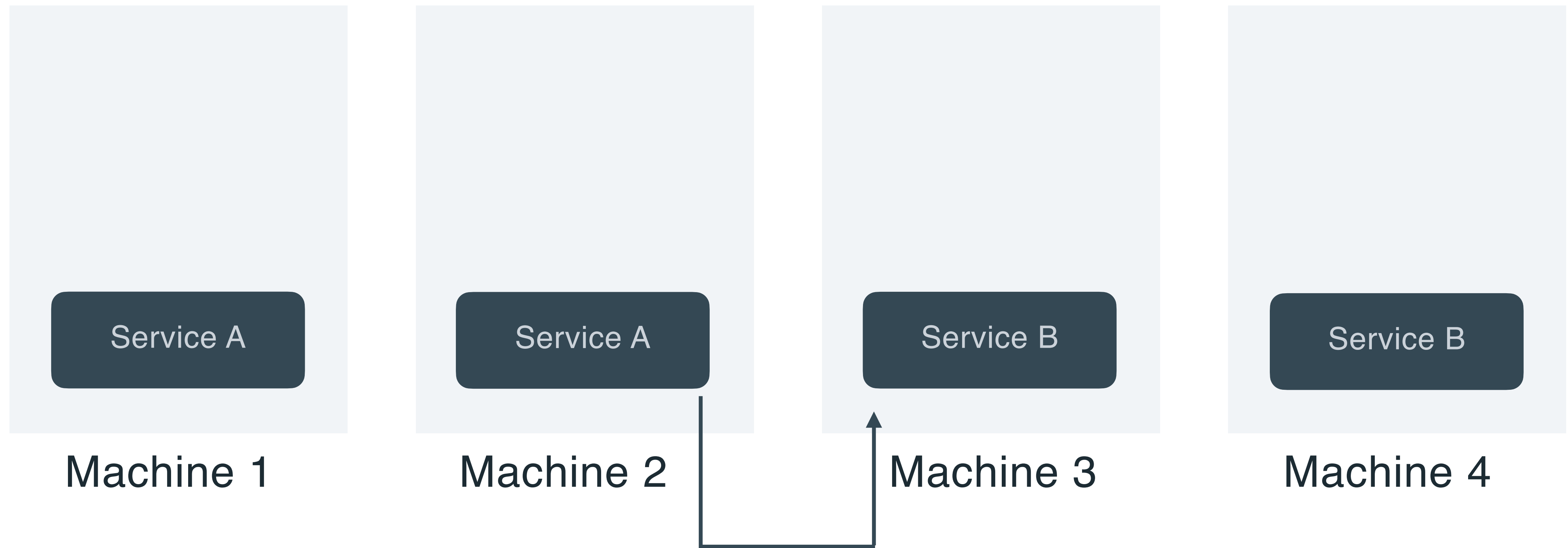# ServiceRouter

**HYPERSCALE AND MINIMAL COST SERVICE MESH AT META**

Harshit Saokar [1]

Soteris Demetriou [1,2]

Nick Magerko [1]

Max Kontorovich [1]

Josh Kirstein [1]

Margot Leibold [1]

Dimitrios Skarlatos [1,3]

Hitesh Khandelwal [1]

Chunqiang Tang [1]

[1] Meta

[2] Imperial College London

[3] Carnegie Mellon University Computer Science Department

# 01 Background & Motivation

Machine 1              Machine 2              Machine 3              Machine 4

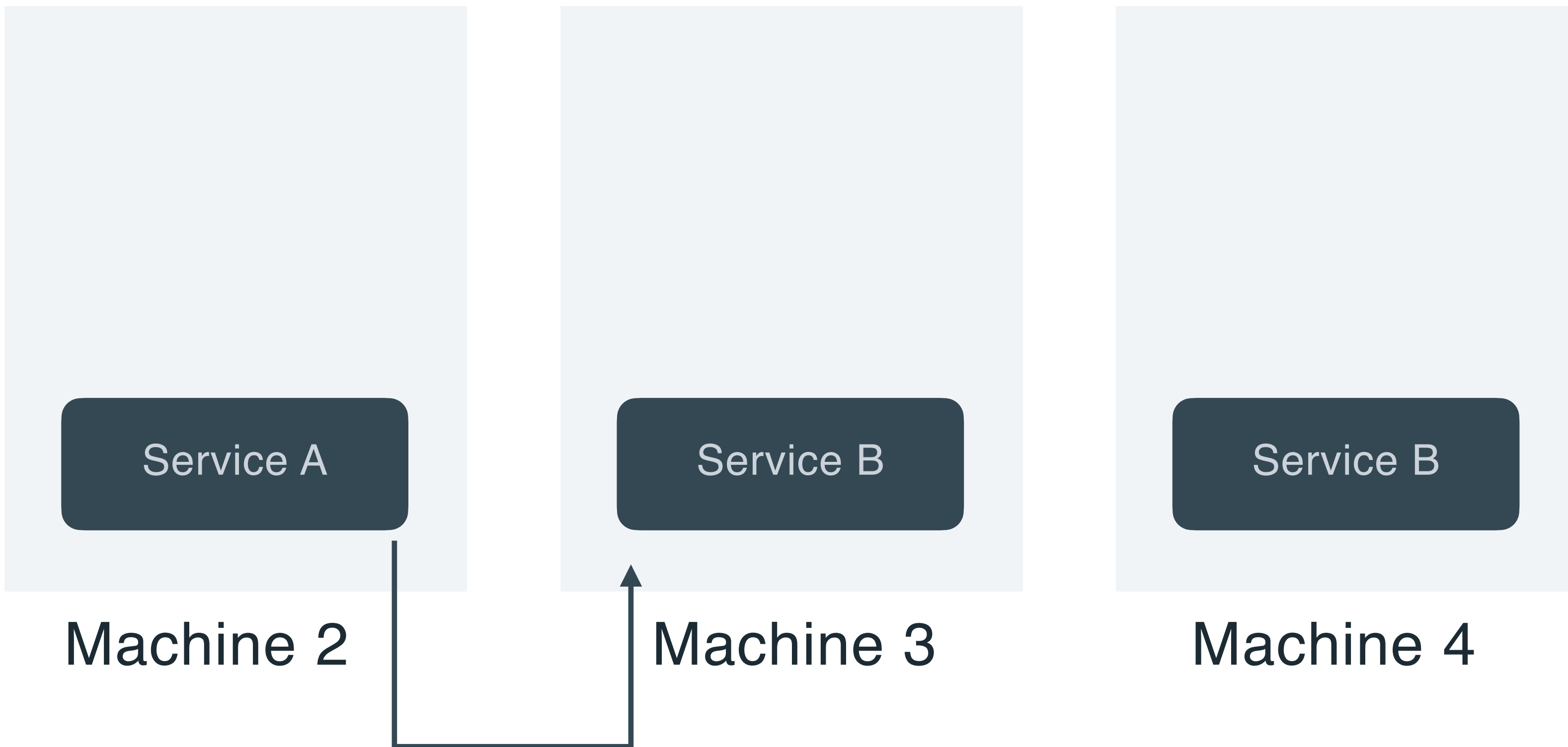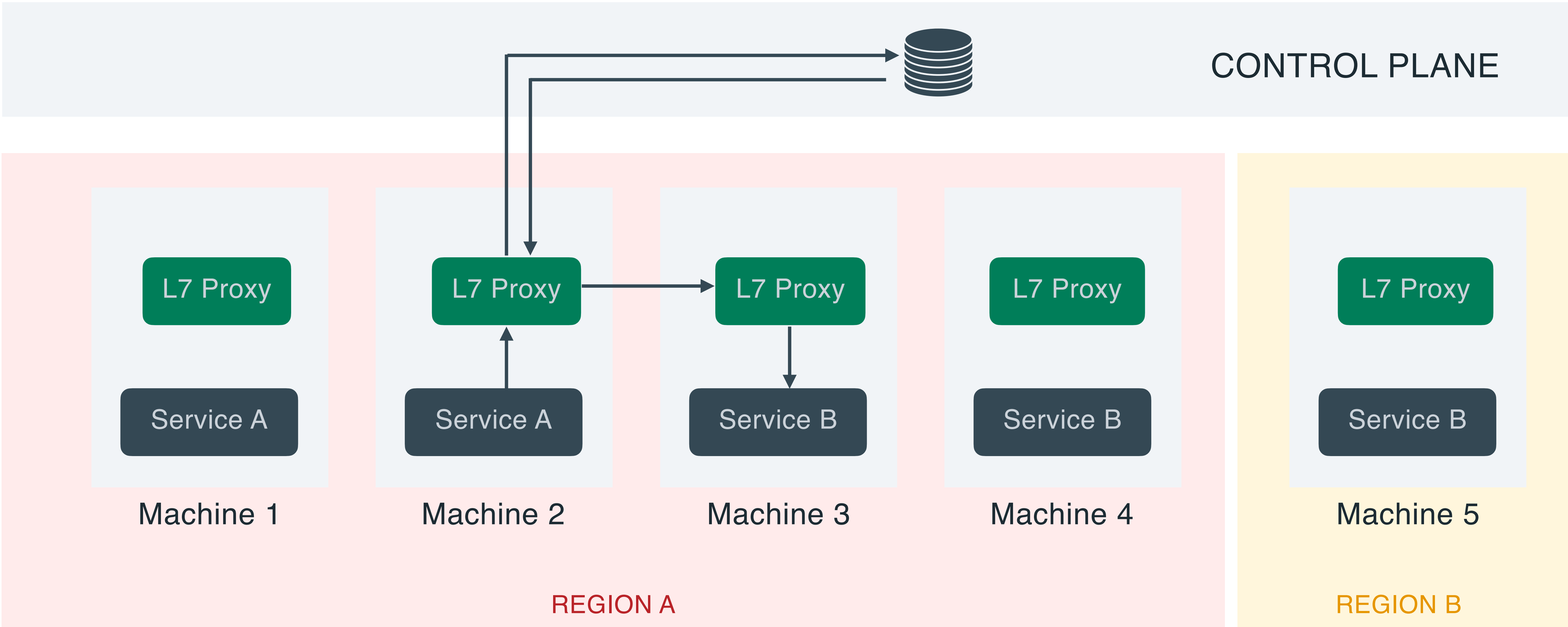Huye et al.  Lifting the veil on Meta's microservice architecture: Analyses of topology and request workflows. USENIX ATC '23

# RPC Frameworks

- No Advanced Load Balancing
- Need external support for service discovery
- Examples: gRPC, Thrift

