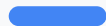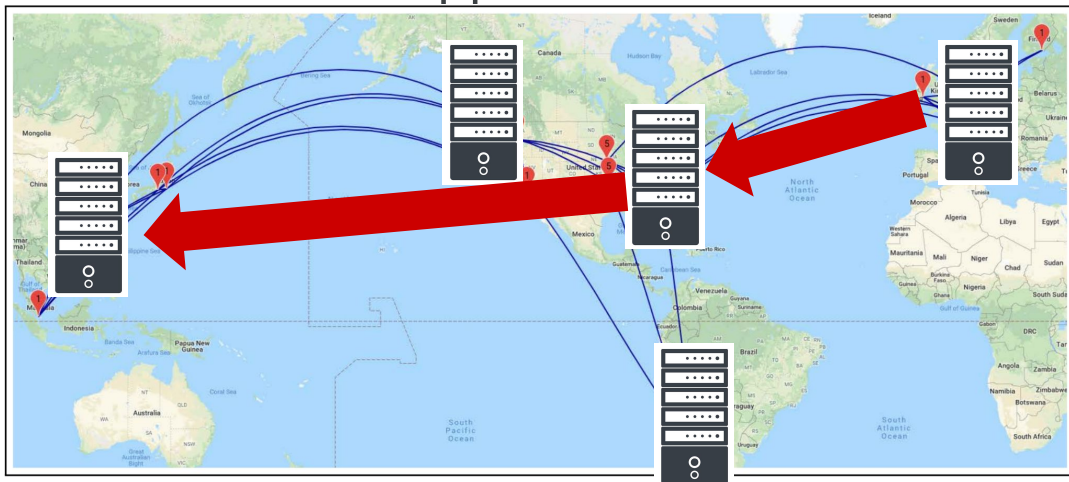Google

# A Cloud-Scale Characterization of Remote Procedure Calls

**Korakit Seemakhupt**, Brent Stephens, Samira Khan, Sihang Liu, Hassan Wessel,

Soheil Hassas Yeganeh, Alex C. Snoeren, Arvind Krishnamurthy, David Culler, Henry M. Levy

# Motivation

A cloud-scale application must be:



- Scalable
- Highly-available
- Failure tolerant
- Easy to maintain

*Understanding RPC is a key to understanding global-scale distributed services*

Google

# RPC Study: Cloud-scale Workloads

This is a study of RPC at Google Scale

We include
- Google's **first party** web services
  - Search, Gmail, Youtube, ….
- Google's **internal services**
  - Spanner, Bigtable, F1, …

We do **not** include
- RPCs serving **Cloud customers** (GCP)
- **RDMA** and software-based **RMA** communication (Snap/Pony Express)

Google

# RPC Study: Measurements

We collected and processed data using three Google internal monitoring systems

Overall we examined:
- Over **700 billion** RPC traces
- **10,000** different RPC methods from over **100** production clusters
- System statistics collected every **30 minutes for ~2 years**

Aggregated statistics include:
- Latency Components, Payload Size, Call Structure
- CPU Utilization, Memory Bandwidth, Scheduling Latency
- Requests/Second, Growth rate, ….

Google

# **Agenda**

- Motivations for studying RPC in the cloud

- RPC Study Results

  - What is the **source of RPCs**? Where do they go?
  - What is the **timescale** of RPC?
  - Which **latency component** affects RPC latency?
  - Which **RPC Latency Tax** component is the bottleneck?
  - How does **utilization** vary across datacenters?
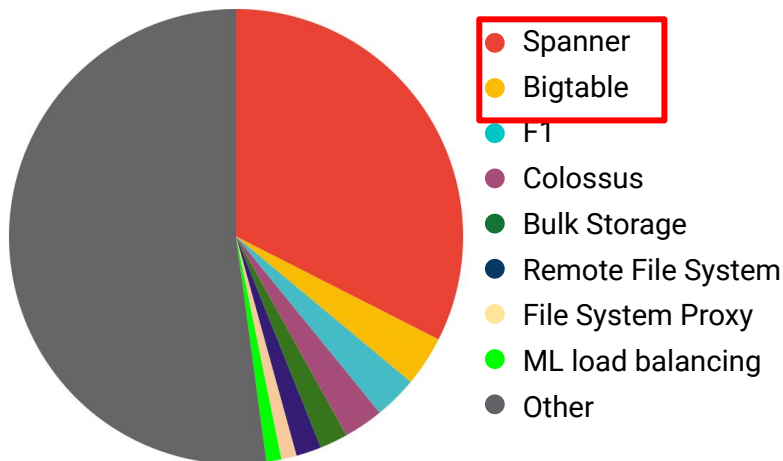- Implications from the study

Google

# Agenda

- Motivations for studying RPC in the cloud

- RPC Study Results

  - What is the **source of RPCs**? Where do they go?
  - What is the **timescale** of RPC?
  - Which **latency component** affects RPC latency?
  - RPC Latency Tax
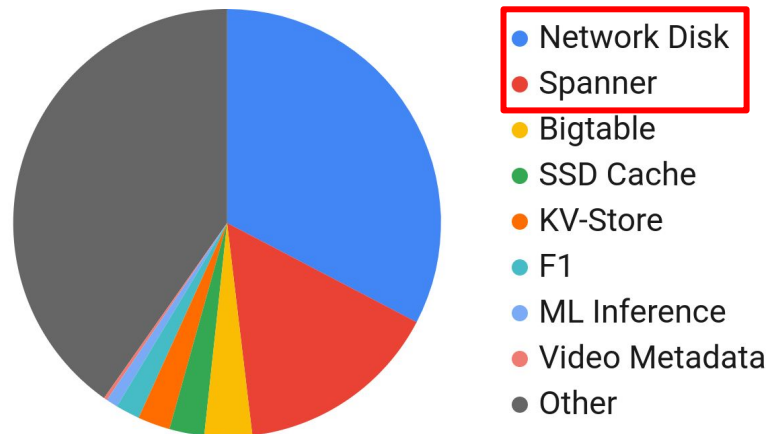  - How does **utilization** vary across datacenters?
- Implications from the study

Google

# Google's RPC Environment
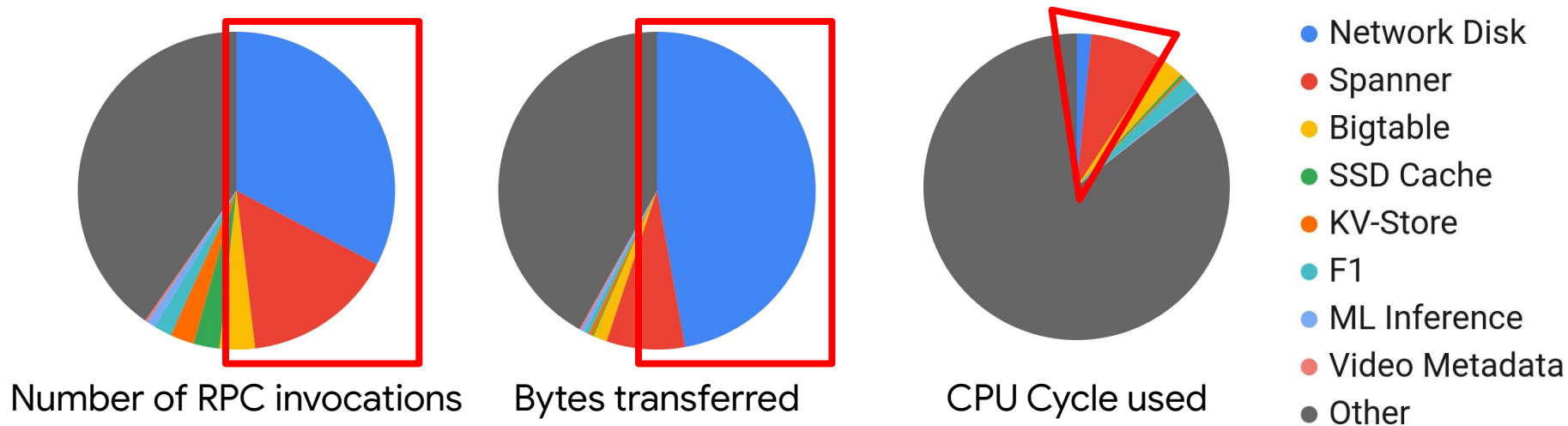
## RPC Sources and Destinations

### RPC Sources



- Spanner
- Bigtable
- F1
- Colossus
- Bulk Storage
- Remote File System
- File System Proxy
- ML load balancing
- Other

### RPC Destinations



- Network Disk
- Spanner
- Bigtable
- SSD Cache
- KV-Store
- F1
- ML Inference
- Video Metadata
- Other

Google's Internal RPC is dominated by communication between **storage services**.

