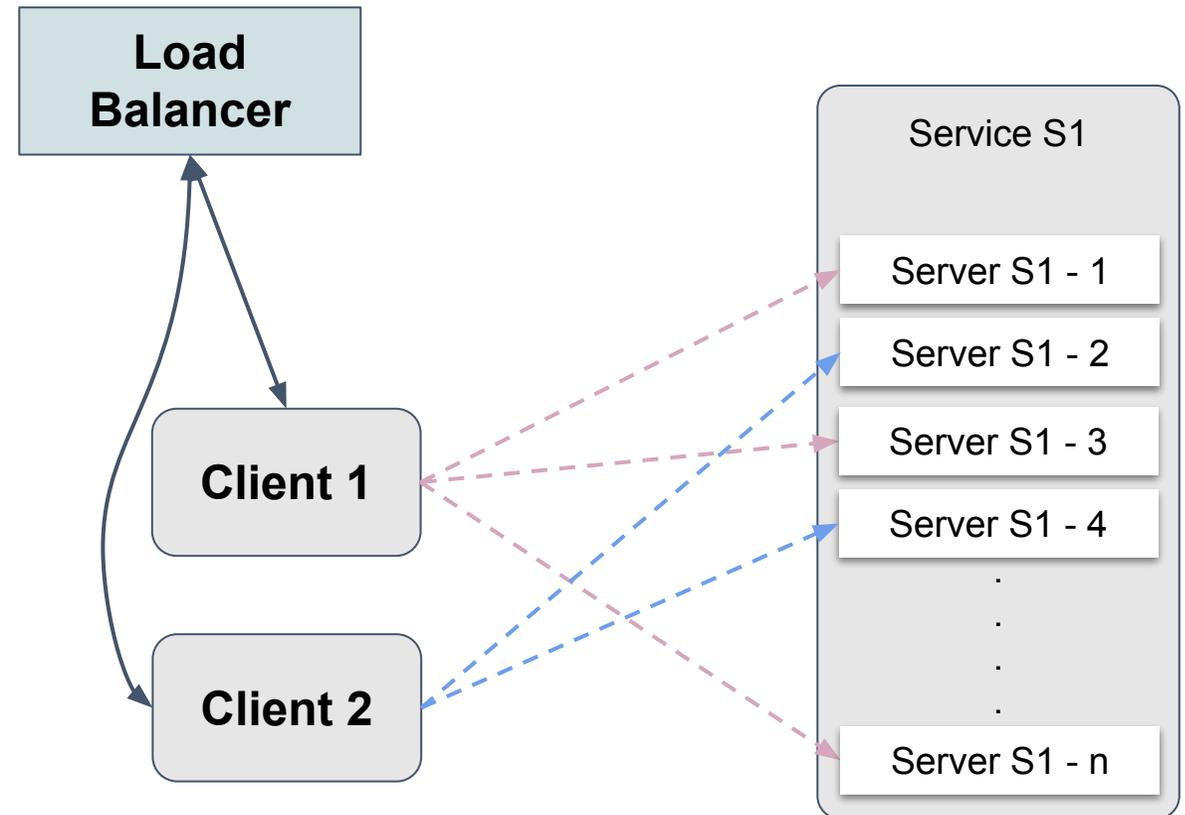
KubeCon



CloudNativeCon

Europe 2019

- Load balancer takes global decision to protect endpoints from client local decision.