Service A

Service B
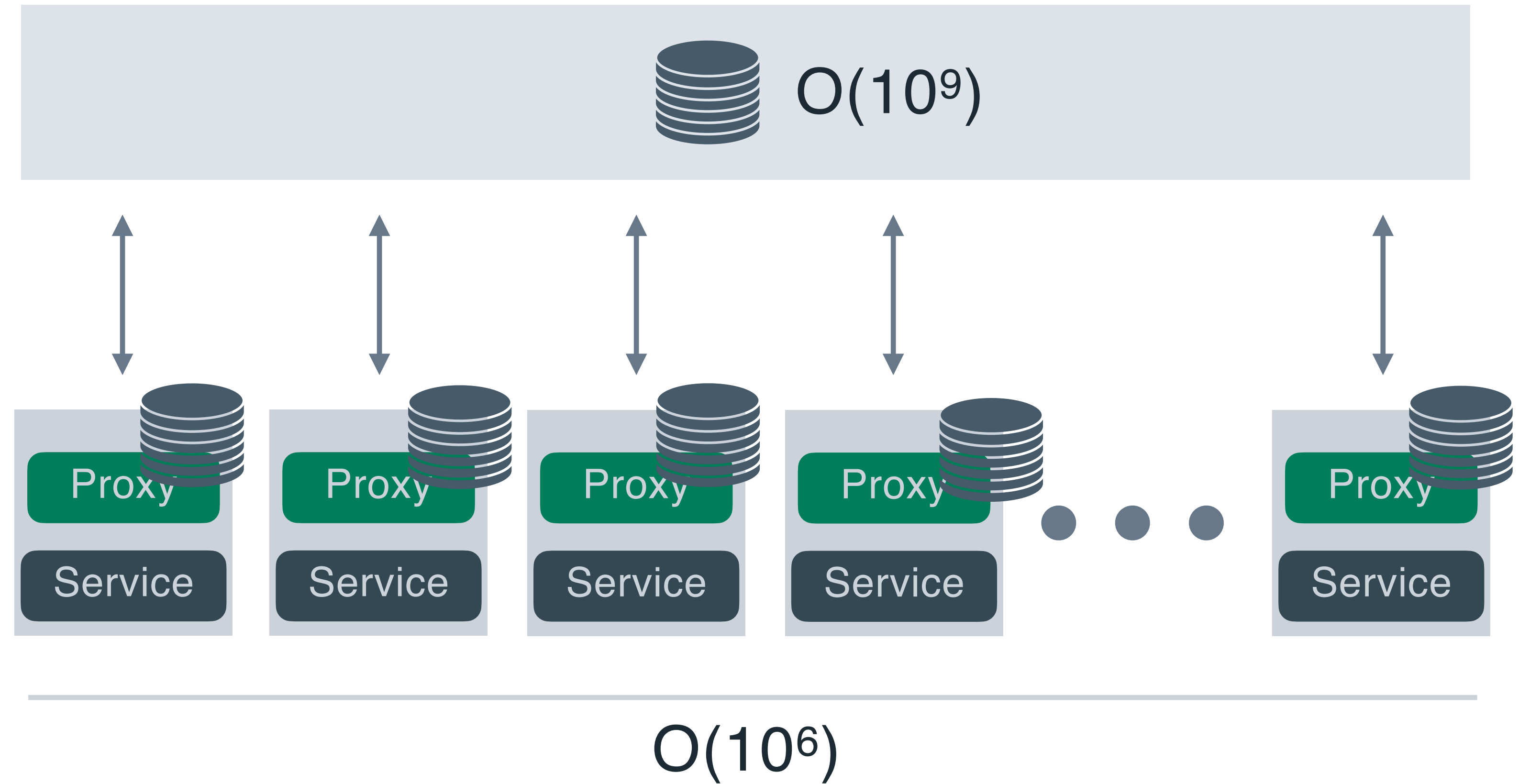
Service B

Machine 2          Machine 3          Machine 4

# Service Mesh Challenges

- **[SCALABILITY]** How can we scale service discovery to $O(10^6)$ clients and proxies?



$O(10^9)$

Proxy · Service · Proxy · Service · Proxy · Service · Proxy · Service · Proxy · Service

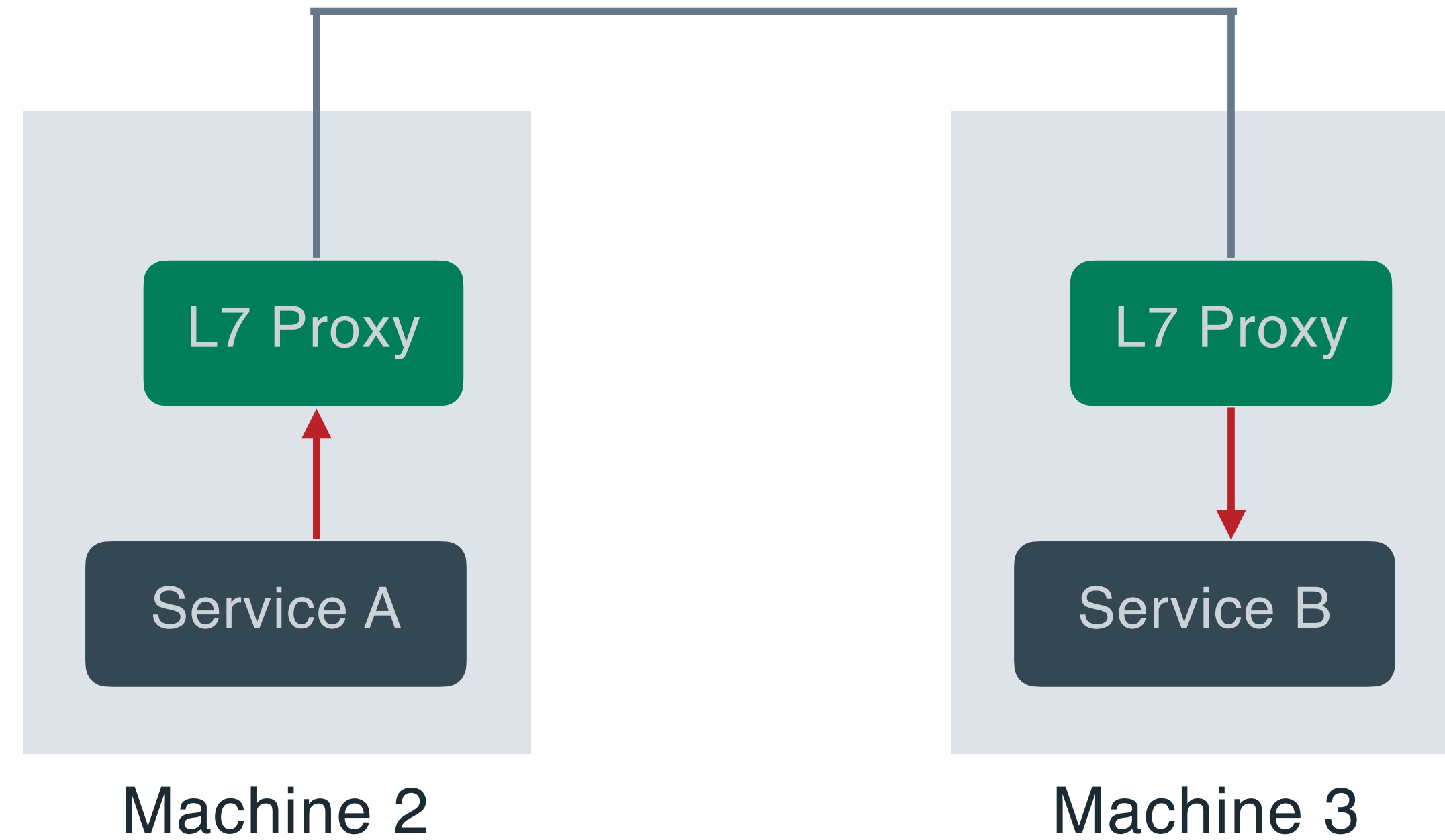$O(10^6)$

# Service Mesh Challenges

- **[SCALABILITY]** How can we scale service discovery to $O(10^6)$ clients and proxies?

- **[HW COST]** How to minimize HW cost?

Istio: 0.35vCPU for $O(10^3)$ rps

**1,750,000 AWS t4g.small VMs for 10B rps**

Proxy

Service

Proxy

Service

Proxy

Service

# $O(10^9)$ RPS

Proxy

Service

Proxy

Service

# Service Mesh Challenges

- [SCALABILITY] How can we scale service discovery to $O(10^6)$ clients and proxies?

- [HW COST] How to minimize HW cost?

- **[RPC LATENCY & LB]** How to simultaneously minimize RPC latency and load balance across geo-distributed hosts?
  - Sidecars add extra latency