Google

# Google's RPC Environment
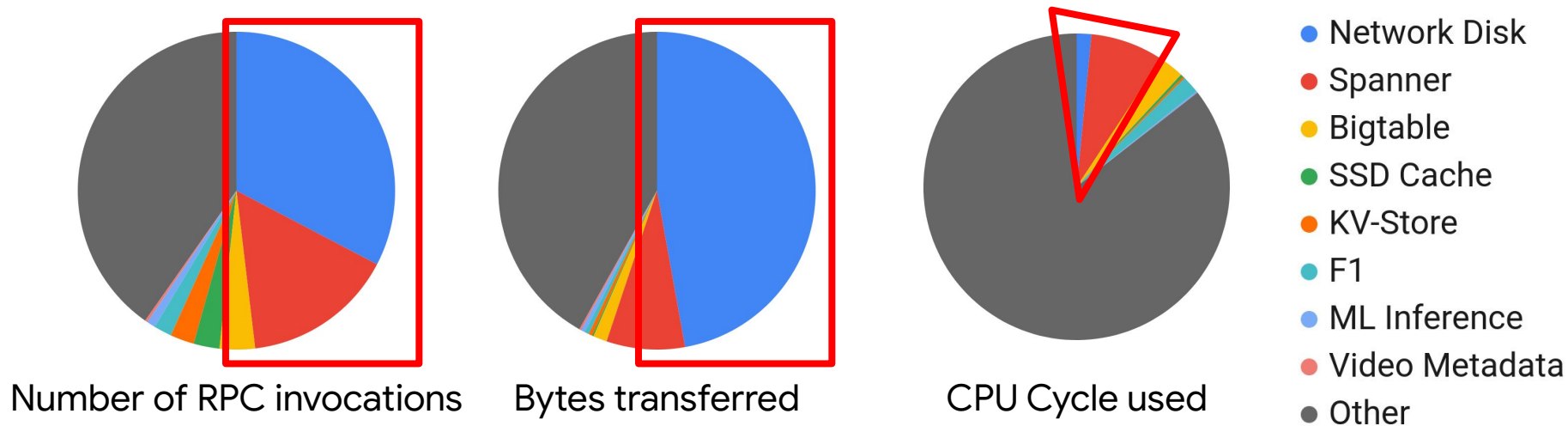
## RPC Popularity and Resource Utilization (by destination)



Number of RPC invocations     Bytes transferred     CPU Cycle used

- Network Disk
- Spanner
- Bigtable
- SSD Cache
- KV-Store
- F1
- ML Inference
- Video Metadata
- Other

Half of RPC invocations and data transferred are from **Spanner** and **Network Disk**

8

Google

# Google's RPC Environment

## RPC Popularity and Resource Utilization (by destination)



Number of RPC invocations

Bytes transferred

CPU Cycle used

- Network Disk
- Spanner
- Bigtable
- SSD Cache
- KV-Store
- F1
- ML Inference
- Video Metadata
- Other

**Takeaway**: Storage RPC is by far the largest contributor to fleet-wide RPC and bytes transfer in the network.
This motivates for research on **data-movement acceleration**.
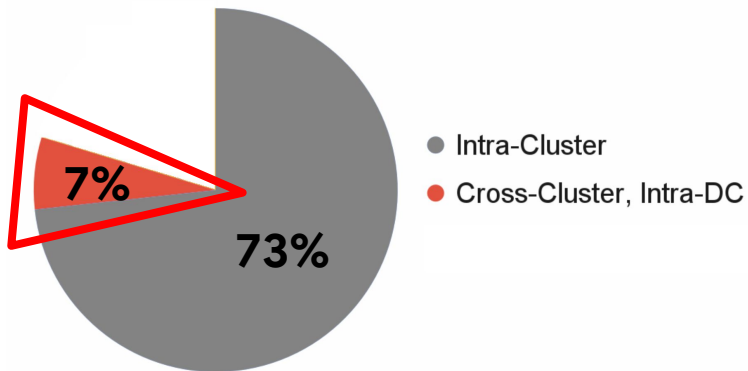
Google

# Google's RPC Environment

## Cross-cluster RPC and WAN



Cluster 1

Cluster 2

**B4 WAN**

Cluster N

Cluster N-1

**Google's geo-distributed datacenters**
Each datacenter can consist of **multiple clusters**
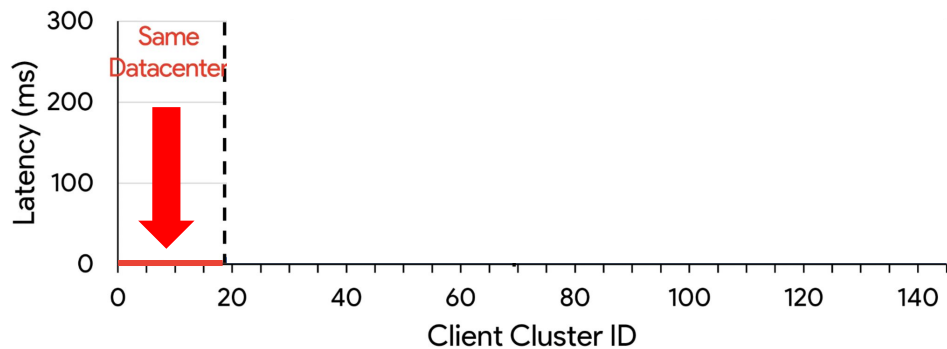Datacenters are connected through **WAN links (B4)**

Google

# Google's RPC Environment

## Cross-cluster RPC and WAN

Cross Cluster Median Latency



- Intra-Cluster
- Cross-Cluster, Intra-DC

7% of RPCs are **cross-cluster**, but in **the same datacenter**

Same Datacenter RTT is under 10 ms

Google

# Google's RPC Environment

## Cross-cluster RPC and WAN



- Intra-Cluster
- Cross-Cluster, Intra-DC
- Cross-Cluster over B4 WAN

20%

7%

73%

### Cross Cluster Median Latency



Same Datacenter

Different Datacenter Same Contitent

Different Continents

Latency (ms)

Cluster Percentile

20% of RPCs are **cross-cluster** over **B4 WAN**   Cross-continent RTT can be over 200 ms