# gRPC LB



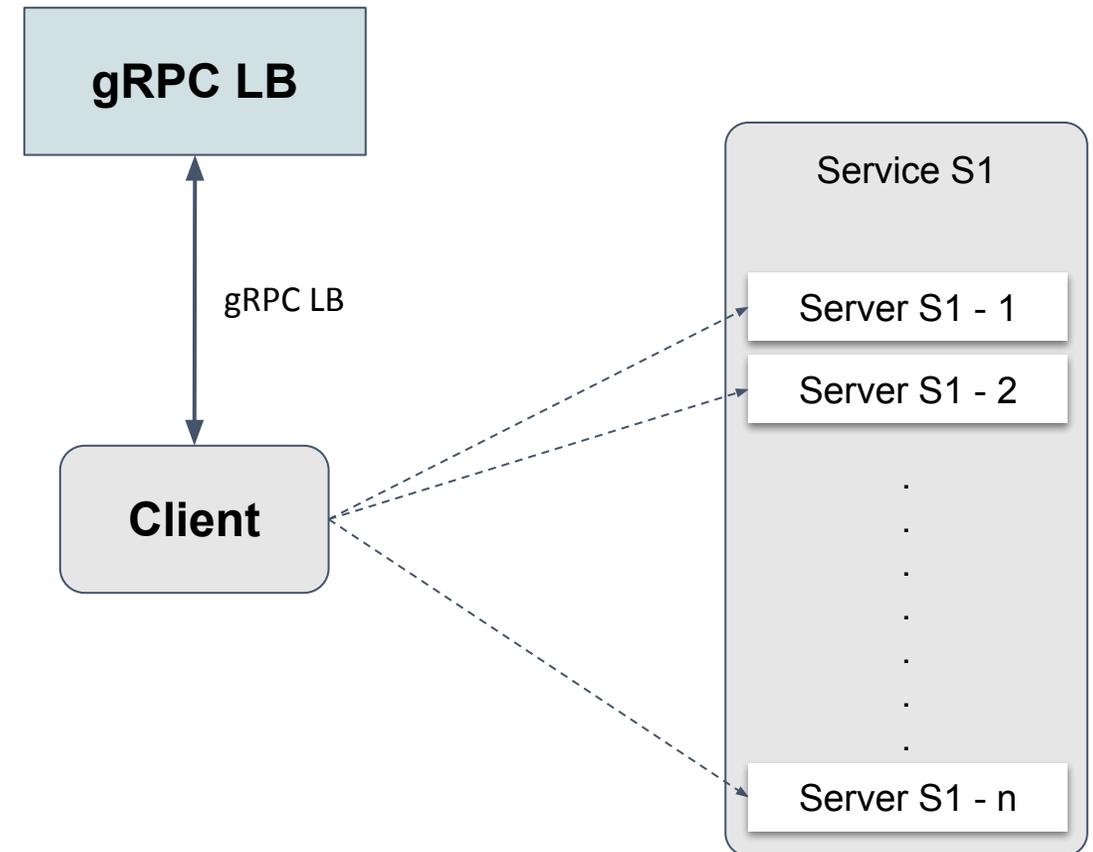
KubeCon



CloudNativeCon

Europe 2019

- Look-aside load balancing with gRPC LB protocol.
- Balancer provides list of endpoints to use.
  - The list encodes weight information that clients use for making request.