Machine 2                    Machine 3

Zhu et al show that Istio

- increases the latency by **185%**

  Zhu et al. Dissecting Service Mesh Overheads. In *arXiv preprint arXiv:2207.00592*, 2022.
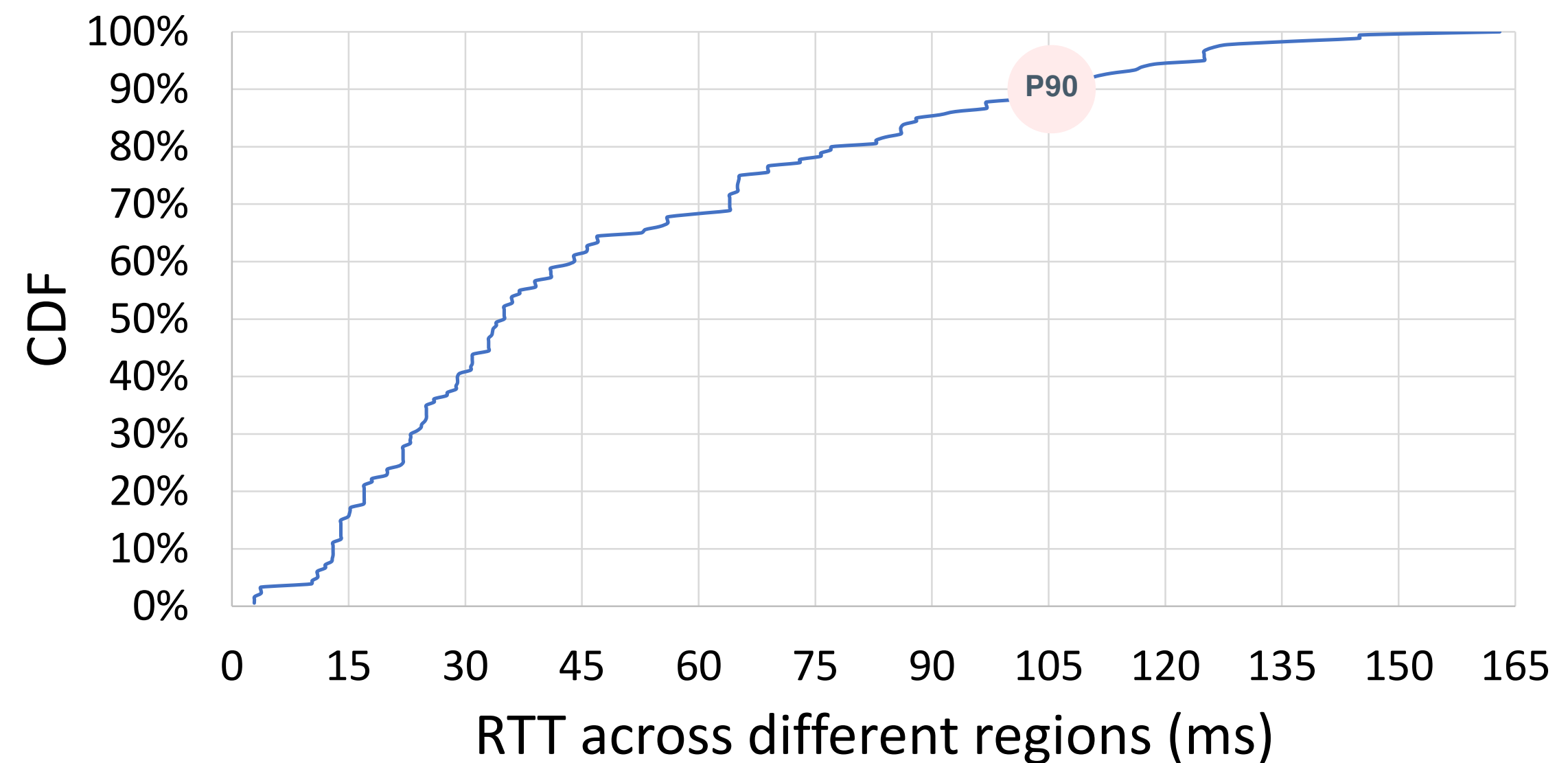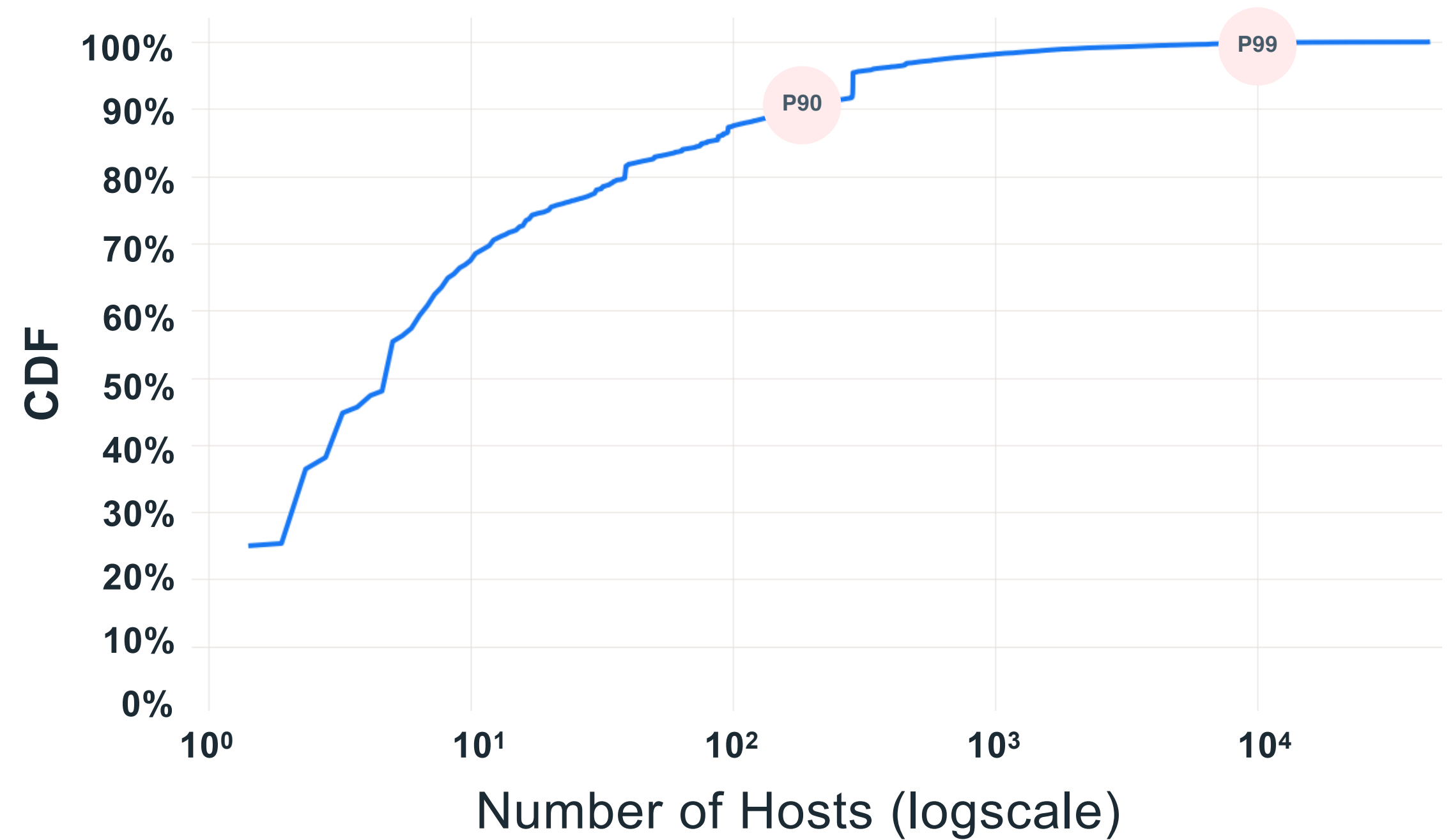
mRPC shows that a sidecar approach:

- increases P99 RPC latency by **180%**

  Chen, et al. Remote procedure call as a managed system service. *NSDI '23*

# Service Mesh Challenges

- [SCALABILITY] How can we scale service discovery to $O(10^6)$ clients and proxies?

- [HW COST] How to minimize HW cost?

- **[RPC LATENCY & LB]** How to simultaneously minimize RPC latency and load balance across geo-distributed hosts?

  - Sidecars add extra latency

  - $O(10\text{-}10^4)$ hosts per service
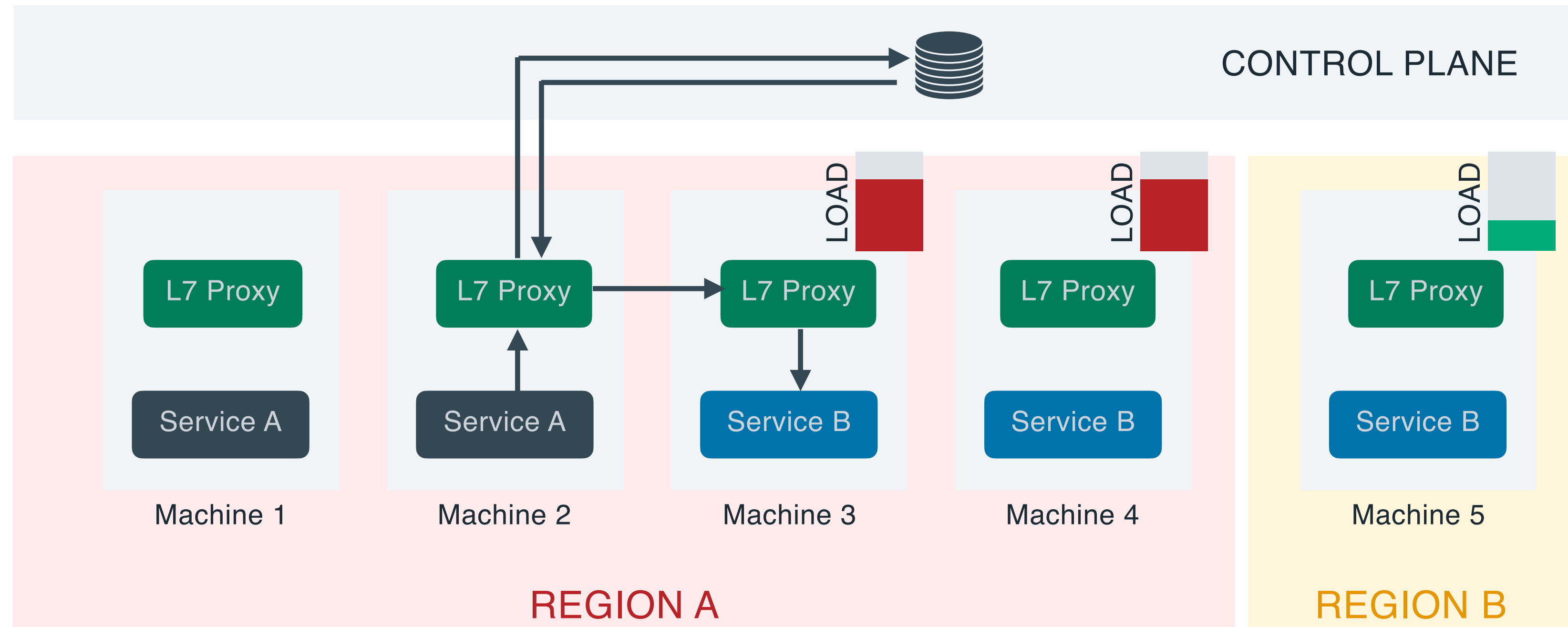
  - P90 cross-region latency: 106ms

# Service Mesh Challenges

- [SCALABILITY] How can we scale service discovery to $O(10^6)$ clients and proxies?

- [HW COST] How to minimize HW cost?

- **[RPC LATENCY & LB]** How to simultaneously minimize RPC latency and load balance across geo-distributed hosts?