**Takeaway**: RPC locality significantly affects the latency. Cross-cluster RPCs over WAN introduces significant overhead
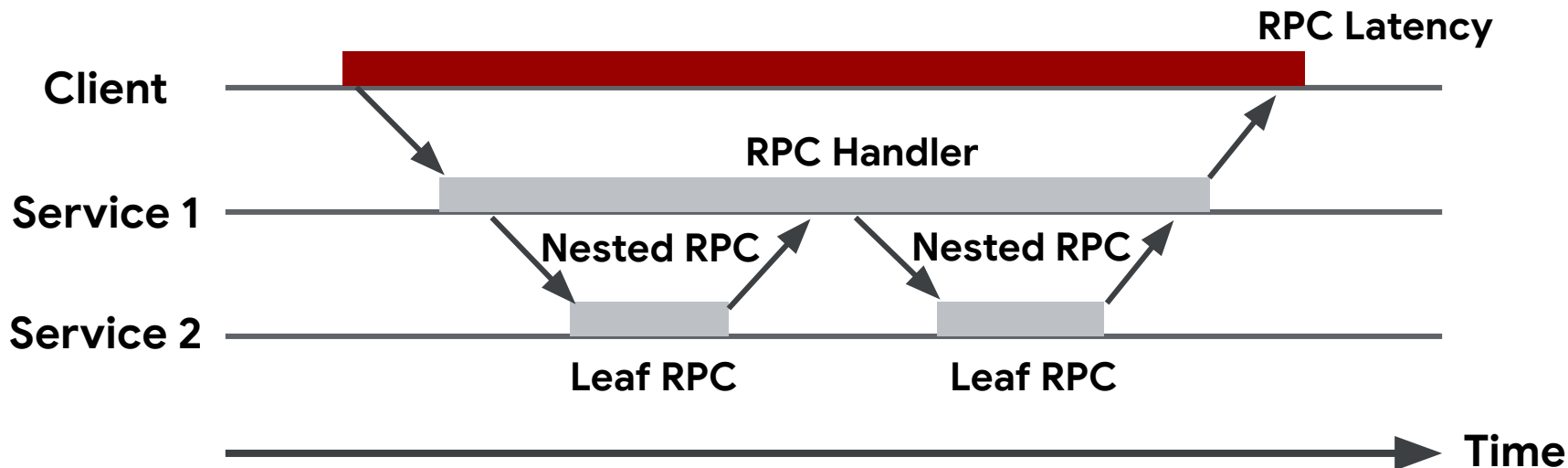
Google

# RPC Study Results:
# **Takeaways**

What is the **source of RPCs**? Where do they go?

- Storage RPCs are the largest contributor to fleet-wide RPCs
  - Motivates research on **data-movement acceleration**

- RPC locality significantly affects latency
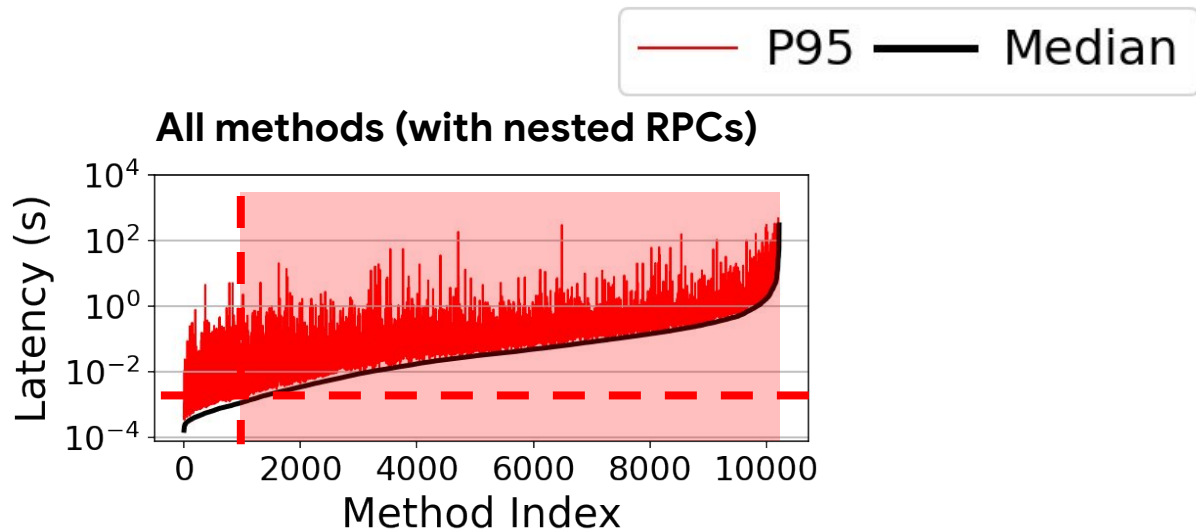  - Motivates research on **locality-aware scheduling**

Google

# Agenda

- Motivations for studying RPC in the cloud

- RPC Study Results

    - What is the **source of RPCs**? Where do they go?
    - What is the **timescale** of RPC?
    - Which **latency component** affects RPC latency?
    - Which **RPC Latency Tax** component is the bottleneck?
    - How does **utilization** vary across datacenters?

- Implications from the study

Google

# RPC completion time includes nested RPC calls

RPC Latency

Client

RPC Handler

Service 1

Nested RPC          Nested RPC

Service 2

Leaf RPC              Leaf RPC

Time

RPC Latency includes **RPC handler** and **nested RPC calls**.
We also show **leaf RPC latency**

15

Google

# What is the **timescale** of RPC?



**All methods (with nested RPCs)**

Legend: P95 (red), Median (black)

90% of RPC methods have median latency **over a millisecond**.

Google

# What is the **timescale** of RPC?



**All methods (with nested RPCs)**

**Leaf RPC methods**

90% of RPC methods have median latency **over a millisecond.**

53% of leaf RPC methods have median latency **over a millisecond**.

**Takeaway**: Majority of RPC methods in this environment
are **millisecond**,  not microsecond scale

Google

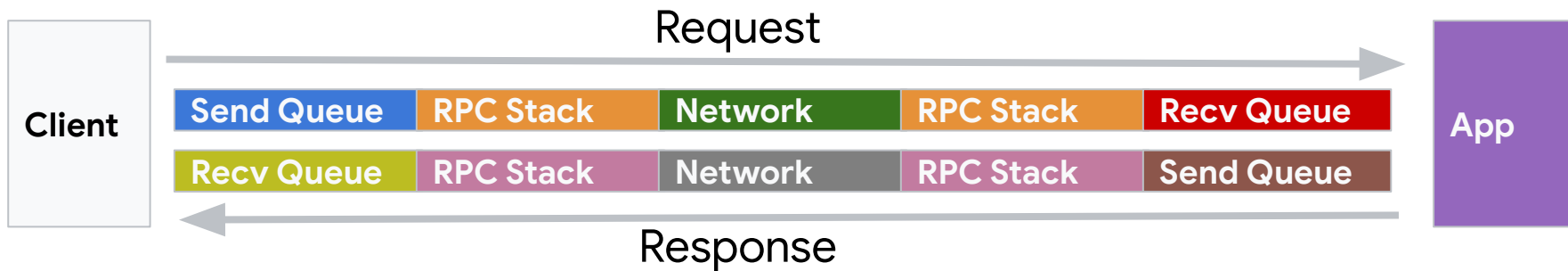# RPC Study Results:
# **Takeaways**

What is the **timescale** of RPC?

- Majority of RPC methods in this environment are **millisecond** scale


- But half of the **leaf RPC** methods have **sub-millisecond** latency
    - Optimizing for **latency** is still important for median **leaf RPCs**
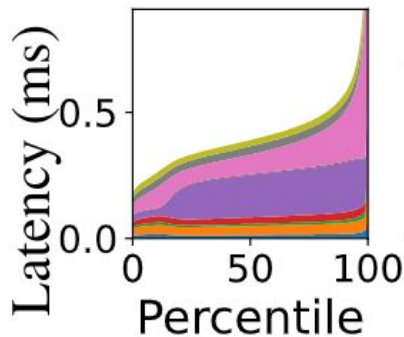
Google

# Agenda

- Motivations for studying RPC in the cloud

- RPC Study Results

  - What is the **source of RPCs**? Where do they go?

  - What is the **timescale** of RPC?

  - Which **latency component** affects RPC latency?

  - Which **RPC Latency Tax** component is the bottleneck?

  - How does **utilization** vary across datacenters?