# Proxy load balancing (envoy)



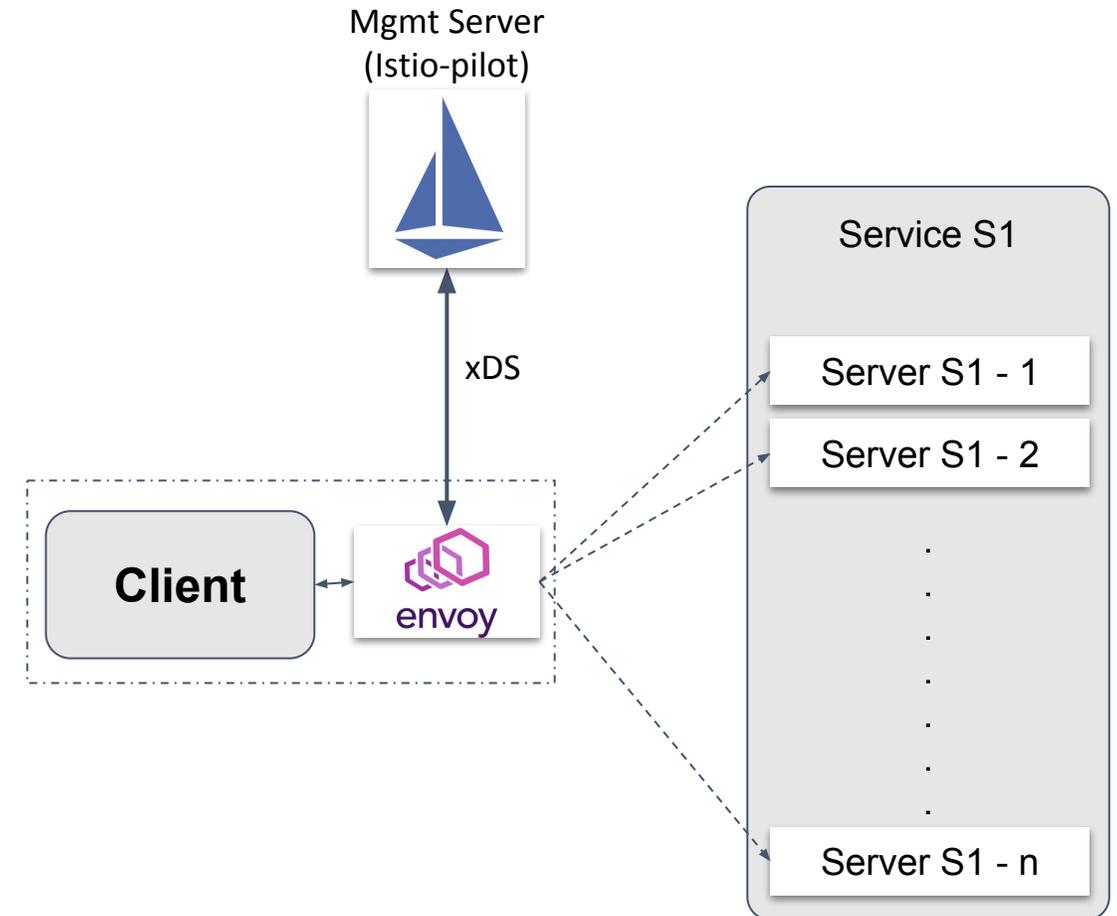
KubeCon



CloudNativeCon

Europe 2019

- Middle proxy or sidecar deployment.
- xDS server provides endpoints and load balancing configuration for envoy.



# Deployment considerations



KubeCon



CloudNativeCon

Europe 2019

- Deployments are heterogenous with endpoints differing in
  - e.g. capacity, location
- Individual endpoints health and capacity can change.
  - e.g. service upgrade, hardware failure.
- For better balancing, central load balancer needs to know about all of this.