  - Sidecars add extra latency

  - $O(10\text{-}10^4)$ hosts per service
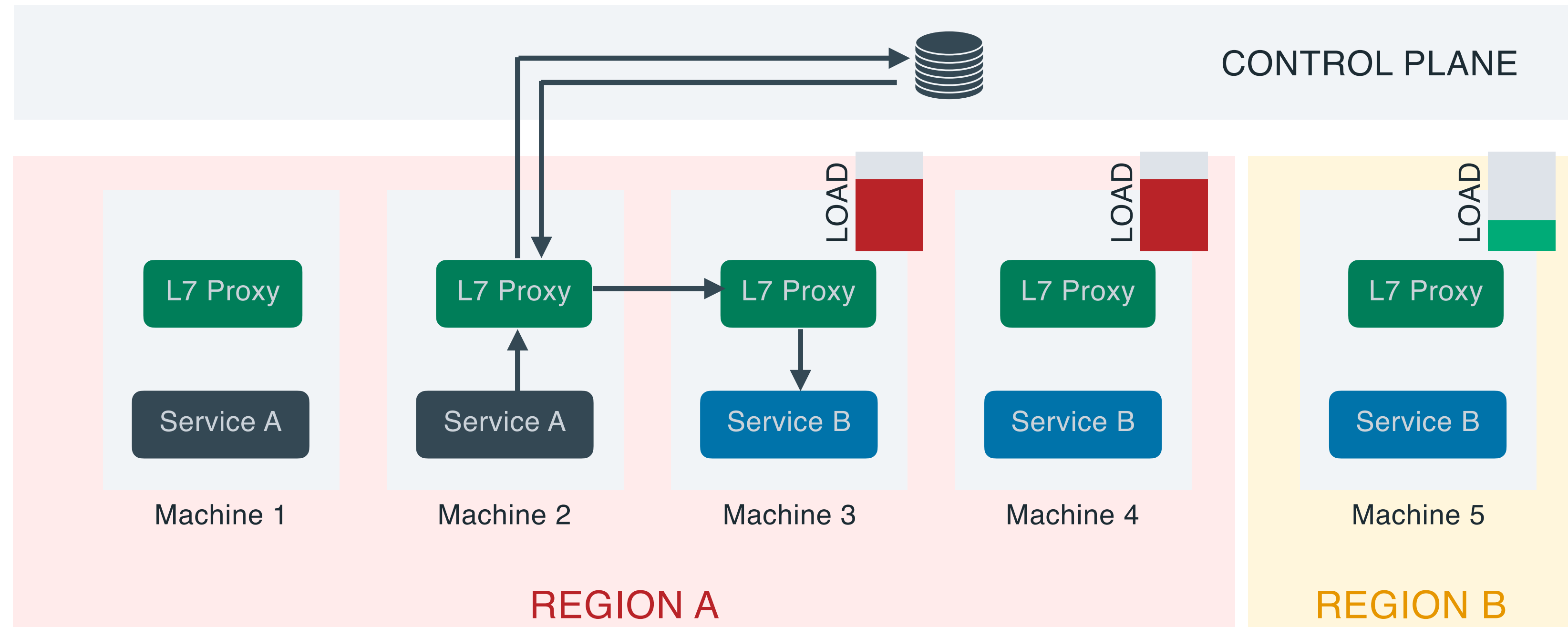
  - P90 cross-region latency: 106ms

# Service Mesh Challenges

- [SCALABILITY] How can we scale service discovery to $O(10^6)$ clients and proxies?

- [HW COST] How to minimize HW cost?

- **[RPC LATENCY & LB]** How to simultaneously minimize RPC latency and load balance across geo-distributed hosts?

  - Sidecars add extra latency

  - $O(10\text{-}10^4)$ hosts per service

  - P90 cross-region latency: 106ms

- [**SHARDED SERVICES**] Support for shared services **NOT COVERED**



9

# 03   ServiceRouter
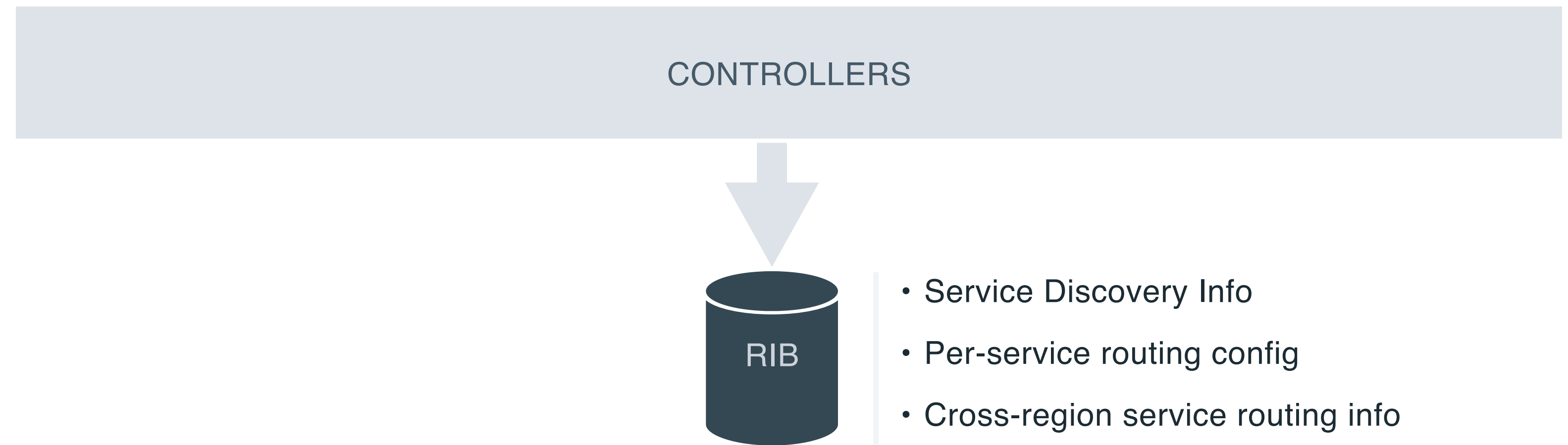
**KEY DESIGN CONCEPTS**

# RIB

Routing Information Base

Decentralize the unscalable part of the control plane in order to scale out.

- Independent controllers execute different functions such as registering services and generating a per-service cross-region routing table.

CONTROLLERS

RIB

- Service Discovery Info
- Per-service routing config
- Cross-region service routing info

11

# RIB

Routing Information Base

Decentralize the unscalable part of the control plane in order to scale out.

- Independent controllers execute different functions such as registering services and generating a per-service cross-region routing table.

- The data distribution layer massively replicates the RIB so that there are sufficient RIB replicas to handle read traffic from millions of proxies.

CONTROLLERS

RIB

- Service Discovery Info
- Per-service routing config
- Cross-region service routing info

Data Distribution Layer

# RIB
Routing Information Base

Decentralize the unscalable part of the control plane in order to scale out.