- Implications from the study

19

Google

# What are the **Latency components**?



| | Request → | | | | |
|---|---|---|---|---|---|
| **Client** | Send Queue | RPC Stack | Network | RPC Stack | Recv Queue |
| | Recv Queue | RPC Stack | Network | RPC Stack | Send Queue |
| | | | ← Response | | |

**App**

We measure time spent on **queues**, **RPC stack**, **network**, and the **application processing time** in leaf RPCs.
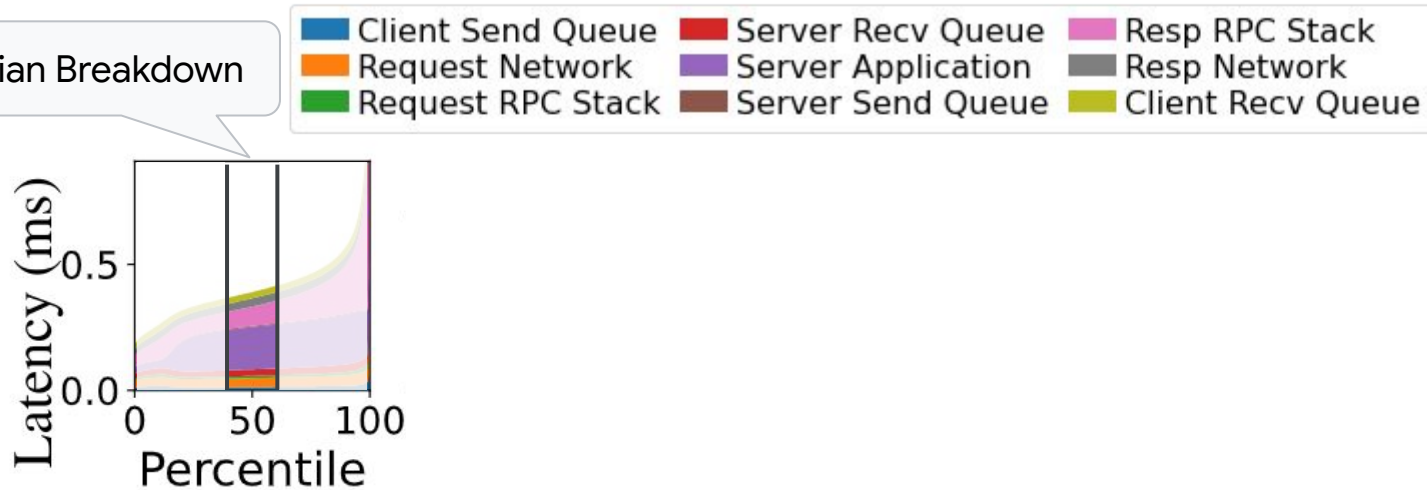
Google

# What are the different causes of latency?



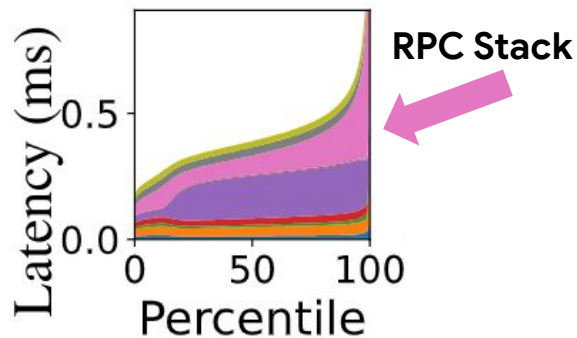**RPC Stack dominated:**
**(e.g., K/V-Store)**

Google

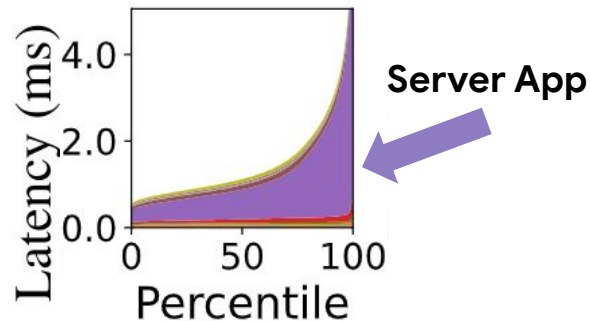# What are the different causes of latency?

Median Breakdown

| Client Send Queue | Server Recv Queue | Resp RPC Stack |
| Request Network | Server Application | Resp Network |
| Request RPC Stack | Server Send Queue | Client Recv Queue |

**RPC Stack dominated:**
**(e.g., K/V-Store)**

Google

# What are the different causes of latency?



Legend: Client Send Queue, Request Network, Request RPC Stack, Server Recv Queue, Server Application, Server Send Queue, Resp RPC Stack, Resp Network, Client Recv Queue

**RPC Stack dominated: (e.g., K/V-Store)**

**RPC method dominated: (e.g., ML Inference)**

**Queueing dominated: (e.g., SSD Cache)**

Different dominant component for each application
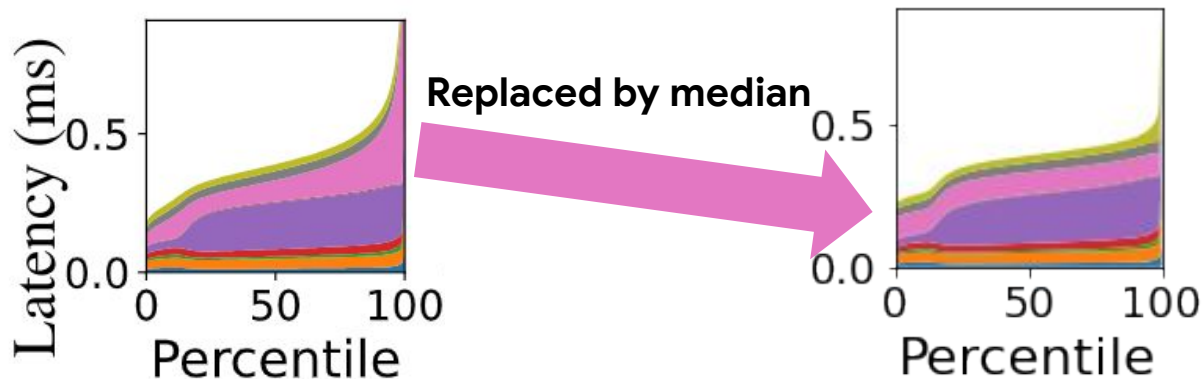
Google

# What-if analysis with causal modeling

Research question:

What is the latency component that is most responsible for tail latency?

How much latency can we improve by optimizing that component?
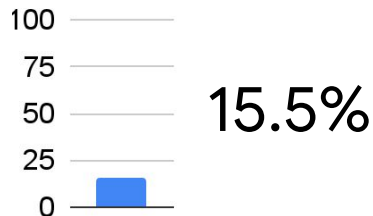
Methodology:



**Replaced by median**

Goal: Understand impact of reducing the variation
caused by that component

Google

# Potential latency improvement

Potential Tail Latency Improvement

| 100 |
|---|
| 75 |
| 50 |
| 25 |
| 0 |

15.5%

Potential Tail Latency Improvement

| 100 |
|---|
| 75 |
| 50 |
| 25 |
| 0 |

68.0%

Potential Tail Latency Improvement

| 100 |
|---|
| 75 |
| 50 |
| 25 |
| 0 |

33.6%

**RPC bottlenecked: K/V-Store**

**RPC optimization** can improve tail latency by 15.5%
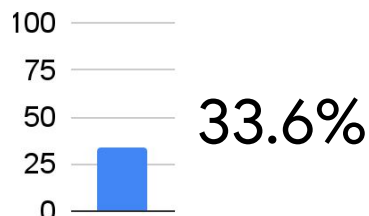
**App bottlenecked: ML Inference**