# Informed balancing - Capacity



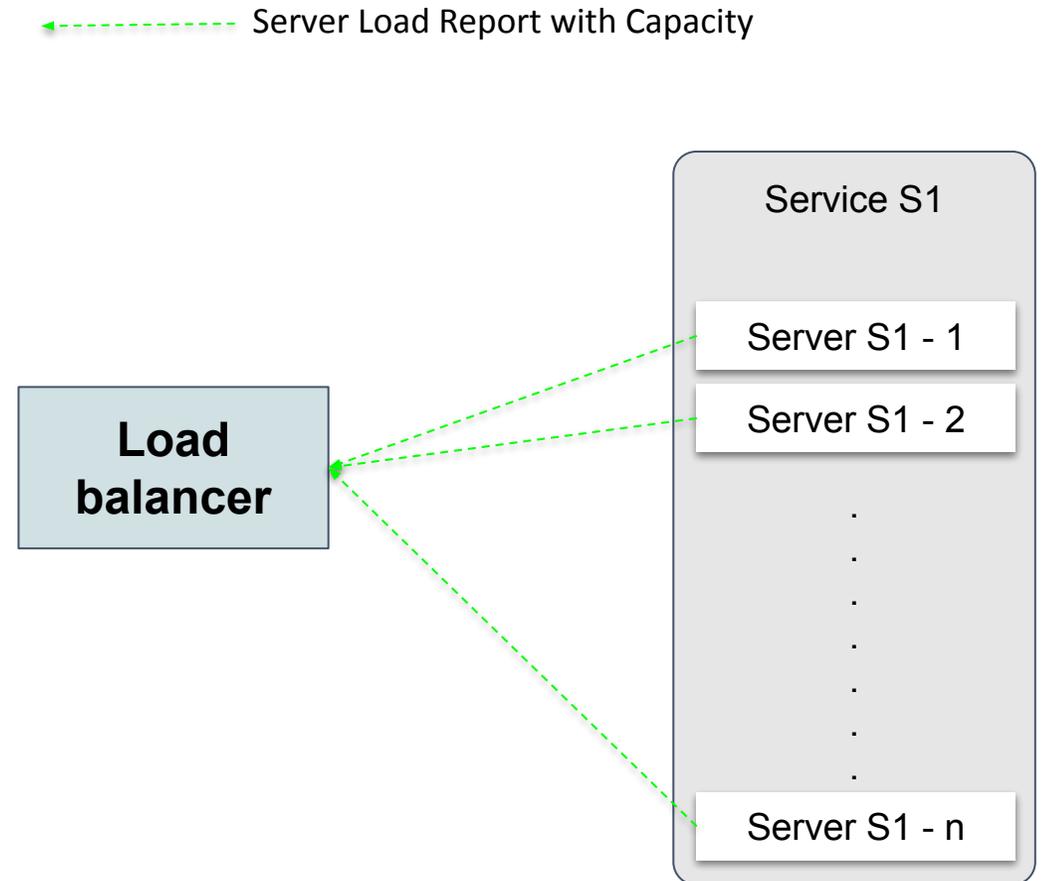
KubeCon



CloudNativeCon

Europe 2019

- Health can be determined by connecting to the endpoint.
- Capacity is service specific, and can be configured or reported by endpoints.
  - e.g compute, memory.