- Independent controllers execute different functions such as registering services and generating a per-service cross-region routing table.
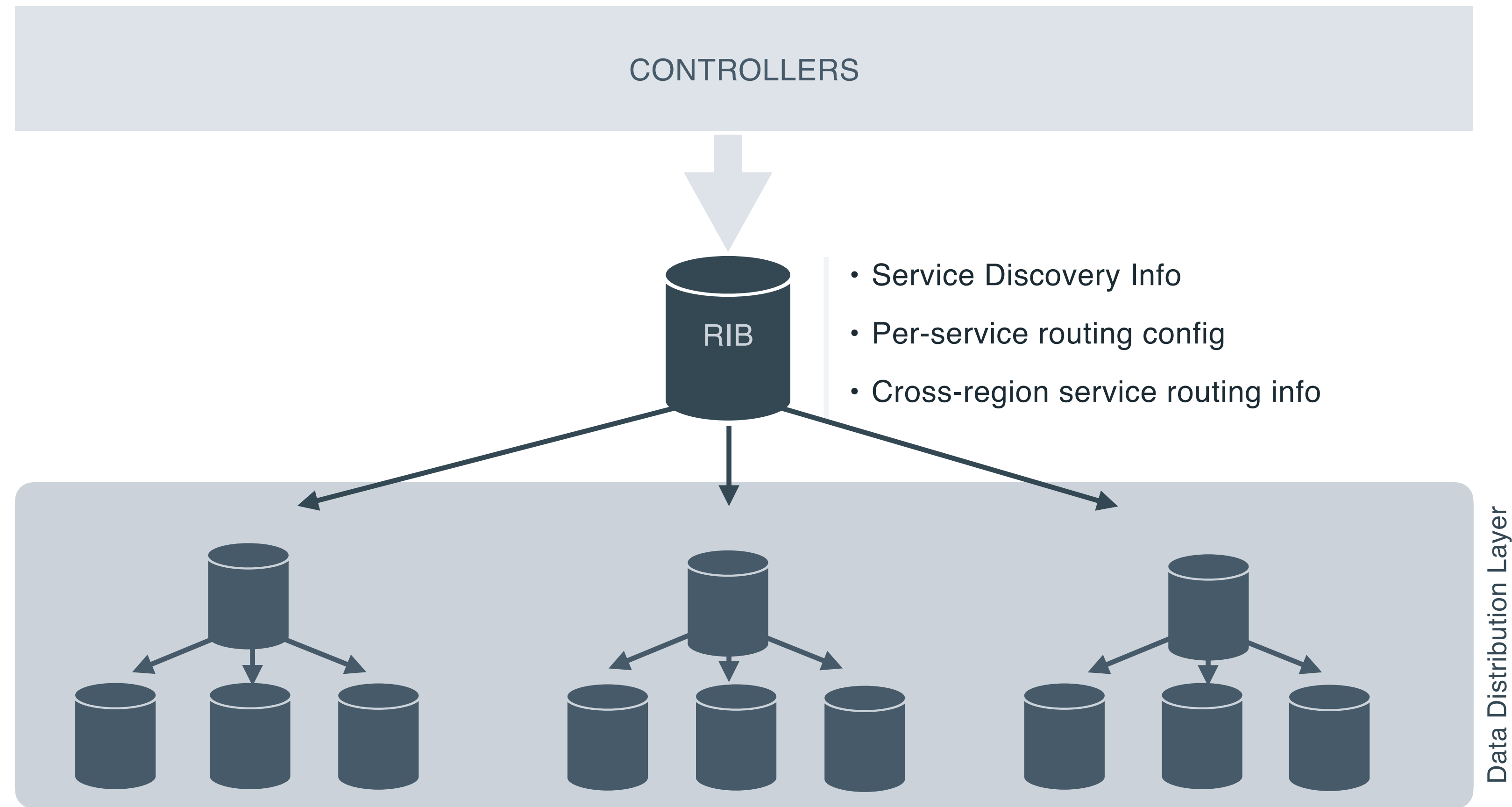
- The data distribution layer massively replicates the RIB so that there are sufficient RIB replicas to handle read traffic from millions of proxies.
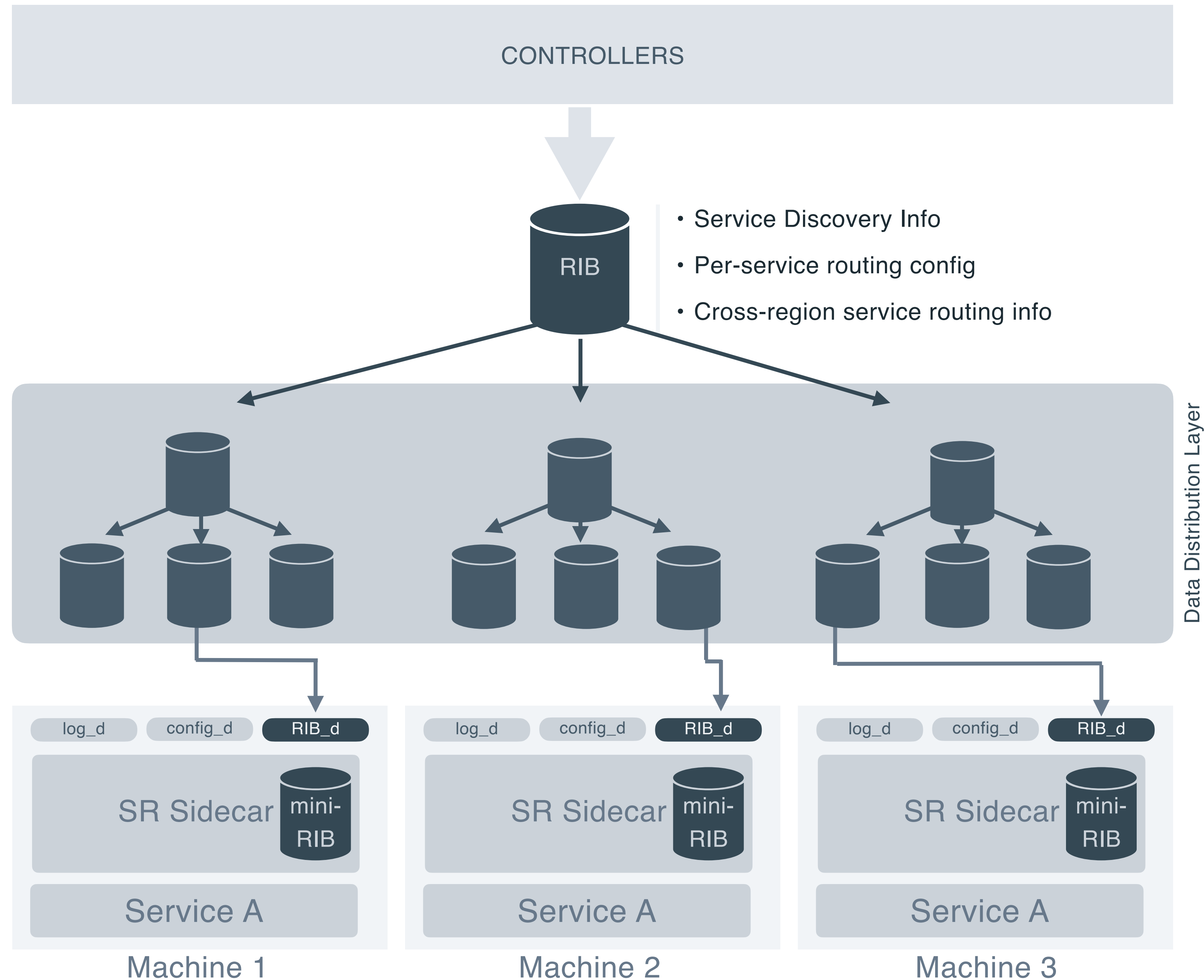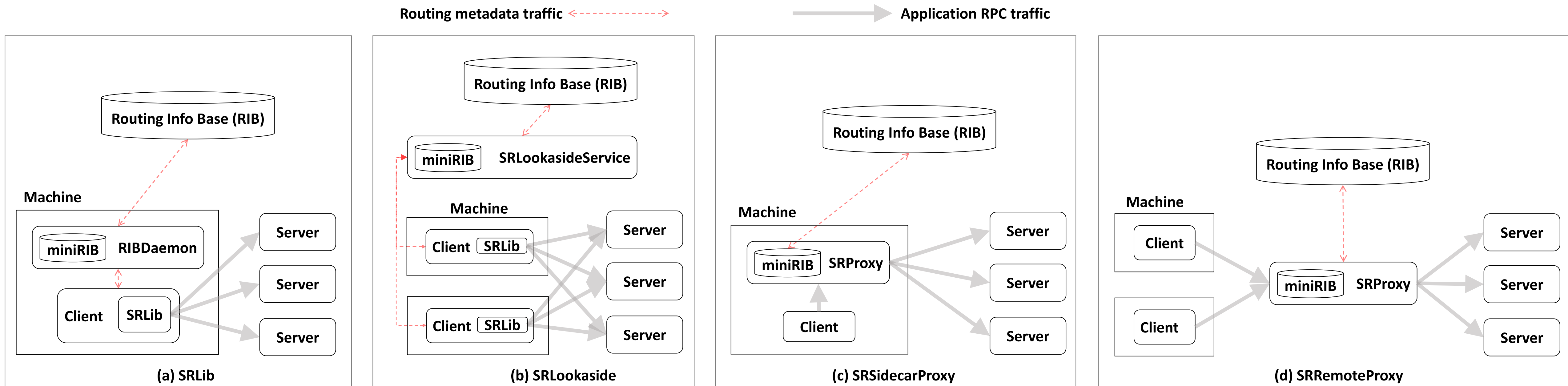
- Each proxy self-configures and self-manages without the control plane's direct involvement.

CONTROLLERS

RIB

- Service Discovery Info
- Per-service routing config
- Cross-region service routing info

Data Distribution Layer

| log_d | config_d | RIB_d |

SR Sidecar   mini-RIB

Service A

**Machine 1**

| log_d | config_d | RIB_d |

SR Sidecar   mini-RIB

Service A

**Machine 2**

| log_d | config_d | RIB_d |

SR Sidecar   mini-RIB

Service A

**Machine 3**

11

# Versatility

Controllers are agnostic to the L7 architecture.



Routing metadata traffic ← - - - - → 　　Application RPC traffic →

(a) SRLib

(b) SRLookaside

(c) SRSidecarProxy

(d) SRRemoteProxy

# 99% RPC traffic routed through SRLib.
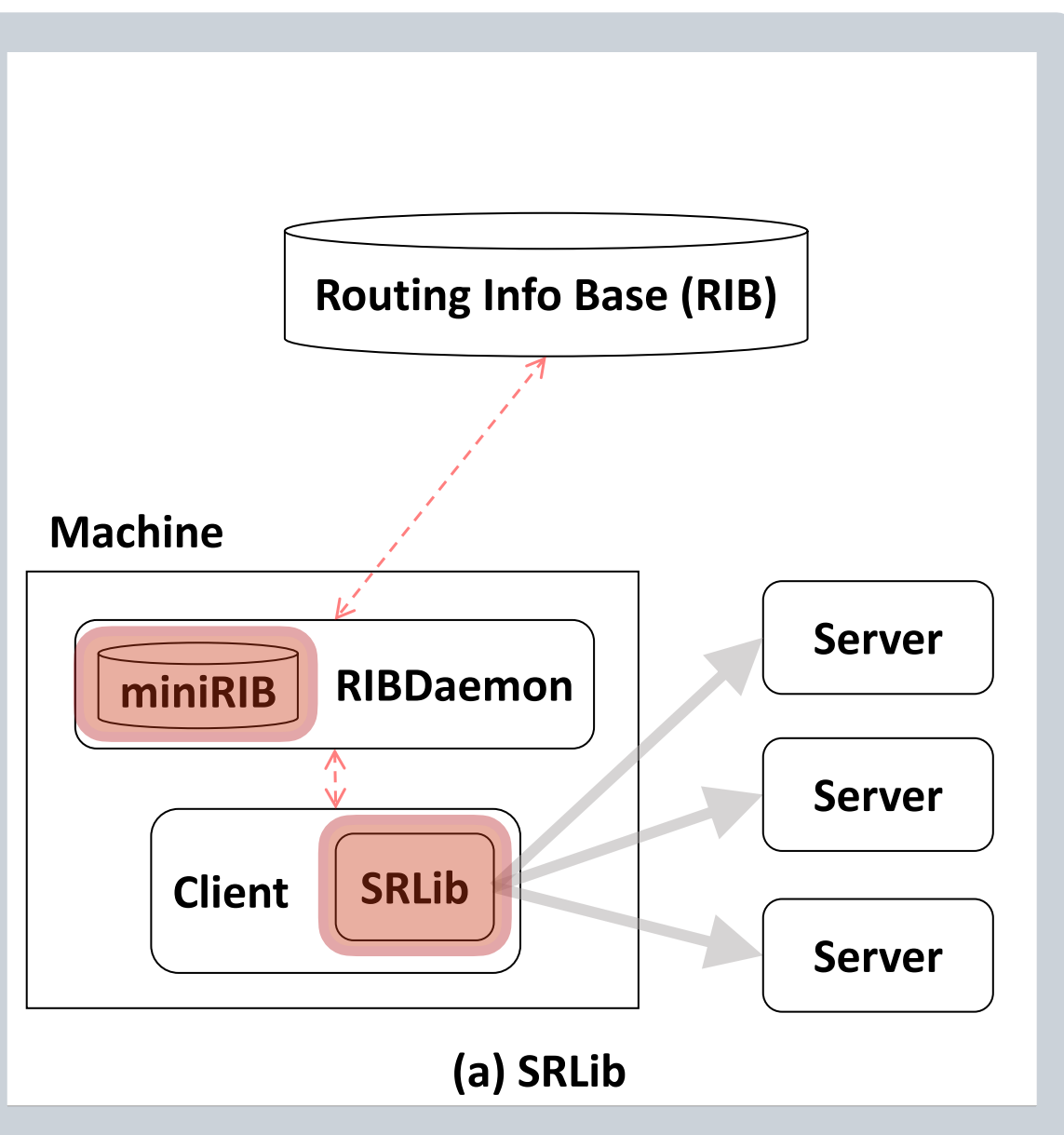
**Routing metadata traffic** ◀- - - - - - - -▶       **Application RPC traffic** ━━━▶



(a) SRLib

# SRLib

Provide the service-mesh functions out of a library that is directly linked into the RPC client's executable