**Accelerated application processing** can improve tail latency by 68.0%
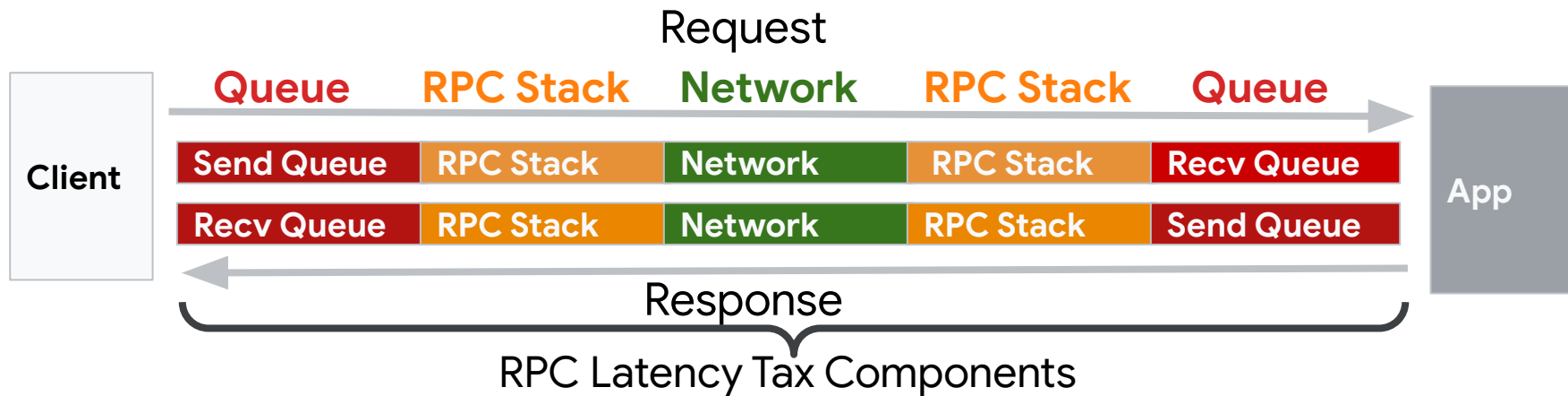
**Queueing bottlenecked: SSD Cache**

**Scheduling or resource management** can improve tail latency by 33.6%

**Takeaway**: There is a potential for a significant reduction to tail latency by **eliminating the variation** caused by the **dominant component**.

Google

# Agenda

- Motivations for studying RPC in the cloud

- RPC Study Results

  ○ What is the **source of RPCs**? Where do they go?

  ○ What is the **timescale** of RPC?

  ○ Which **latency component** affects RPC latency?

  ○ Which **RPC Latency Tax** component is the bottleneck?

  ○ How does **utilization** vary across datacenters?

- Implications from the study

26

Google

# RPC Latency Tax



Request

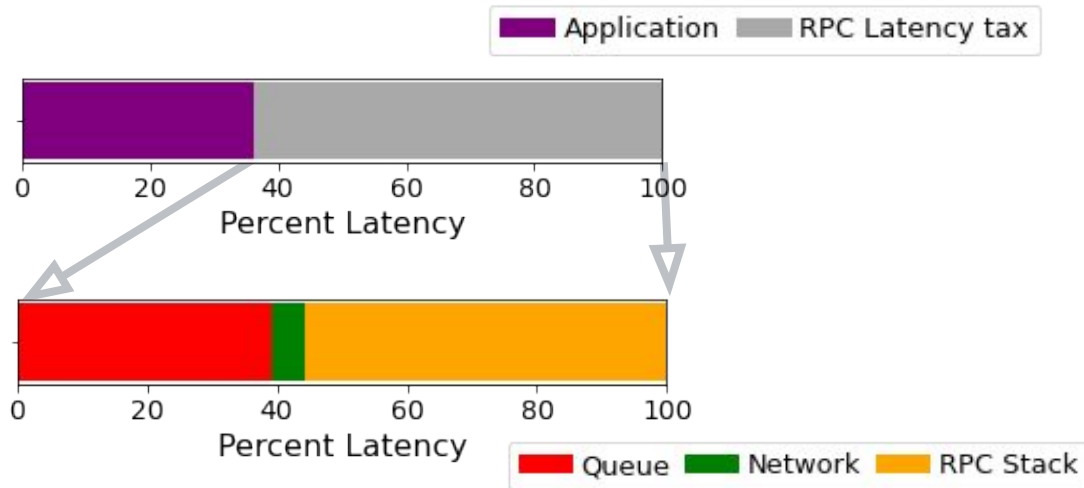| Queue | RPC Stack | Network | RPC Stack | Queue |
|---|---|---|---|---|
| Send Queue | RPC Stack | Network | RPC Stack | Recv Queue |
| Recv Queue | RPC Stack | Network | RPC Stack | Send Queue |

Client

App

Response

RPC Latency Tax Components

We define RPC Latency Tax as the overhead of running application over RPC, all latency components excluding the application processing time

How significant is RPC Latency Tax for **within** and **across clusters**?

Google
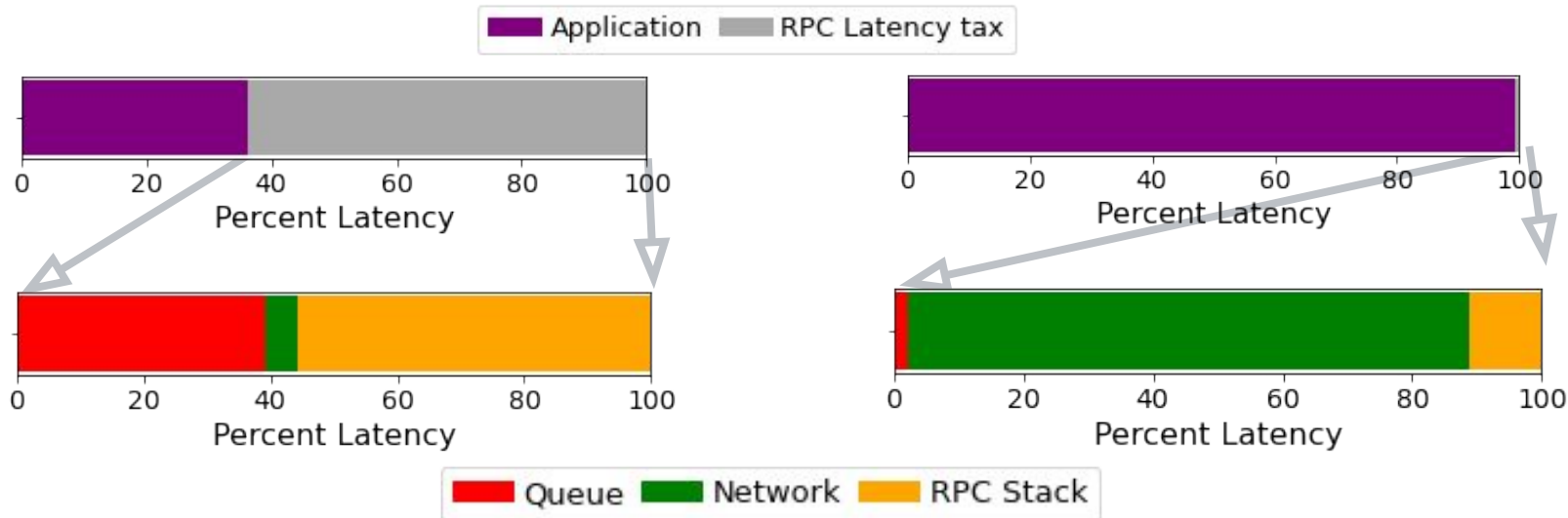
# RPC Latency Tax: Variation at Tail