# Informed balancing - Locality



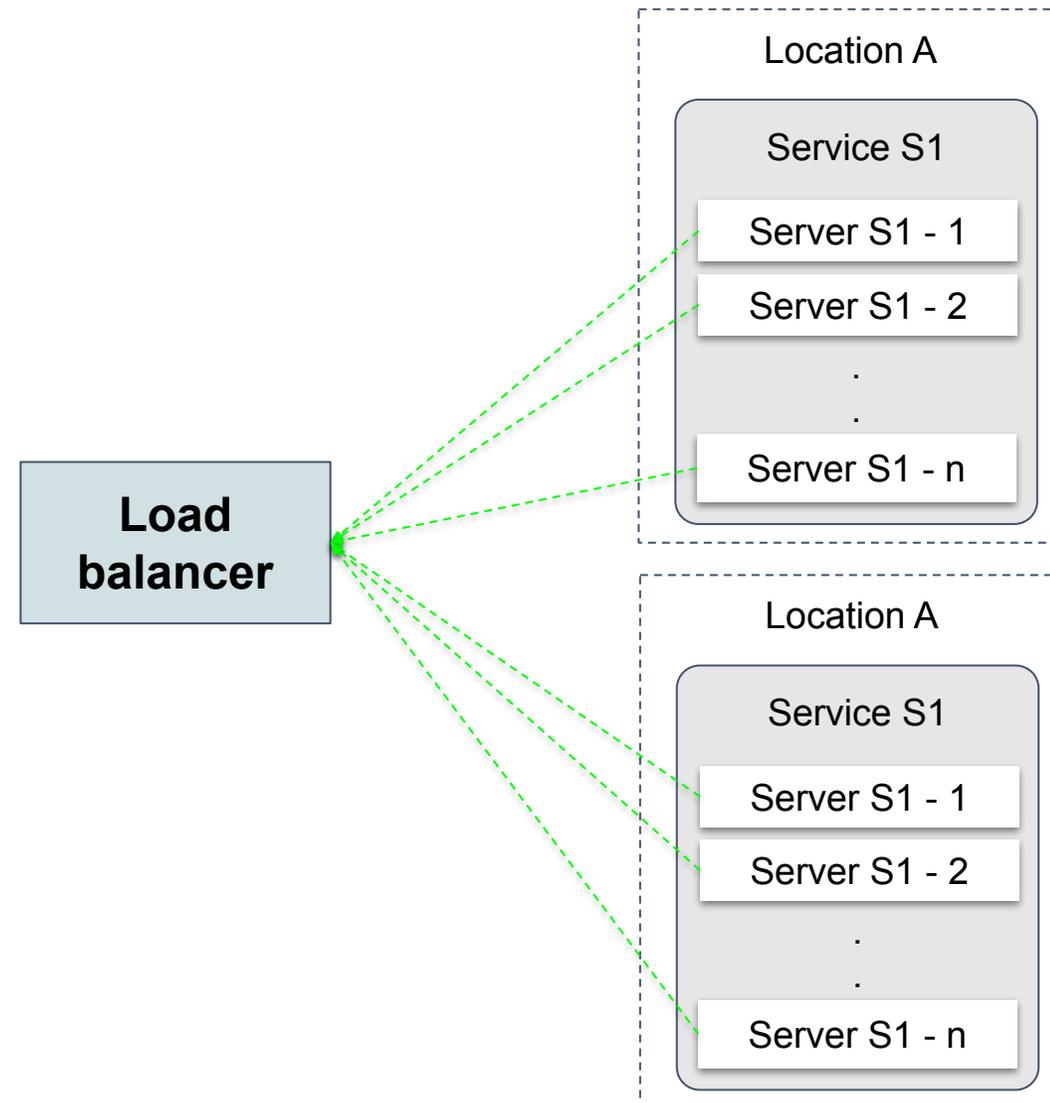
KubeCon



CloudNativeCon

Europe 2019

- Routing requests closer to the client has advantages.
- Both endpoints and clients locality needs to be known by the load balancer.
- Locality capacity can be used for balancing decisions.



