# CARD: A Congestion-Aware Request Dispatching Scheme for Replicated Metadata Server Cluster

Shangming Cai, Dongsheng Wang,

Zhanye Wang and Haixia Wang
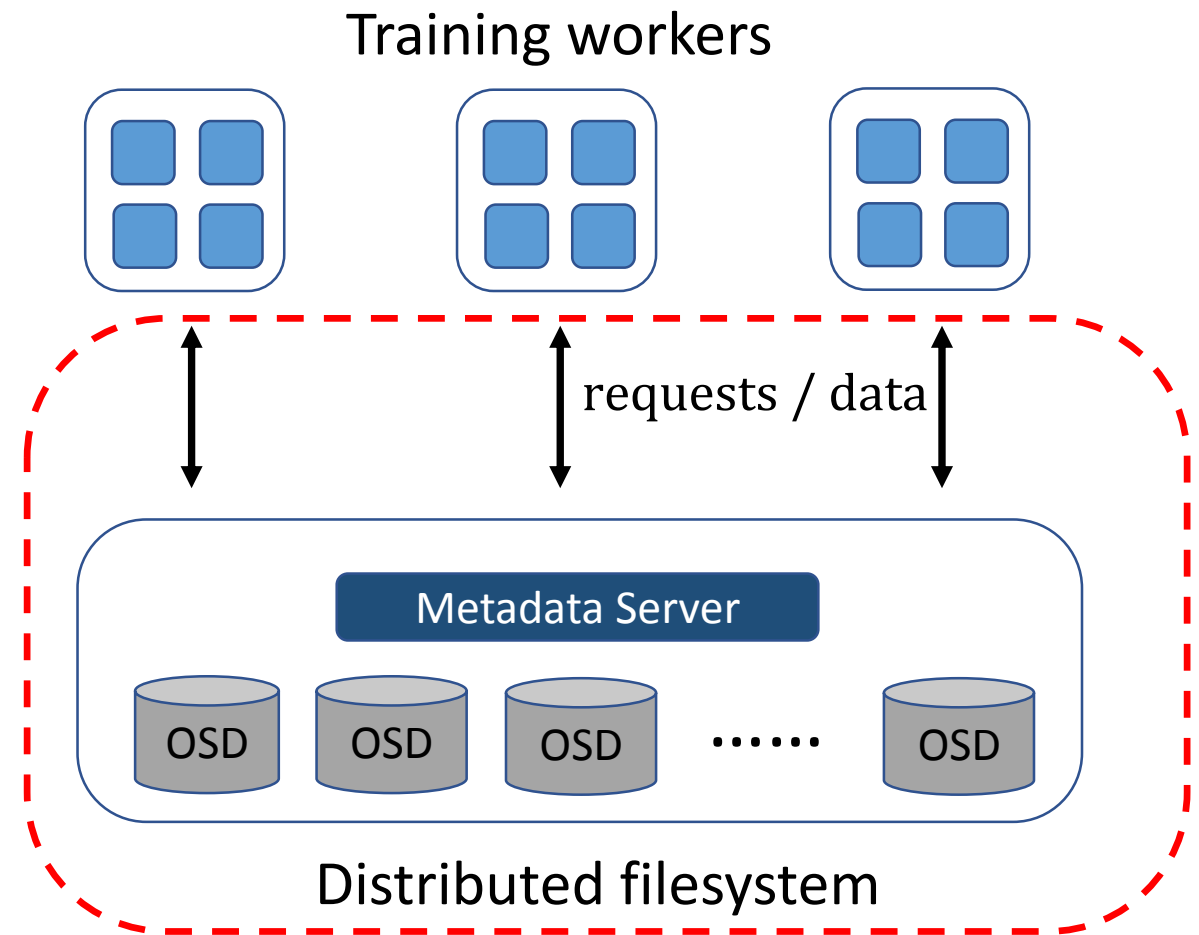
*Tsinghua University*

# Background: Massive-scale ML in product environments

- Datasets updated hourly or daily

  - data collected and stored in an HDFS-like distributed filesystem
  - periodically offline training for online inference