- Eliminates side car latency overhead

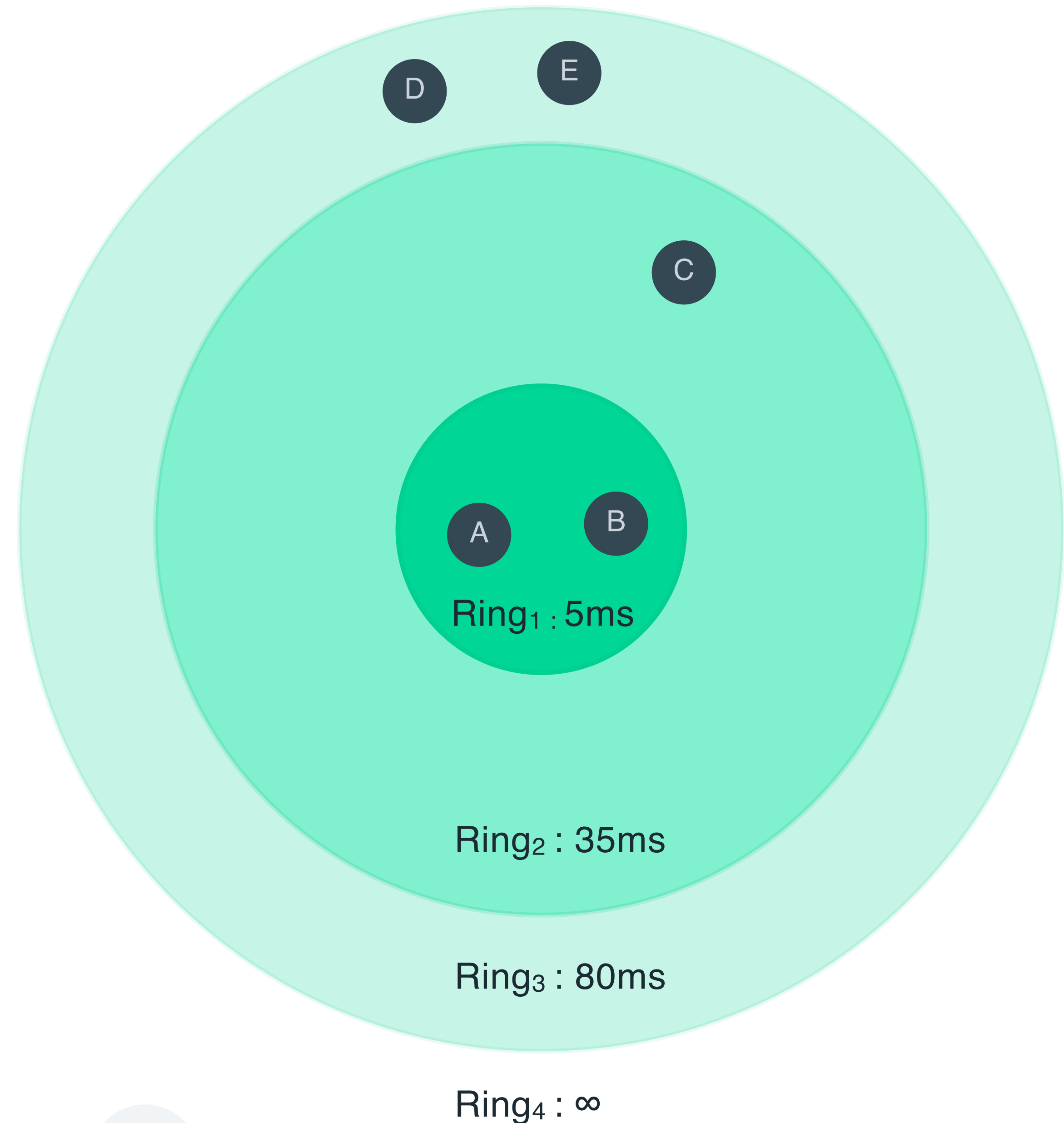Run a separate RIBDaemon on the client machine to cache **miniRIB**.

# LATENCY RINGS AND CROSS-REGION ROUTING

SR strives to simultaneously minimize RPC latency and balance load across global regions.

- SR introduces the concept of latency rings to minimize latency.

- SR collects per-service global traffic and load information, computes a per-service cross-region routing table, and disseminate it to L7 routers to guide their local routing decisions.

$Ring_1 : 5ms \mid Ring_2 : 35ms \mid Ring_3 : 80ms \mid Ring_4 : \infty$



D    E

C

A    B

$Ring_1 : 5ms$

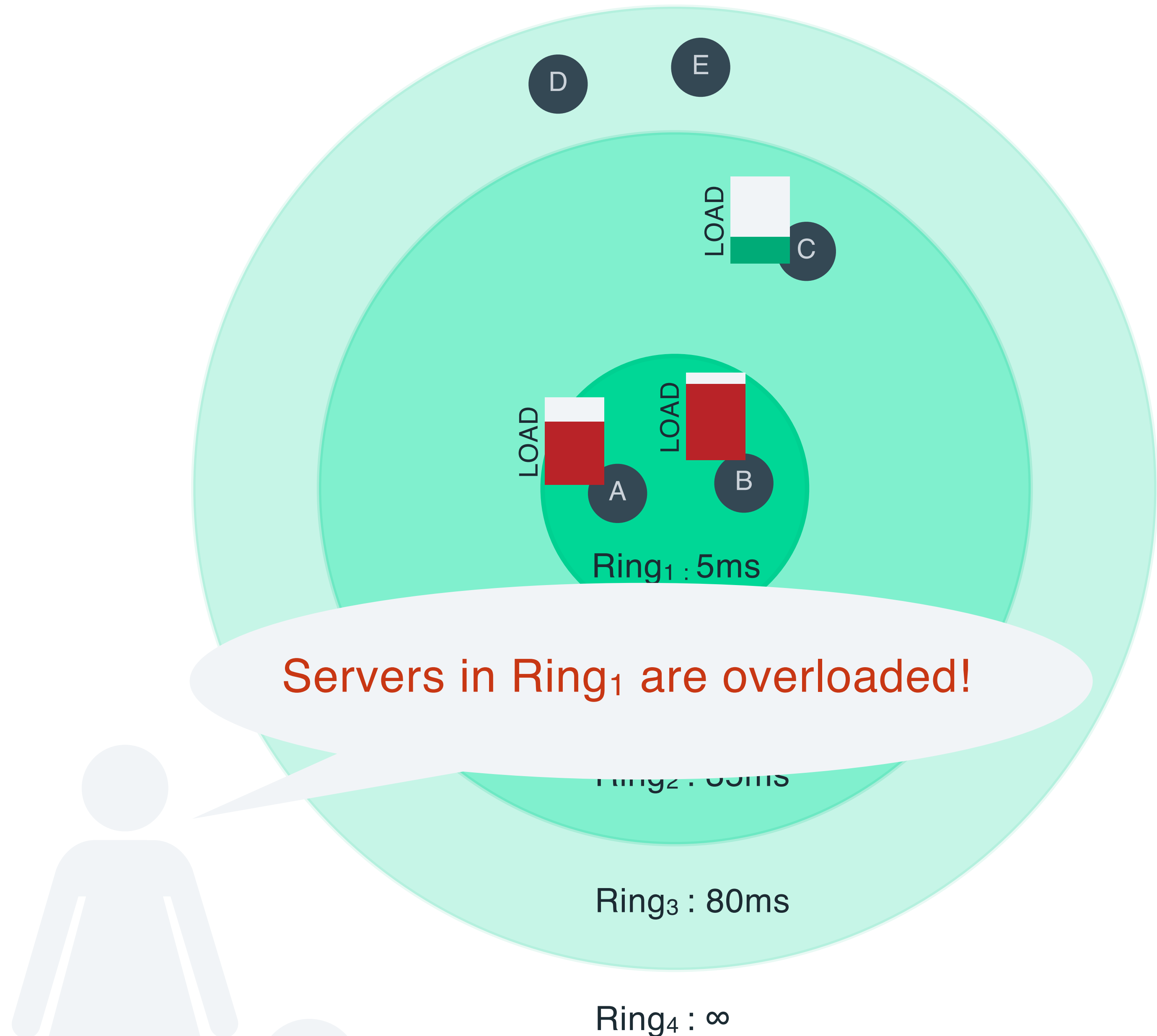$Ring_2 : 35ms$

$Ring_3 : 80ms$

$Ring_4 : \infty$

# LATENCY RINGS AND CROSS-REGION ROUTING

SR strives to simultaneously minimize RPC latency and balance load across global regions.

- SR introduces the concept of latency rings to minimize latency.

- SR collects per-service global traffic and load information, computes a per-service cross-region routing table, and disseminate it to L7 routers to guide their local routing decisions.
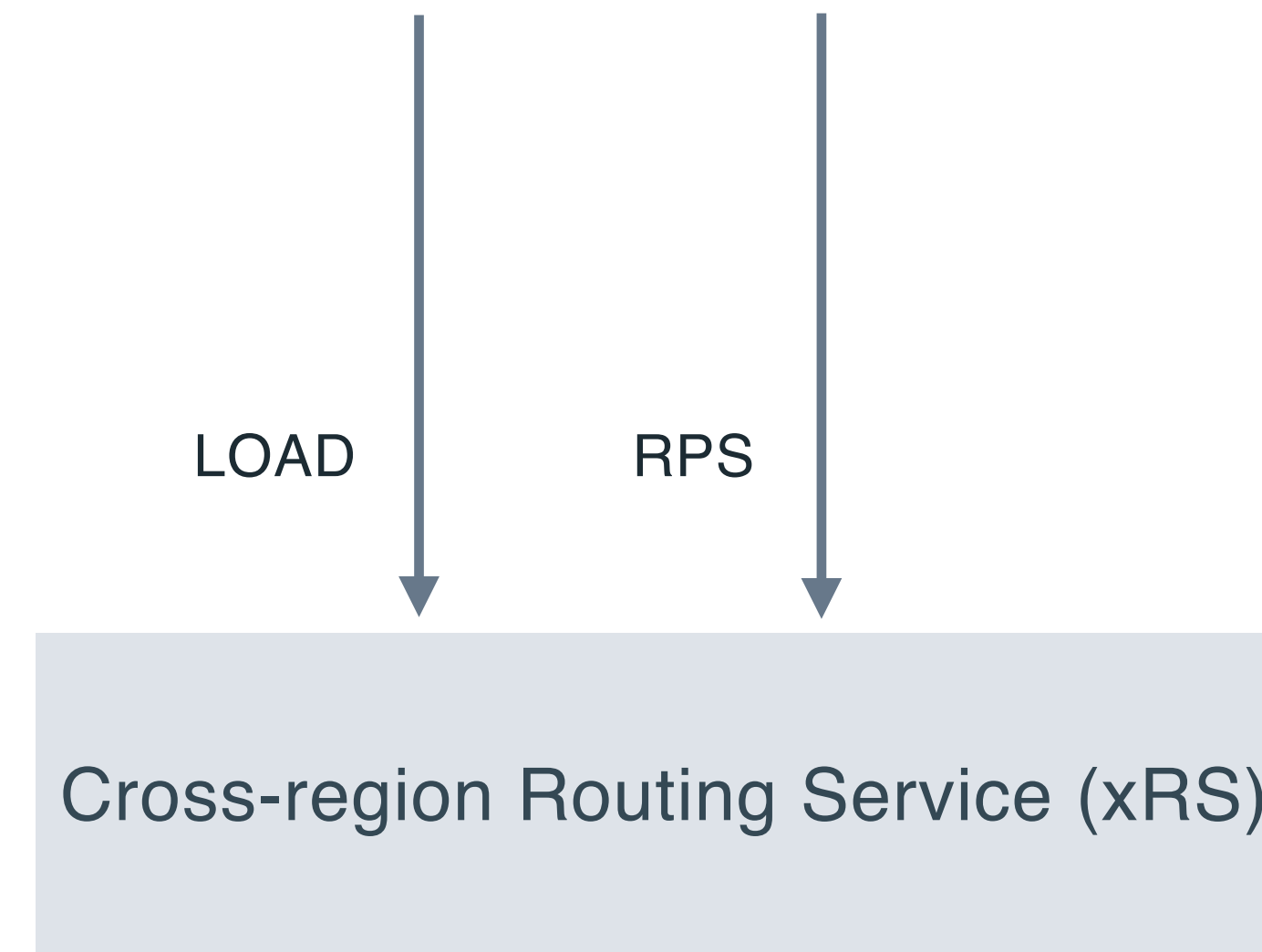
$Ring_1 : 5ms \mid Ring_2 : 35ms \mid Ring_3 : 80ms \mid Ring_4 : \infty$