# Load balancing at scale



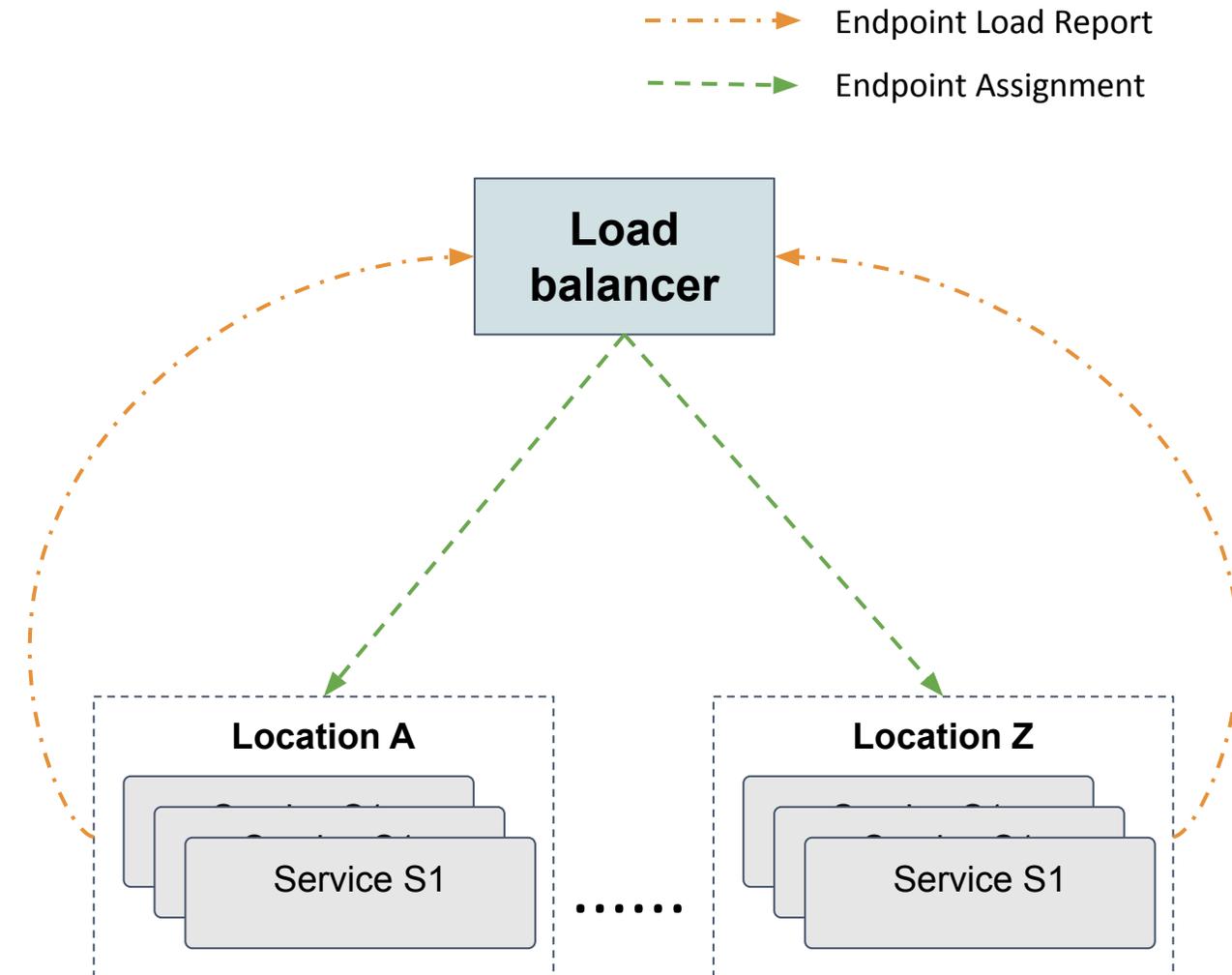
KubeCon



CloudNativeCon

Europe 2019

- Collecting information and making global decisions for each client is expensive.
- Change propagation to clients is slower (~ seconds).



# Load balancing at scale



KubeCon



CloudNativeCon

Europe 2019

- Take global decisions based on locality capacity and proximity.
  - consider client load on each locality.
  - provide enough information for clients to react quickly.
- Have clients take local decisions based on most recent information.

# Service mesh and load balancing

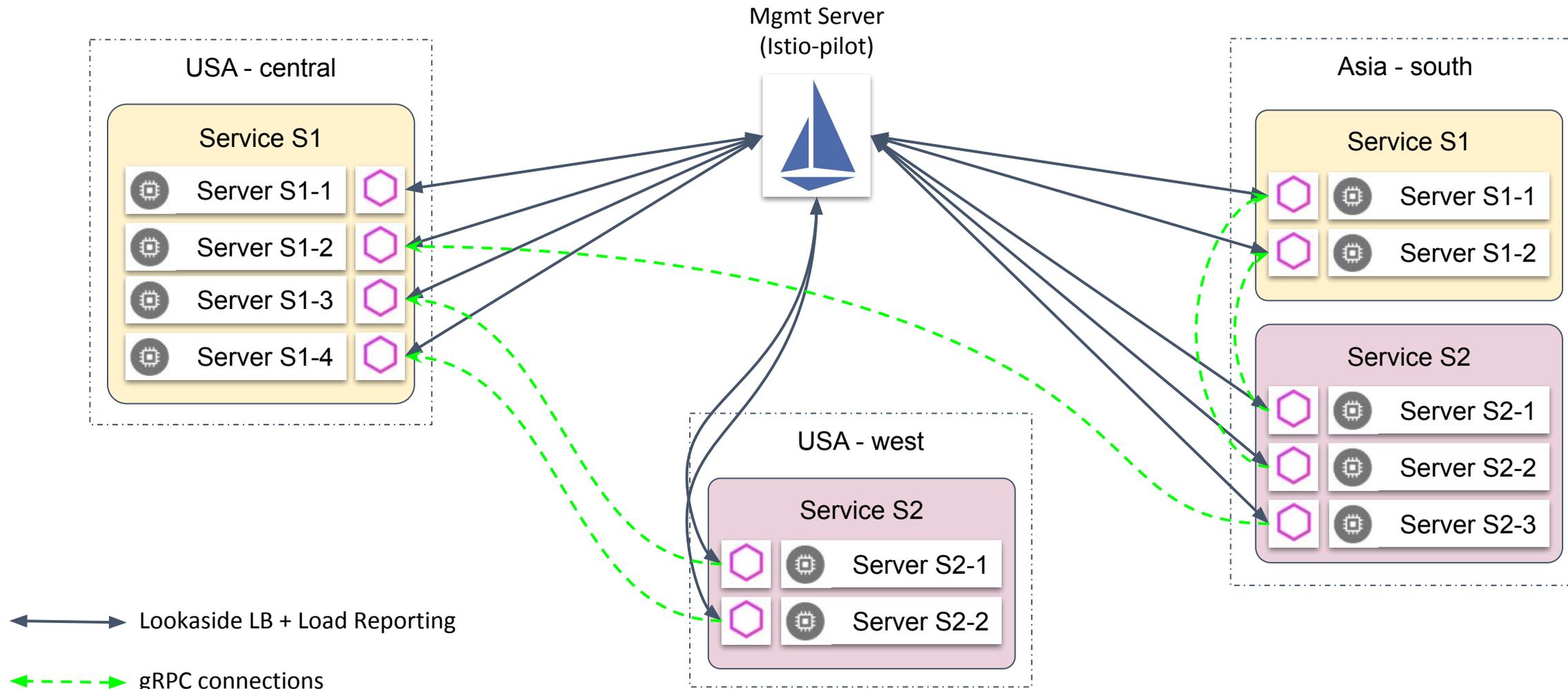


KubeCon



CloudNativeCon

Europe 2019



# gRPC in service mesh



KubeCon



CloudNativeCon

Europe 2019

- xDS provides construct to achieve Load balancing flexibility.
- Information and controls are available for clients side balancing.
  - Weights at Locality and Endpoints
  - Proximity information at locality.