- Challenges of the data-reader pipeline while training

  - extremely heavy read workloads: millions to billions of files per epoch
  - random access pattern: up-level shuffling for convergence speed
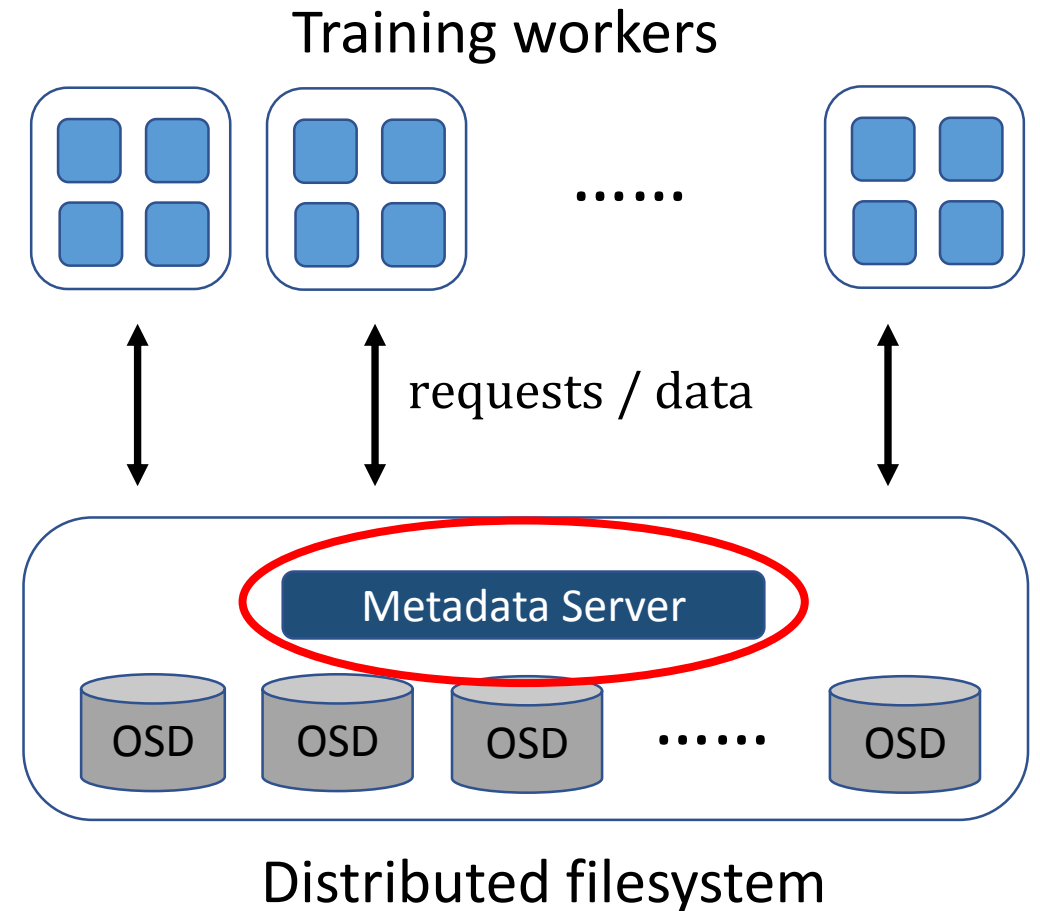
# Background: Massive-scale ML in product environments

- Workers interact with a DFS

- Metadata request
  -> metadata server (MDS)

- File I/O
  -> object storage devices (OSD)

Training workers

requests / data

Metadata Server

OSD  OSD  OSD  ......  OSD

Distributed filesystem

3

# When the number of training workers grows...

- Extremely stressed workloads

- Metadata access step bottlenecks the data-reader pipeline

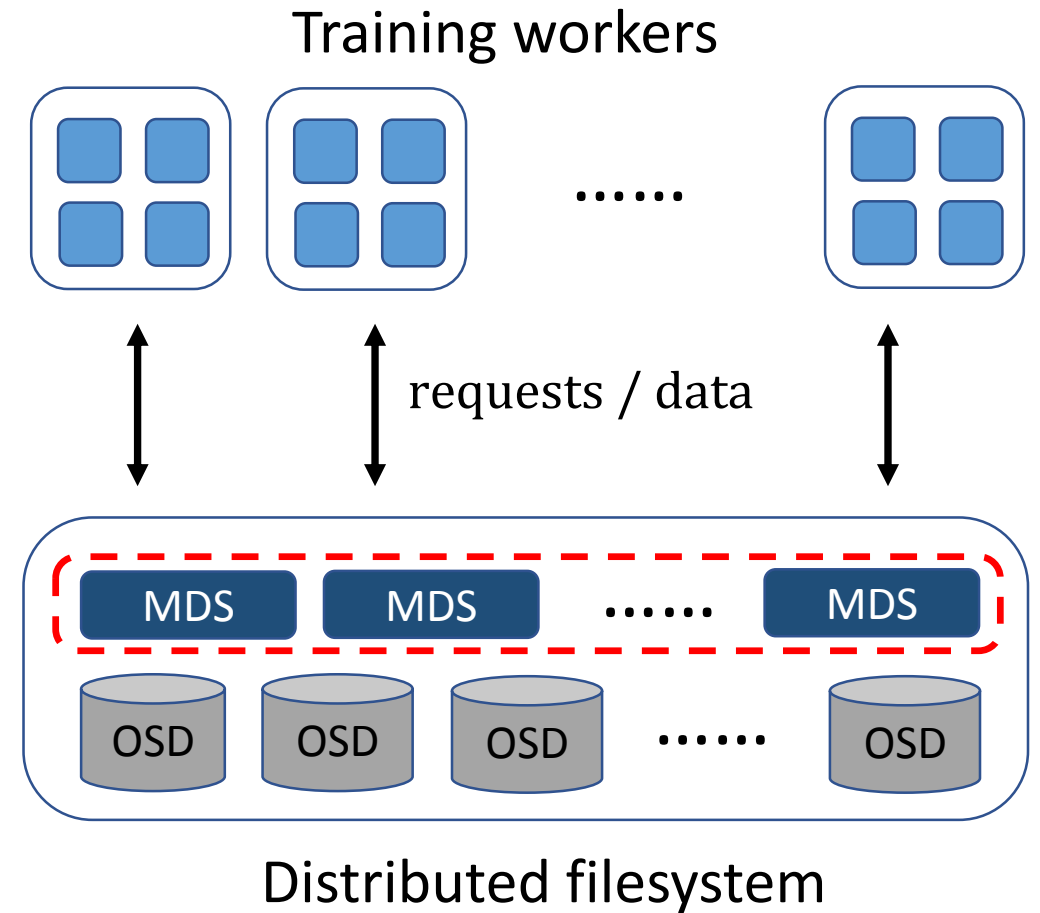- Potential single point of failure on MDS