**Servers in $Ring_1$ are overloaded!**

$Ring_1 : 5ms$

$Ring_2 : 35ms$
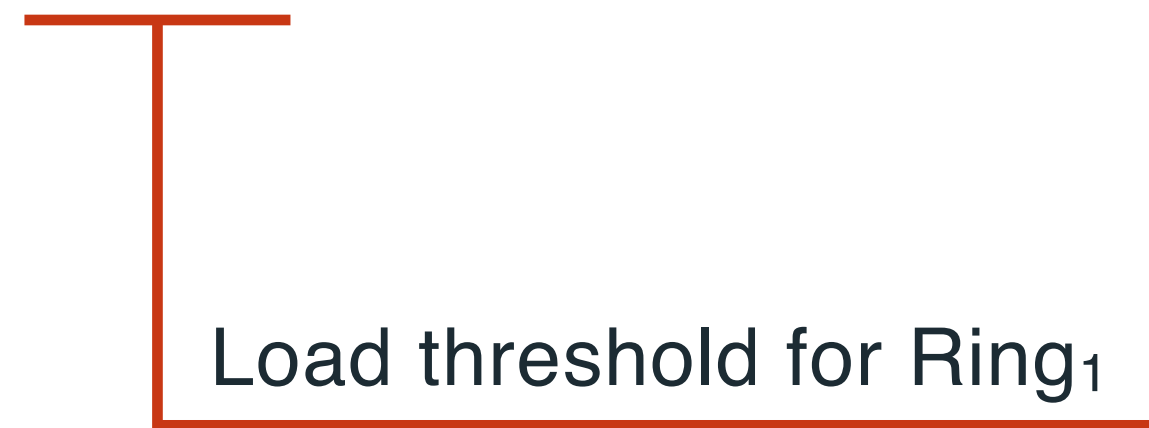
$Ring_3 : 80ms$

$Ring_4 : \infty$

# LATENCY RINGS AND CROSS-REGION ROUTING

SR strives to simultaneously minimize RPC latency and balance load across global regions.

- SR introduces the concept of latency rings to minimize latency.

- SR collects per-service global traffic and load information, computes a per-service cross-region routing table, and disseminate it to L7 routers to guide their local routing decisions.
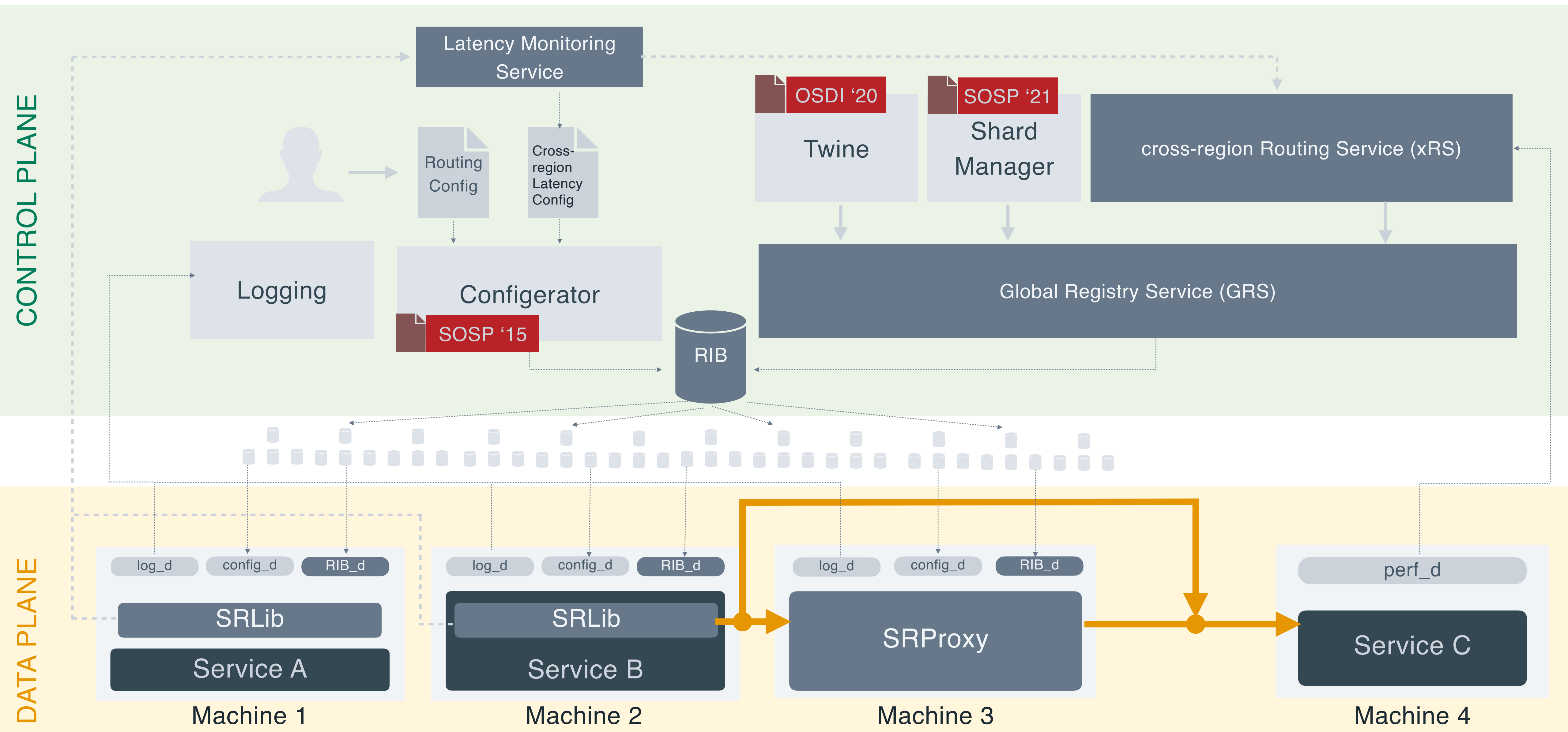
LOAD          RPS

Cross-region Routing Service (xRS)

$Ring_1 : 5ms : 55\%$ | $Ring_2 : 35ms : 65\%$ | $Ring_3 : 80ms : 80\%$ | $Ring_4 : \infty : \infty$