## Contribution of RPC Latency Tax



**Intra-Cluster P95+**

Google

# RPC Latency Tax: Variation at Tail
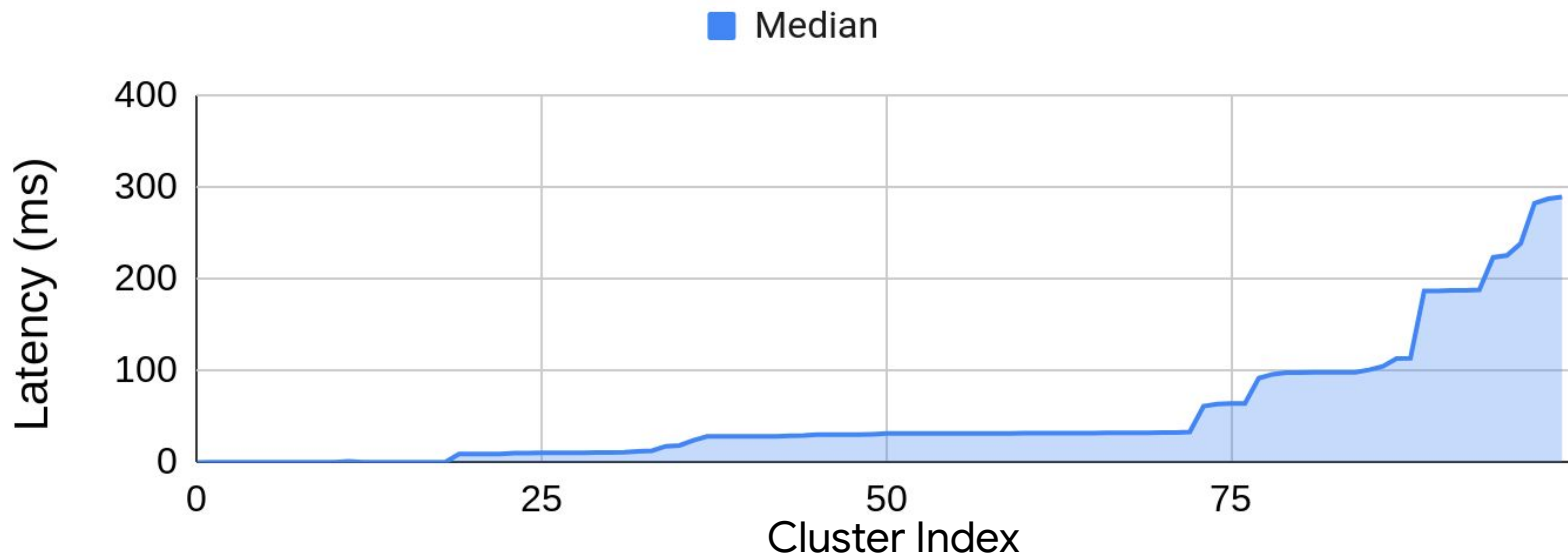
## Contribution of RPC Latency Tax
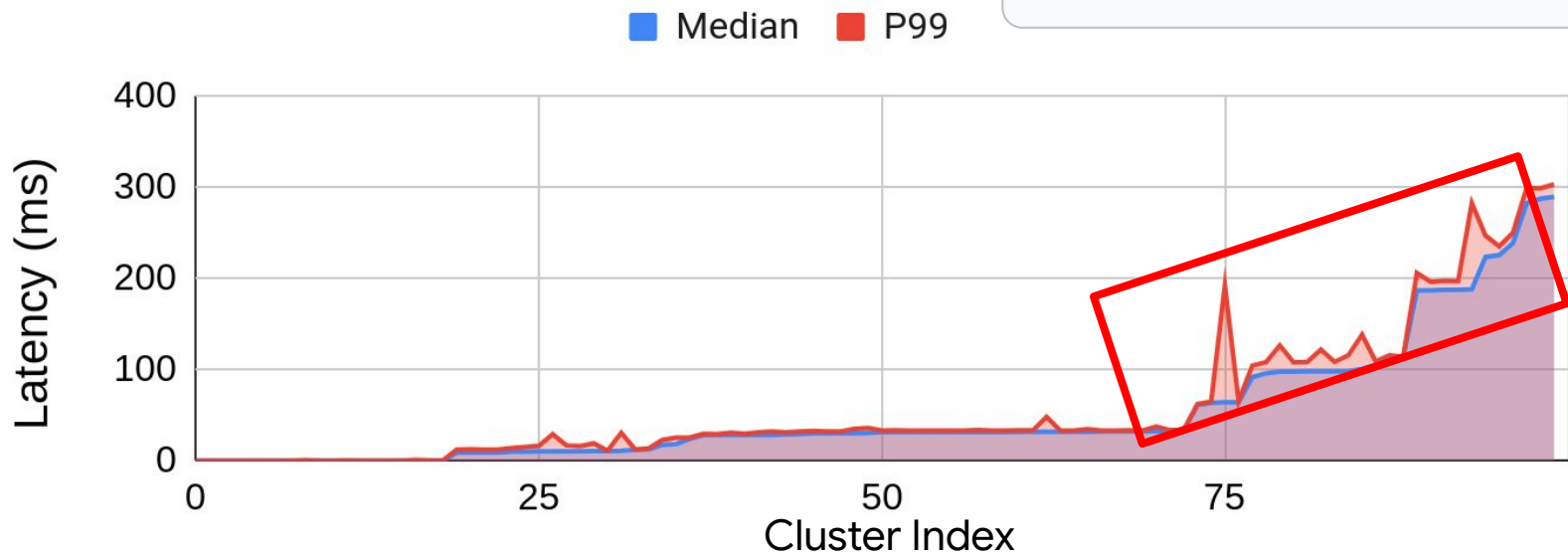


**Intra-Cluster P95+**

**Inter-Cluster P95+**

**Takeaway**: RPC Latency tax is **significant at tail**

Google

# RPC Latency Tax: Inter-Cluster Variation at Tail

Google

# RPC Latency Tax: Inter-Cluster Variation at Tail

Latency spike due to **congestion**



**Takeaway**: **Congestion on WAN** can have an impact on **tail latency across clusters**

Google

# RPC Study Results:
# **Takeaways**

Which **latency component** affects RPC latency?

- **Dominant latency component** is different for each service
  - Optimize RPC stack, queueing or app processing

- **RPC Latency tax is significant at tail**
  - **Queueing matters** for intra-cluster RPCs.
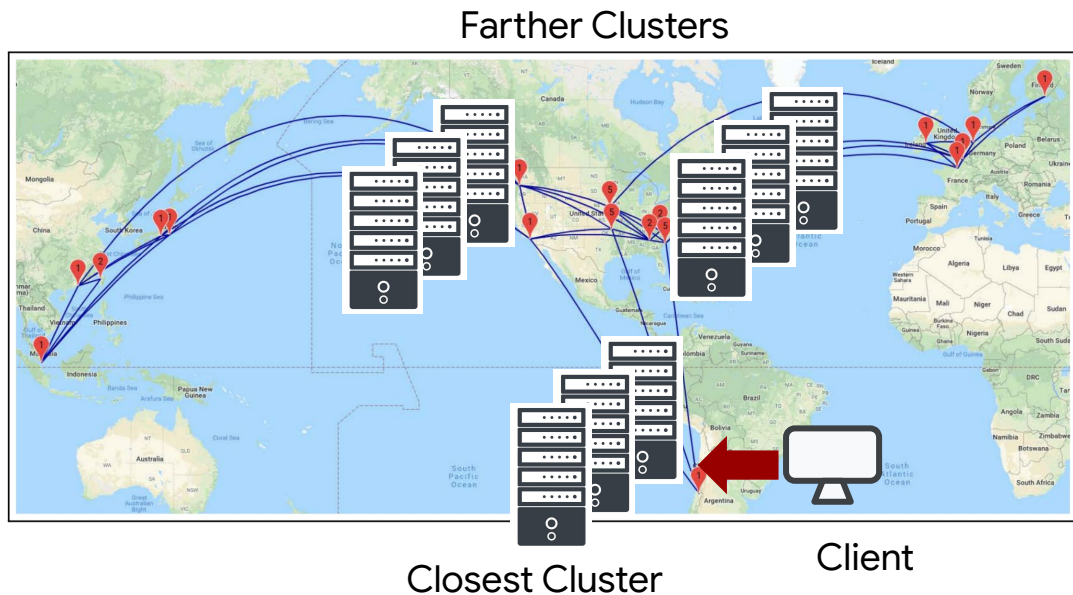  - **WAN congestion** matters for inter-cluster RPCs.