Training workers



requests / data

Distributed filesystem

# Typical industrial response: Scaling out likewise

- Concerns to be addressed:

  - Cost-effectiveness

  - Scalability

  - Run-time stability

Training workers



requests / data

Distributed filesystem
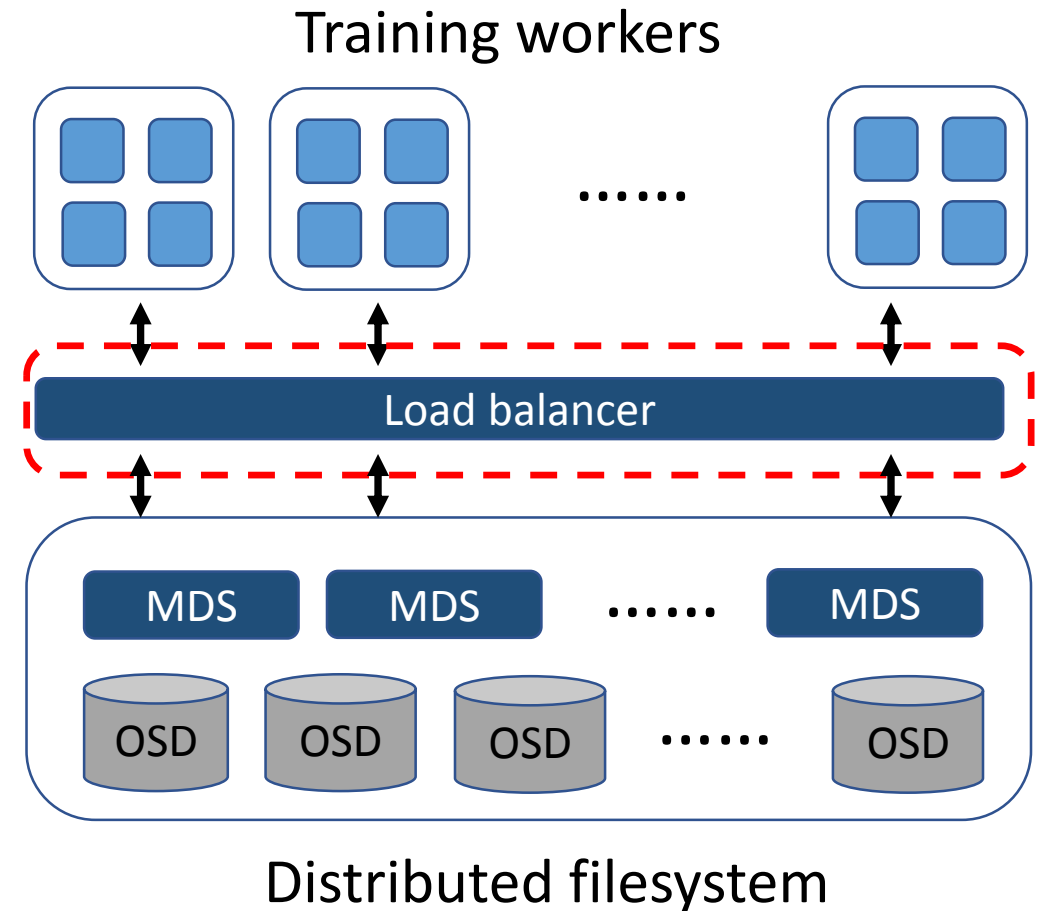
# To achieve load-balance...

- A middle layer load-balancer
  - Pros:
    - good global load balancing
    - more features are optional
  - Cons:
    - load-balancer is stressed
    - reintroduce a potential single point of failure
    - not cost-effective