Load threshold for $Ring_1$

# 04   ServiceRouter

## OVERALL ARCHITECTURE

# 05    ServiceRouter

**EVALUATION**

# Scalability

Overall scale

- Regions
- Routers/Clients/Servers
- Throughput

$O(10)$

Regions

$O(10^6)$

L7 Routers

Clients

$O(10^9)$

RPS

$O(10^{14})$

B/sec

19

# Scalability

RIB - Routing Information Base

- RIB Replicas
- RIB Write bandwidth
- RIB Write throughput



O($10^9$)
Bytes

RIB

O(10)
Paxos
Acceptors

$O(10^2)$
commits/
sec
___
WRITE

$O(10^6)$
B/sec
___
WRITE

O($10^{12}$)
Bytes

Data Distribution Layer
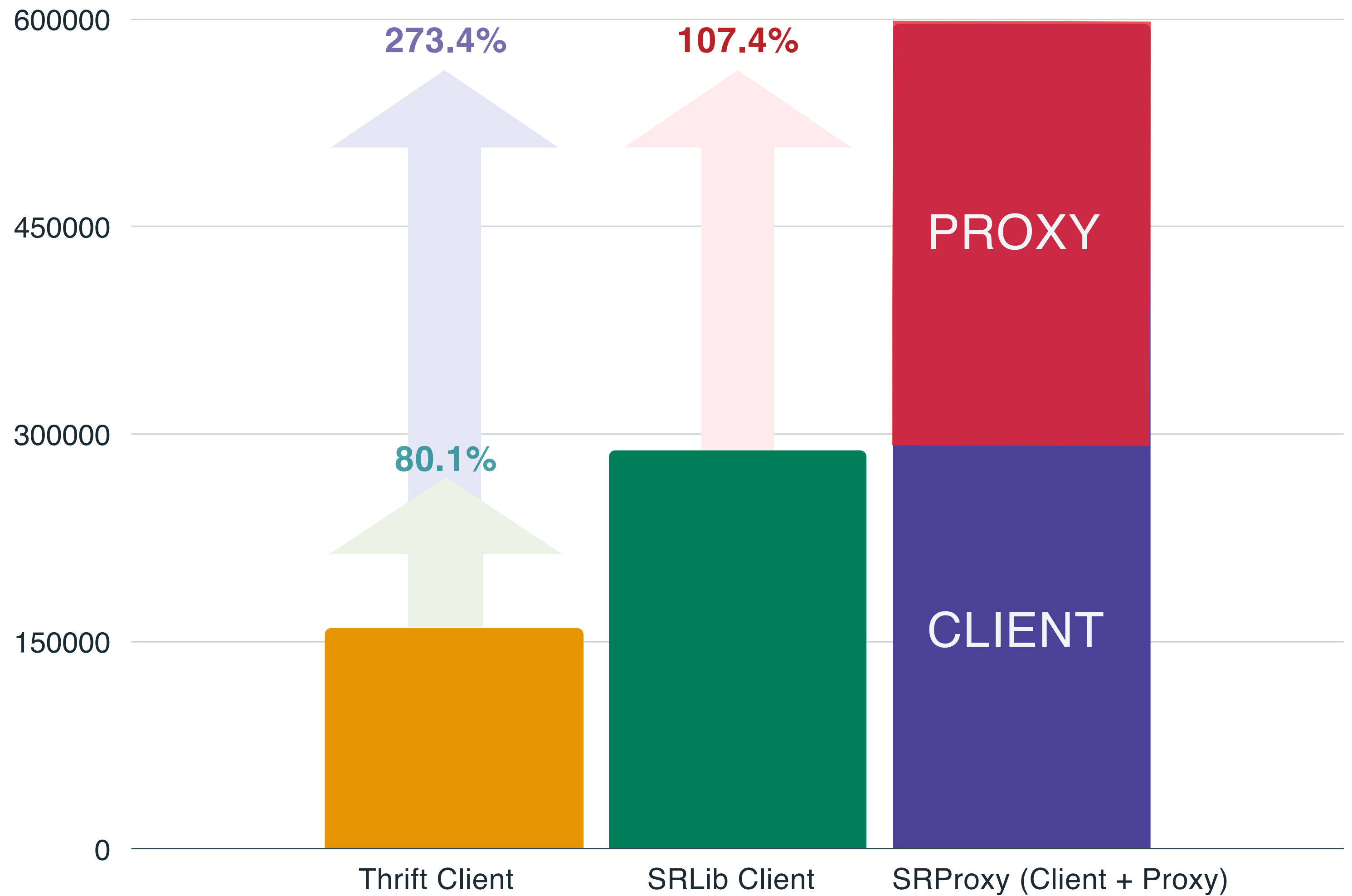
O($10^3$)
Paxos
Learners
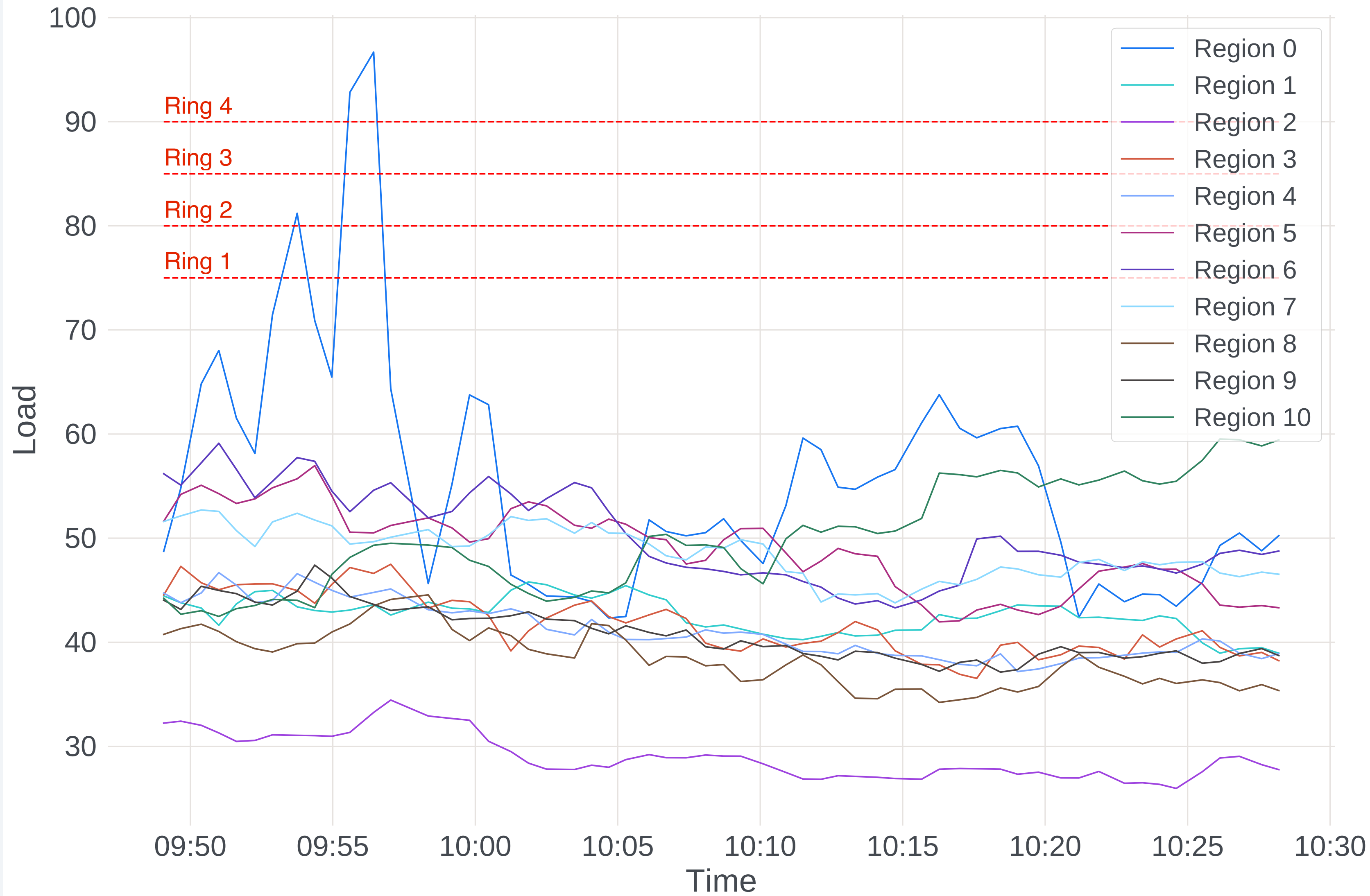
20

# Cost

METHODOLOGY

- Metrics: P50 avg request latency; CPU Instructions per request
- Designs
  - Baseline: Thrift RPC
  - SRLib
  - Remote SRProxy
- Simulated Payload:
  - Production avg request and avg response size
    - $O(10^3)$ B
- 100K requests
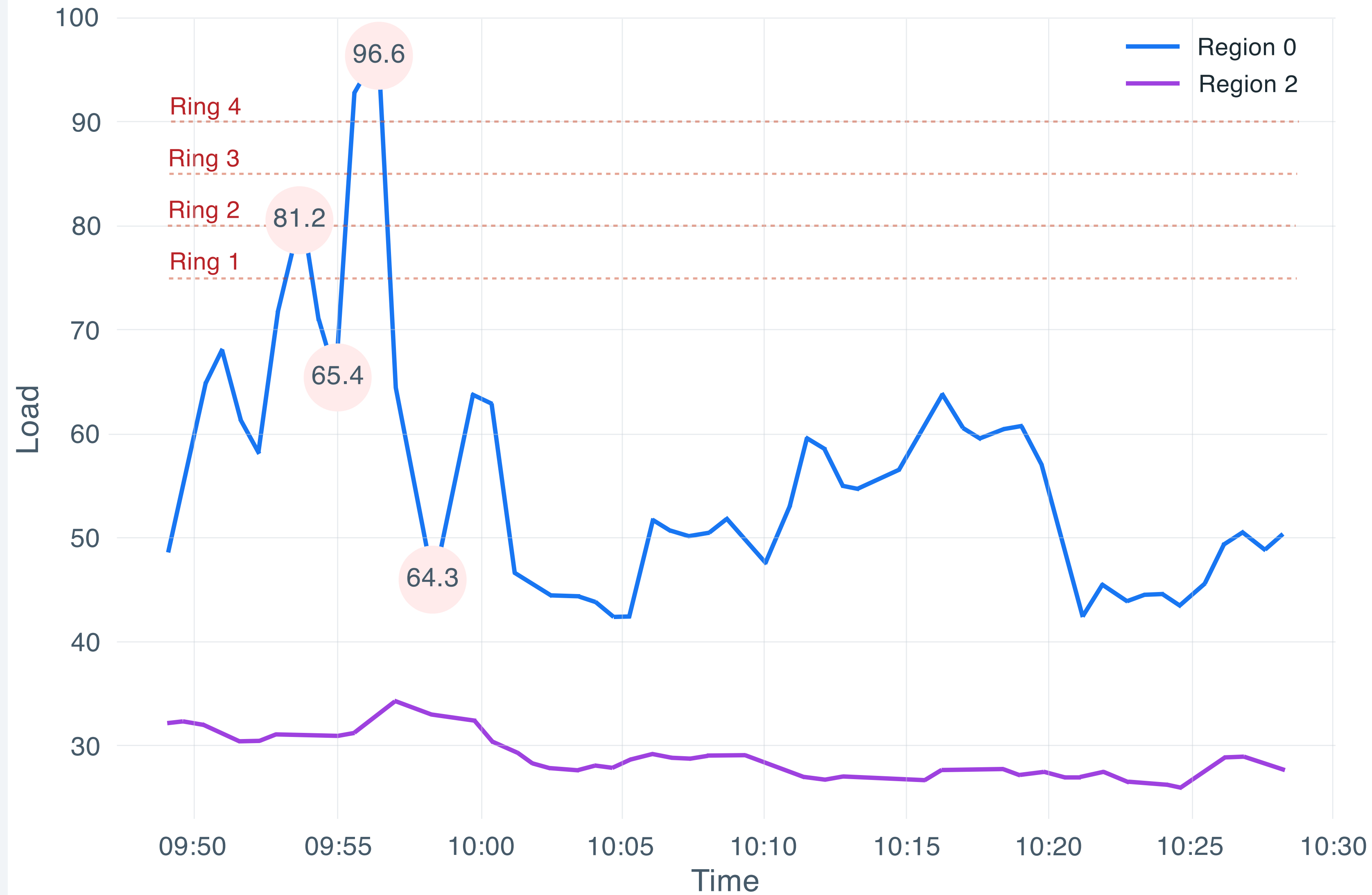- 3 trials per design

## CPU Instructions/Request

# Cross-Region Load Shift

- Real-world Example

# Cross-Region Load Shift

- Real-world Example
- 9:53 —> Region 0 Load = 81.2%
  - xRS Traffic Shifts
    - R0 -5.35%
    - R0 to R2 +5.35%
- 9:54 —> Region 0 Load = 65.47%
- 9:56 —> Region 0 Load = 96.69%
  - xRS Traffic Shifts
    - R0 -25%
    - R0 to R2 +25%
- 9:57 —> Region 0 Load = 64.34%

# ServiceRouter

**HYPERSCALE AND MINIMAL COST SERVICE MESH AT META**

## 06   Summary

ServiceRouter's massive RIB replication allows decentralizing L7 router management and to scale to millions of routers and proxies.

ServiceRouter routes 99% of the traffic with an optimized embedded library approach with astounding HW savings.

ServiceRouter's source-based locality rings and xRS strike a balance between latency wins and load balancing.

Built-in **support for sharded services** which account for 68% of our RPCs **[not covered in this talk]**.

Meta     Imperial College London     Carnegie Mellon University Computer Science Department

Soteris Demetriou | s.demetriou@imperial.ac.uk