Google

# Agenda

- Motivations for studying RPC in the cloud

- RPC Study Results

  - What is the **source of RPCs**? Where do they go?
  - What is the **timescale** of RPC?
  - Which **latency component** affects RPC latency?
  - Which **RPC Latency Tax** component is the bottleneck?
  - How does **utilization** vary across datacenters?
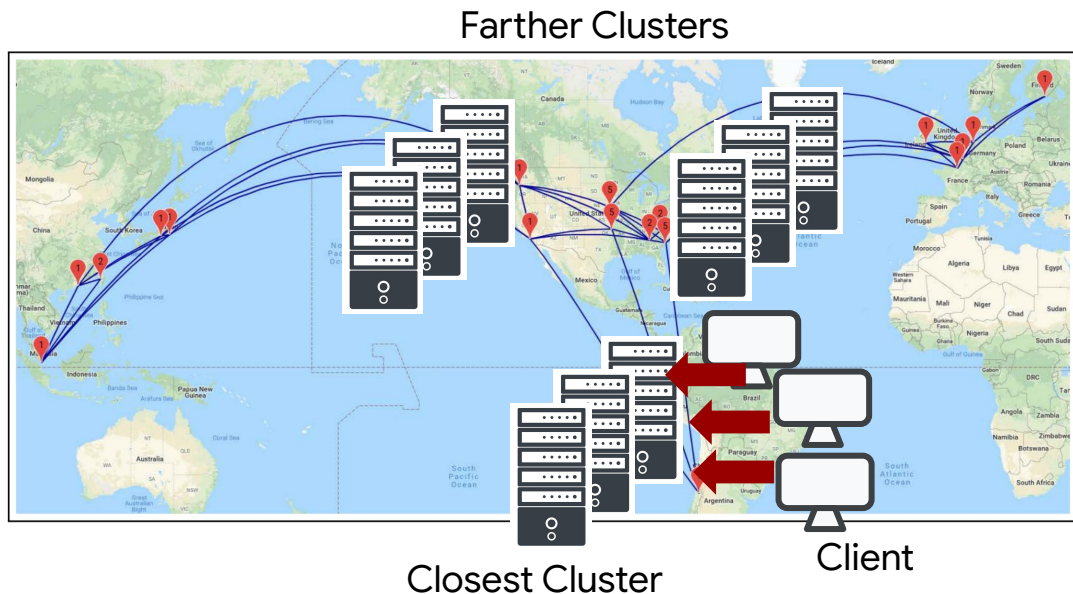- Implications from the study

Google

# Load balancing

To avoid inter-cluster traffic, we could serve in the cluster **closest to the client**



Farther Clusters

Closest Cluster

Client

Google

# Load balancing

To avoid inter-cluster traffic, we could serve in the cluster **closest to the client**

Farther Clusters



Closest Cluster                                    Client

However, serving requests on the closest cluster could unbalance load.

Google

# Load balancing: CPU Utilization Variation

Research question:

Is CPU utilization balanced across **different clusters**?

Is CPU utilization balanced across **different machines** within a cluster?
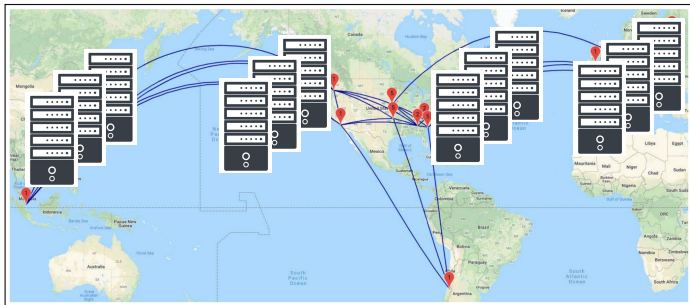
Methodology: Collect CPU Utilization



36

Google

# Load balancing: CPU Utilization Variation

Research question:

Is CPU utilization balanced across **different clusters**?

Is CPU utilization balanced across **different machines** within a cluster?

Methodology: Collect CPU Utilization

**1. Across different clusters**

Google

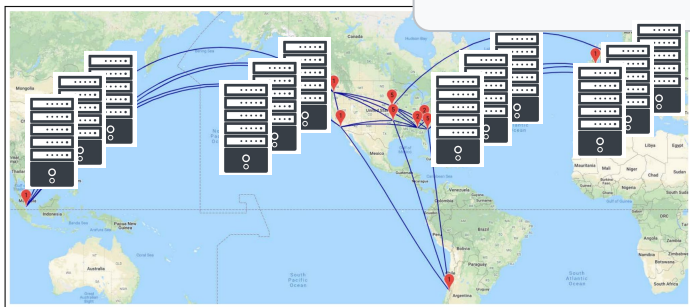# Load balancing: CPU Utilization Variation

Research question:

Is CPU utilization balanced across **different clusters**?

Is CPU utilization balanced across **different machines** within a cluster?

Methodology: Collect CPU Utilization



**1. Across different clusters**

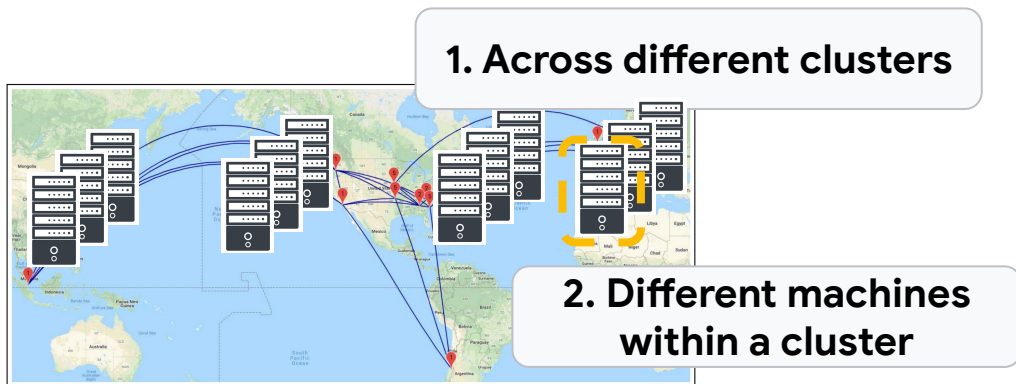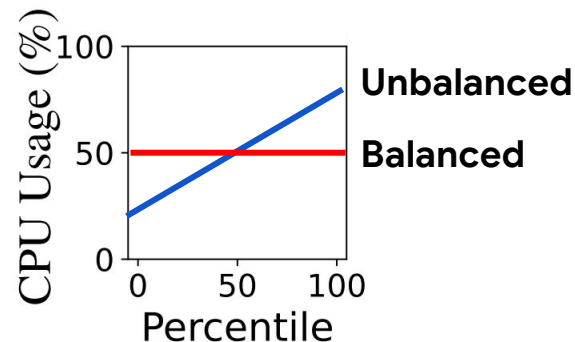**2. Different machines within a cluster**
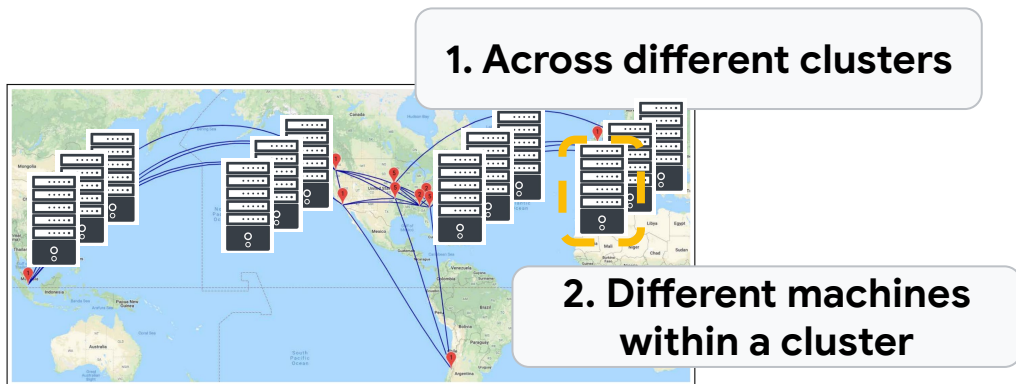
38