Training workers



Load balancer

MDS    MDS    ......    MDS

OSD    OSD    OSD    ......    OSD

Distributed filesystem
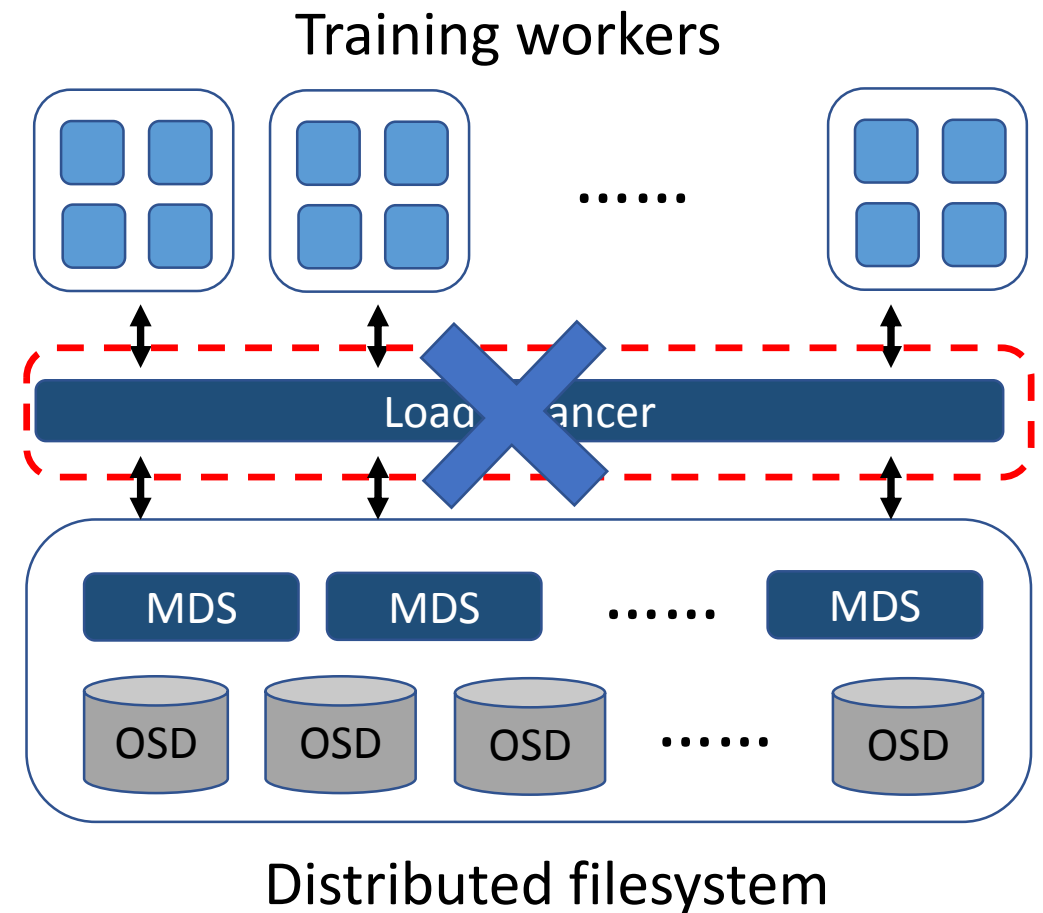
# To achieve load-balance…

- A middle layer load-balancer
  - Pros:
    - good global load balancing
    - more features are optional
  - Cons:
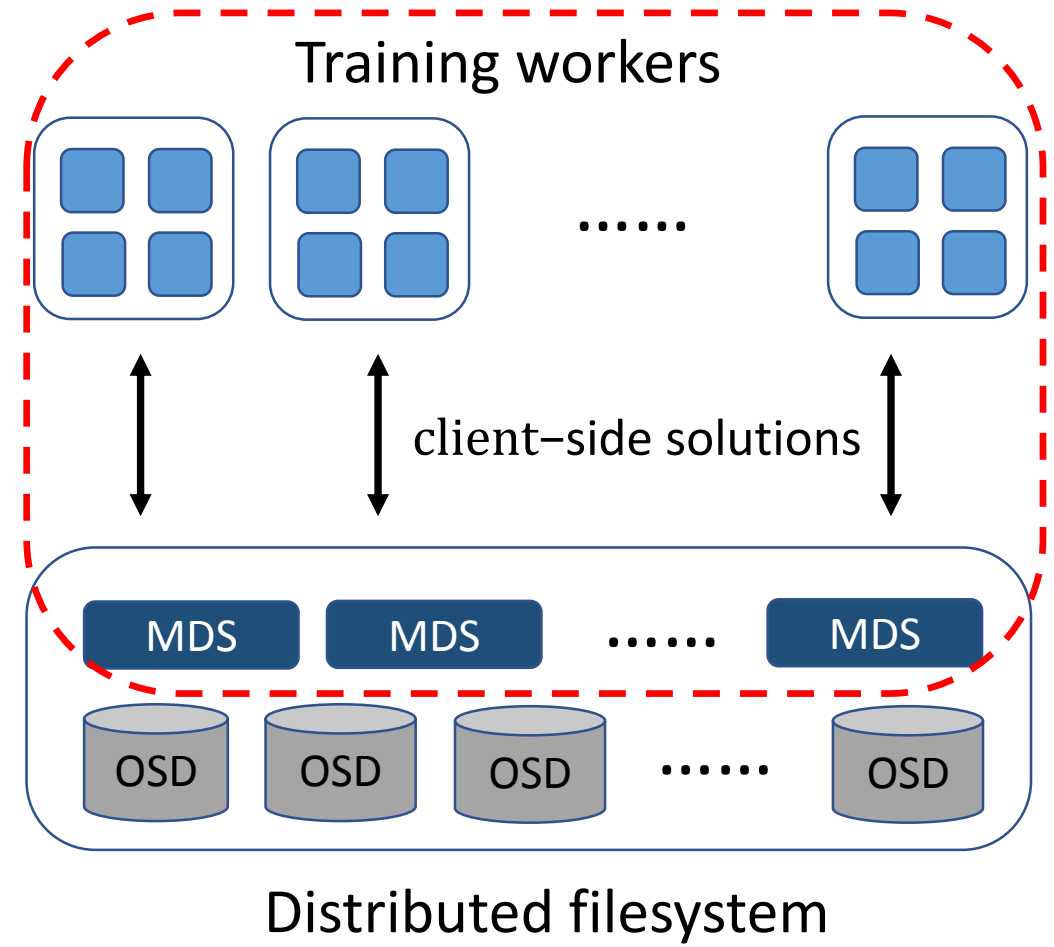    - load-balancer is stressed
    - reintroduce a potential single point of failure
    - not cost-effective