# Demo

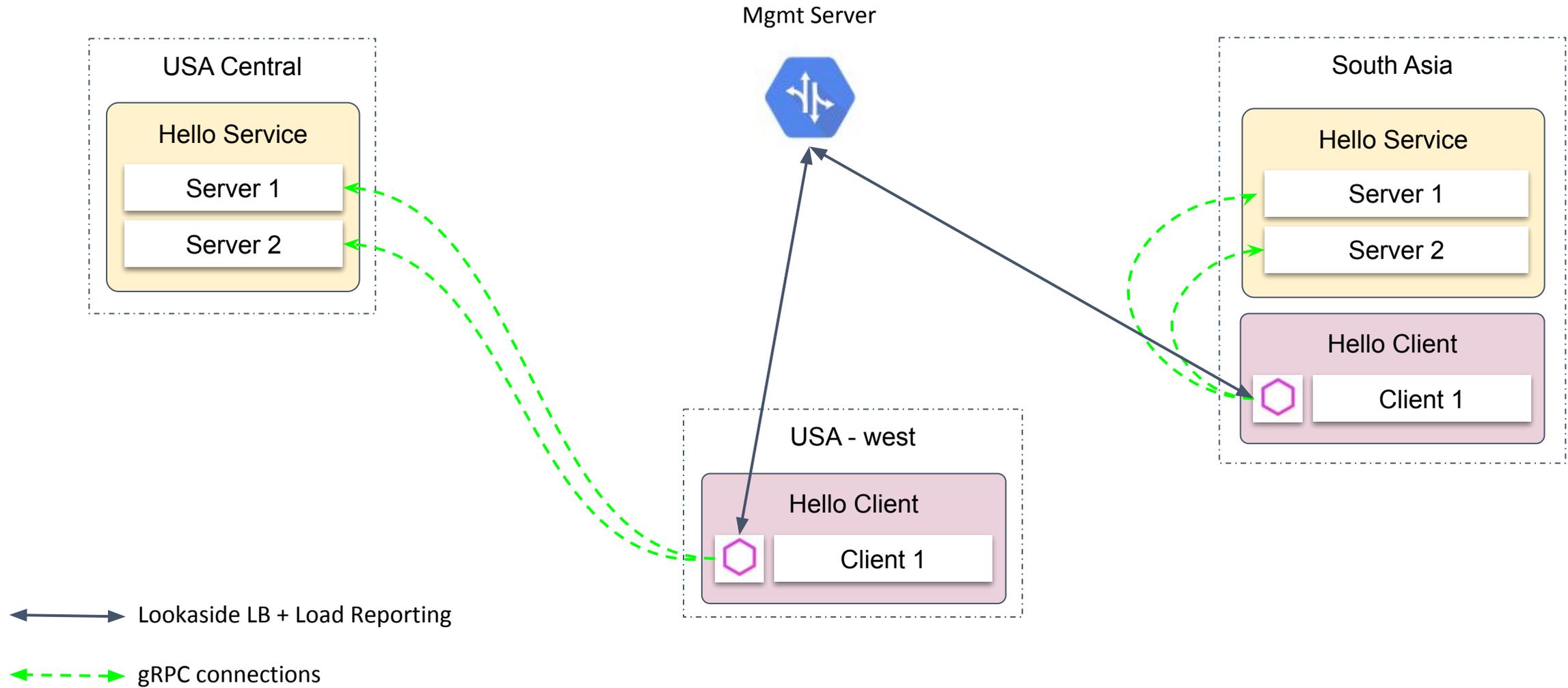


KubeCon



CloudNativeCon

Europe 2019



# Thank You



KubeCon



CloudNativeCon

Europe 2019