Google

# Load balancing: CPU Utilization Variation
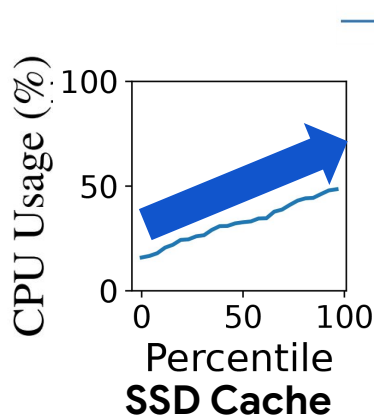
Research question:

Is CPU utilization balanced across **different clusters**?

Is CPU utilization balanced across **different machines** within a cluster?

Methodology: Collect CPU Utilization

**1. Across different clusters**

**2. Different machines within a cluster**

Unbalanced

Balanced

CPU Usage (%)

Percentile

Google

# Load balancing: Cross-Cluster CPU Utilization Variation



Clusters

SSD Cache

Unbalanced

ML Inference

Unbalanced

Unbalanced across clusters

Unbalanced across clusters

**Takeaway**: There are trade-offs between
**network latency** and **load balancing across clusters**.

40

Google

# Load balancing: Same-Cluster CPU Utilization Variation



**Takeaway**: Compute services with **unpredictable** latency are unbalance within a cluster.

# RPC Study Results:
# Takeaways

**How do latency and utilization vary across datacenters?**

- Hard to balance load for services with **varied computation**

- Inter-cluster RPC placement helps load balancing across clusters but **can increase latency**

Google

# Agenda

- Motivations for studying RPC in the cloud

- RPC Study Results

  - What is the **source of RPCs**? Where do they go?

  - What is the **timescale** of RPC?

  - Which **latency component** affects RPC latency?

  - Which **RPC Latency Tax** component is the bottleneck?

  - How does **latency and utilization** vary across datacenters?

- Implications from the study

Google

# Implications

- **Storage data flow optimization is important**
    - Majority of RPC invocations and data transfer are from storage applications
    - Optimizing data movement for storage RPCs can significantly improve resource efficiency
- **Millisecond, not just microsecond timescales**
    - Most RPCs operate in millisecond scale
    - Reducing CPU utilization can be more beneficial than saving a few microseconds
- **Host queuing matters**
    - Client & Server queuing latency are major contributors to the tail latency
    - Improving scheduling and placement is important
- **RPC Latency Tax is significant at tail**
    - Need to optimize RPC overhead at the tail requests
- **Load-balancing needs to account for latency**
    - Need research on predicting latency for RPCs with  varied computational needs
    - Scheduling  across cluster needs to co-optimize for latency and load-balancing

Google

# Implications

- **Storage data flow optimization is important**
  - Majority of RPC invocations and data transfer are from storage applications
  - Optimizing data movement for storage RPCs can significantly improve resource efficiency
- **Millisecond, not just microsecond timescales**
  - Most RPCs operate in millisecond scale
  - Reducing CPU utilization can be more beneficial than saving a few microseconds
- **Host queuing matters**
  - Client & Server queuing latency are major contributors to the tail latency
  - Improving scheduling and placement is important
- **RPC Latency Tax is significant at tail**
  - Need to optimize RPC overhead at the tail requests
- **Load-balancing needs to account for latency**
  - Need research on predicting latency for RPCs with  varied computational needs
  - Scheduling  across cluster needs to co-optimize for latency and load-balancing

45

Google

# Implications

- **Storage data flow optimization is important**
  - Majority of RPC invocations and data transfer are from storage applications
  - Optimizing data movement for storage RPCs can significantly improve resource efficiency
- **Millisecond, not just microsecond timescales**
  - Most RPCs operate in millisecond scale
  - Reducing CPU utilization can be more beneficial than saving a few microseconds
- **Host queuing matters**
  - Client & Server queuing latency are major contributors to the tail latency
  - Improving scheduling and placement is important
- **RPC Latency Tax is significant at tail**
  - Need to optimize RPC overhead at the tail requests
- **Load-balancing needs to account for latency**
  - Need research on predicting latency for RPCs with  varied computational needs
  - Scheduling  across cluster needs to co-optimize for latency and load-balancing

46

Google

# Implications

- **Storage data flow optimization is important**
  - Majority of RPC invocations and data transfer are from storage applications
  - Optimizing data movement for storage RPCs can significantly improve resource efficiency
- **Millisecond, not just microsecond timescales**
  - Most RPCs operate in millisecond scale
  - Reducing CPU utilization can be more beneficial than saving a few microseconds
- **Host queuing matters**
  - Client & Server queuing latency are major contributors to the tail latency
  - Improving scheduling and placement is important
- **RPC Latency Tax is significant at tail**
  - Need to optimize RPC overhead at the tail requests
- **Load-balancing needs to account for latency**
  - Need research on predicting latency for RPCs with varied computational needs
  - Scheduling across cluster needs to co-optimize for latency and load-balancing

Google

# Implications

- **Storage data flow optimization is important**
  - Majority of RPC invocations and data transfer are from storage applications
  - Optimizing data movement for storage RPCs can significantly improve resource efficiency
- **Millisecond, not just microsecond timescales**
  - Most RPCs operate in millisecond scale
  - Reducing CPU utilization can be more beneficial than saving a few microseconds
- **Host queuing matters**
  - Client & Server queuing latency are major contributors to the tail latency
  - Improving scheduling and placement is important
- **RPC Latency Tax is significant at tail**
  - Need to optimize RPC overhead at the tail requests
- **Load-balancing needs to account for latency**
  - Need research on predicting latency for RPCs with varied computational needs
  - Scheduling across cluster needs to co-optimize for latency and load-balancing
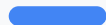
Google

# Conclusion

**RPC Study:**
- First ever study on **Google's fleet-wide RPC characteristics**
- **722 billion** RPC traces over **~2 years** running on **100 production clusters**
- Provides insights on the characteristics of Google's geo-distributed **internal services**

**Key contributions and findings:**
- Storage **data flow optimization** is important
- **Millisecond**-scale RPCs are common – need to **balance** CPU **utilization** vs. **latency**
- **RPC Queuing matters** – need to improve **scheduling** and **load balancing**
- **RPC Latency Tax is significant at tail** – need better optimization within and across clusters to reduce **tail latency variation**
- **Load-balancing needs to account for latency** - co-optimize latency and utilization

Our measurements can influence **future RPC research**.

Google

Google

# A Cloud-Scale Characterization of Remote Procedure Calls

**Korakit Seemakhupt**, Brent Stephens, Samira Khan, Sihang Liu, Hassan Wessel,

Soheil Hassas Yeganeh, Alex C. Snoeren, Arvind Krishnamurthy, David Culler, Henry M. Levy