Training workers

Load Balancer

MDS    MDS    ……    MDS

OSD    OSD    OSD    ……    OSD

Distributed filesystem

# Try client-side solutions

- Easy to implement

- Cost-effective

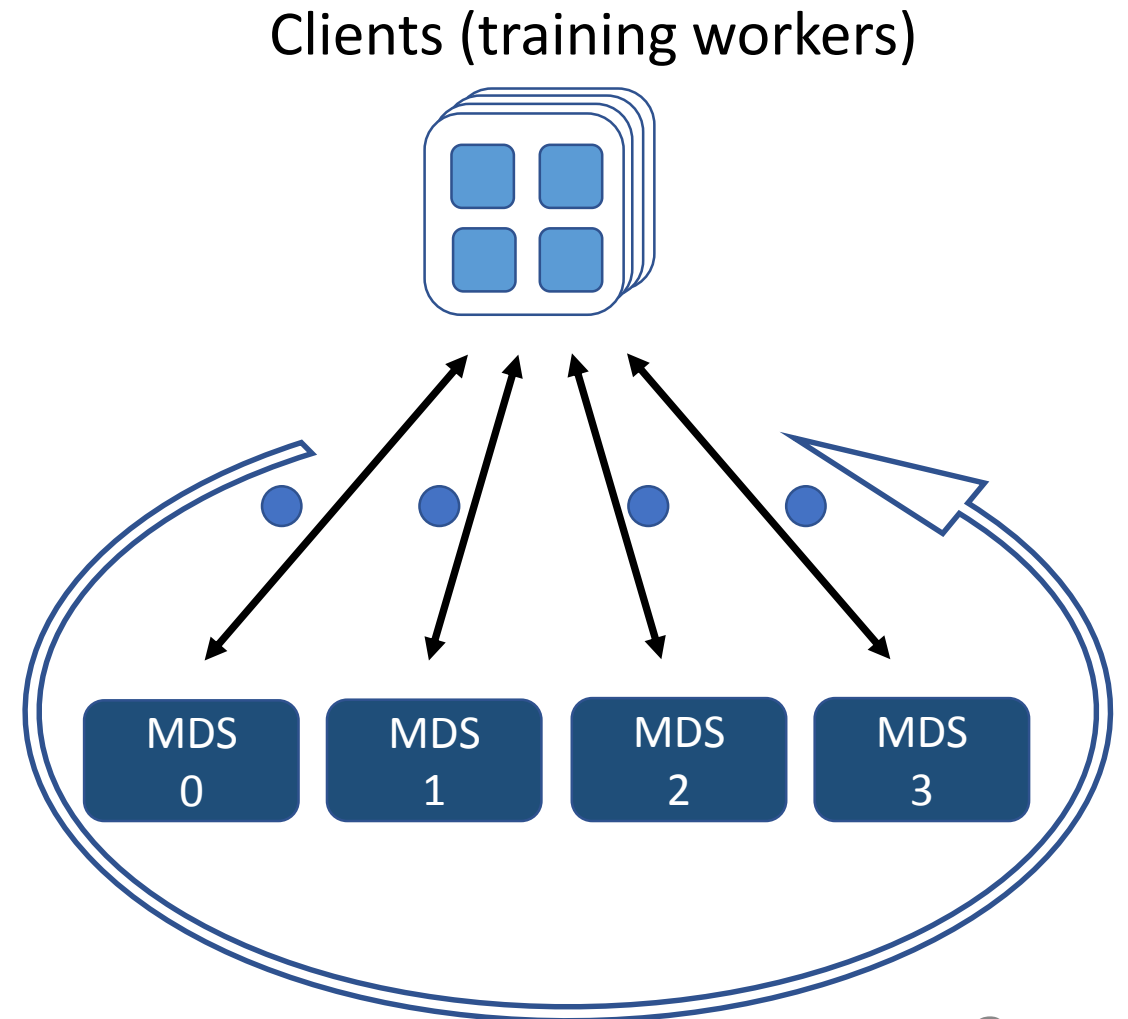# Client-side solution: Round-Robin

- **Round-Robin**

  - Pros:
    - simple yet effective in homogeneous environments
  - Cons:
    - inflexible and inefficient in shifting or heterogeneous environments

Clients (training workers)

MDS 0    MDS 1    MDS 2    MDS 3

# Client-side solution: Heuristic selection

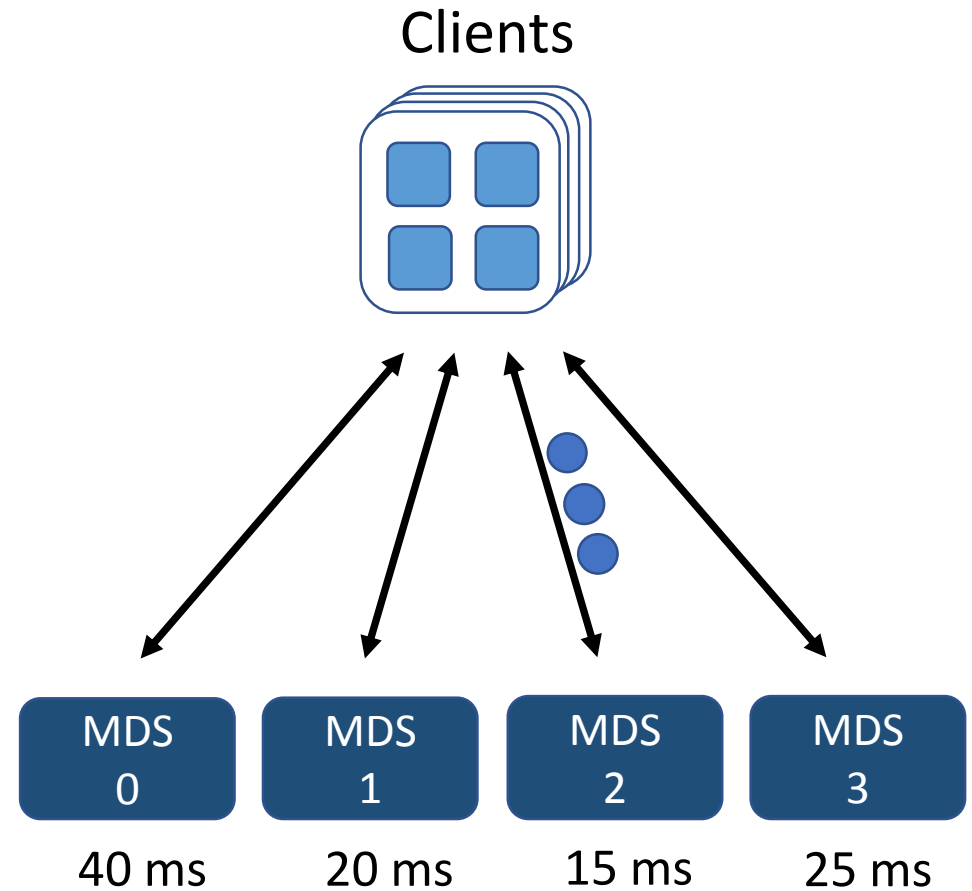- **Heuristic selection**
  - e.g., prefer lowest MART (moving average of response time)

  - Pros:
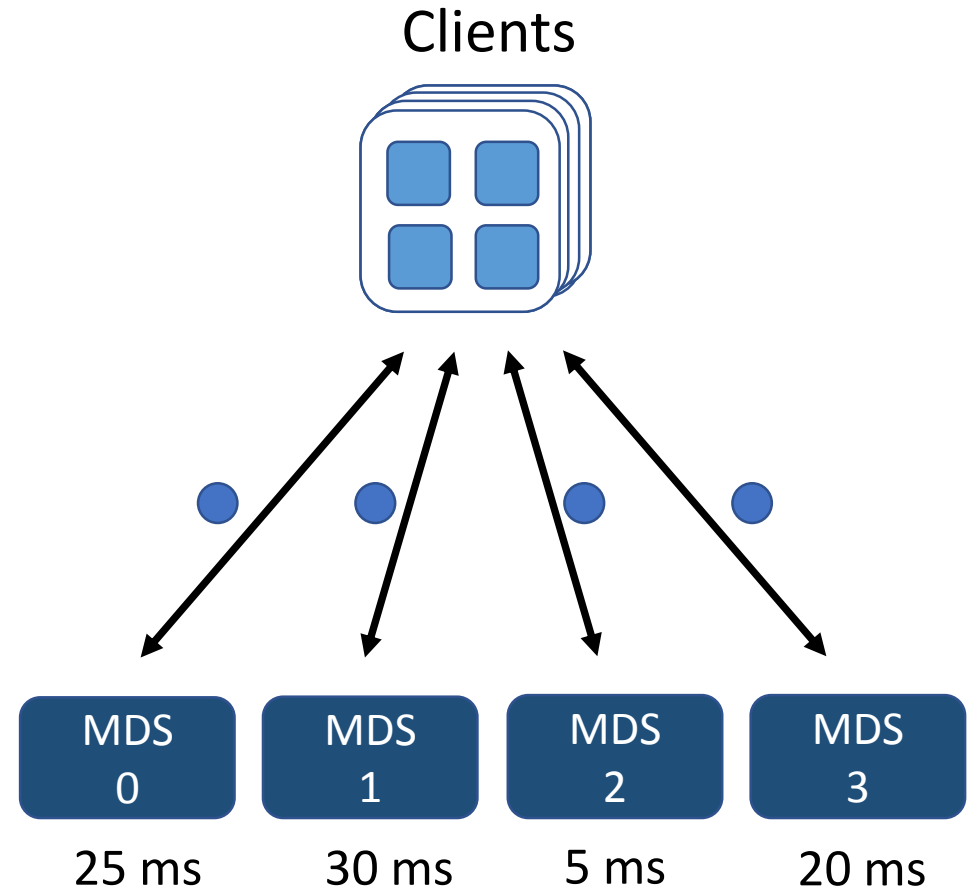    - effective when facing light-weight workloads
  - Cons:
    - cause herd-behavior and load-oscillations

Clients

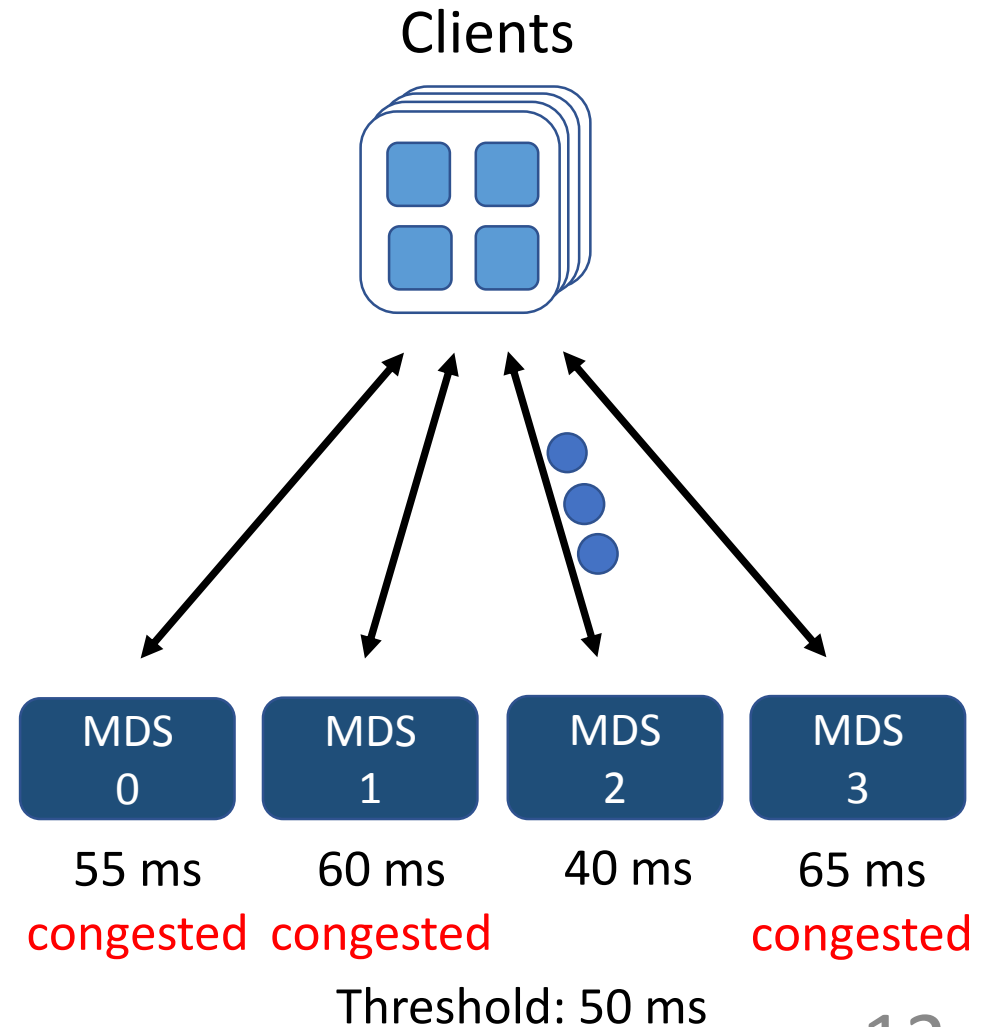| MDS 0 | MDS 1 | MDS 2 | MDS 3 |
|-------|-------|-------|-------|
| 40 ms | 20 ms | 15 ms | 25 ms |

10

# Client-side solution: Round-Robin with Throttling

- Round-Robin with throttling
  - e.g., LADS, preset a MART threshold to mark servers as congested

  - Light-weight workloads
    - = Round-Robin

Clients

| MDS 0 | MDS 1 | MDS 2 | MDS 3 |
|:---:|:---:|:---:|:---:|
| 25 ms | 30 ms | 5 ms | 20 ms |

Threshold: 50 ms

11

# Client-side solution: Round-Robin with Throttling

- Round-Robin with throttling
  - e.g., LADS, preset a MART threshold to mark servers as congested

  - Light-weight workloads
    - = Round-Robin

- Heavy workloads
  - = Heuristic selection
  - herd-behavior and load-oscillations remain

Clients

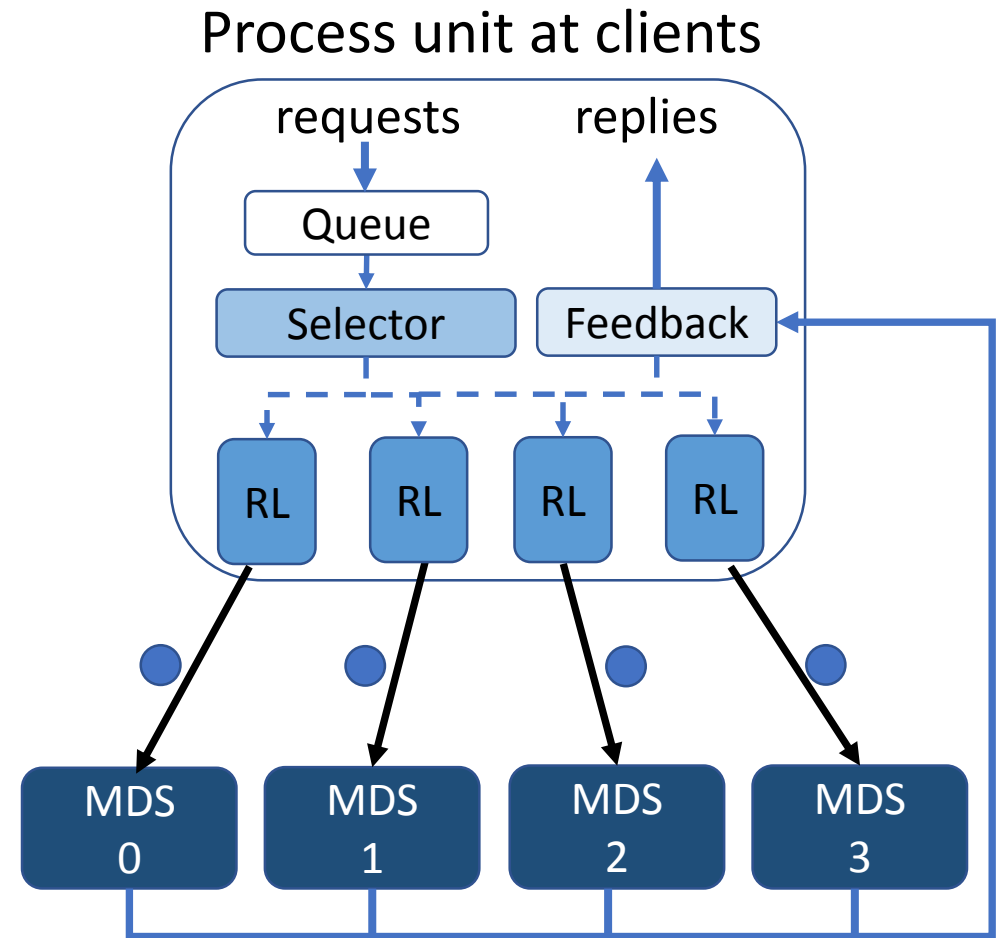| MDS 0 | MDS 1 | MDS 2 | MDS 3 |
|-------|-------|-------|-------|
| 55 ms | 60 ms | 40 ms | 65 ms |
| congested | congested | | congested |

Threshold: 50 ms

12

# CARD: Congestion-Aware Request Dispatching scheme

- Core idea: Round-Robin with adaptive rate-control
  - inspired by CUBIC for TCP protocol
  - counting-based implementation
  - no extra info required from servers

- Light-weight workloads
  - = Round-Robin

- Heavy workloads
  - redirect requests from overloaded MDS to underloaded MDS
  - suppress upcoming requests: if and only if all servers are overloaded
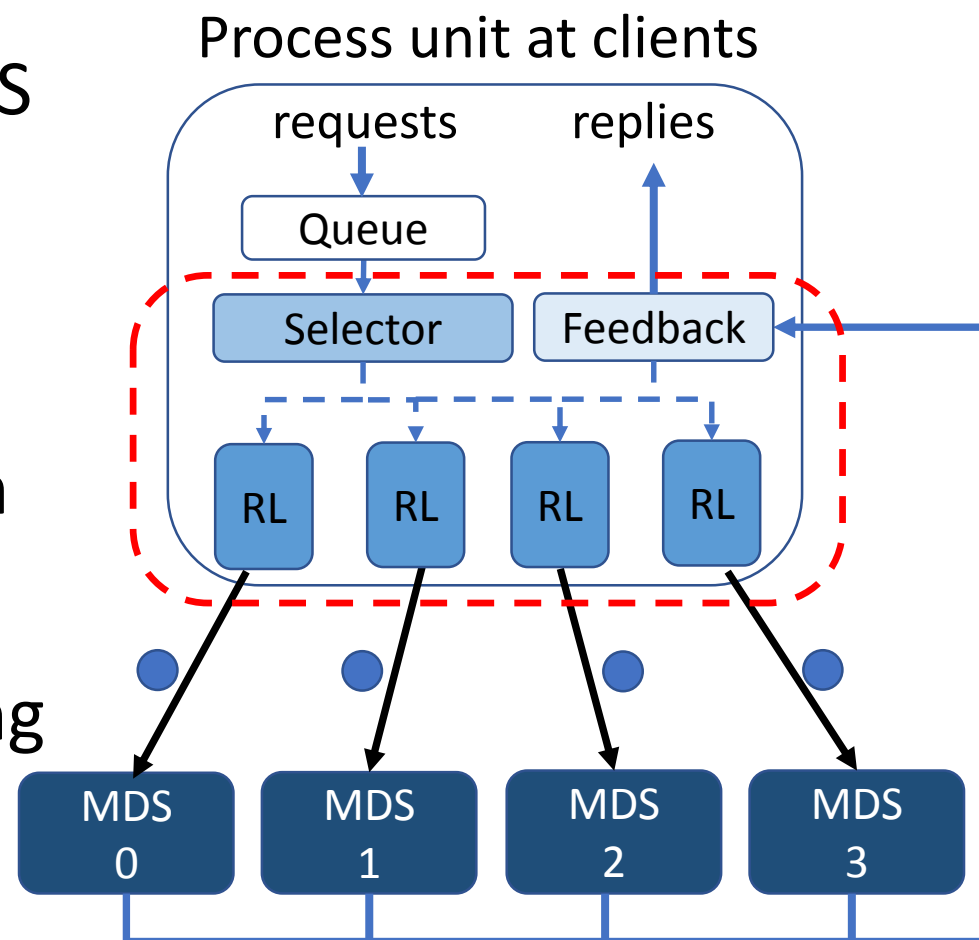
# Congestion-aware rate-control mechanism

- Queue: place pending requests

- Selector: Round-Robin dispatching

- Rate-limiter: rate-control module
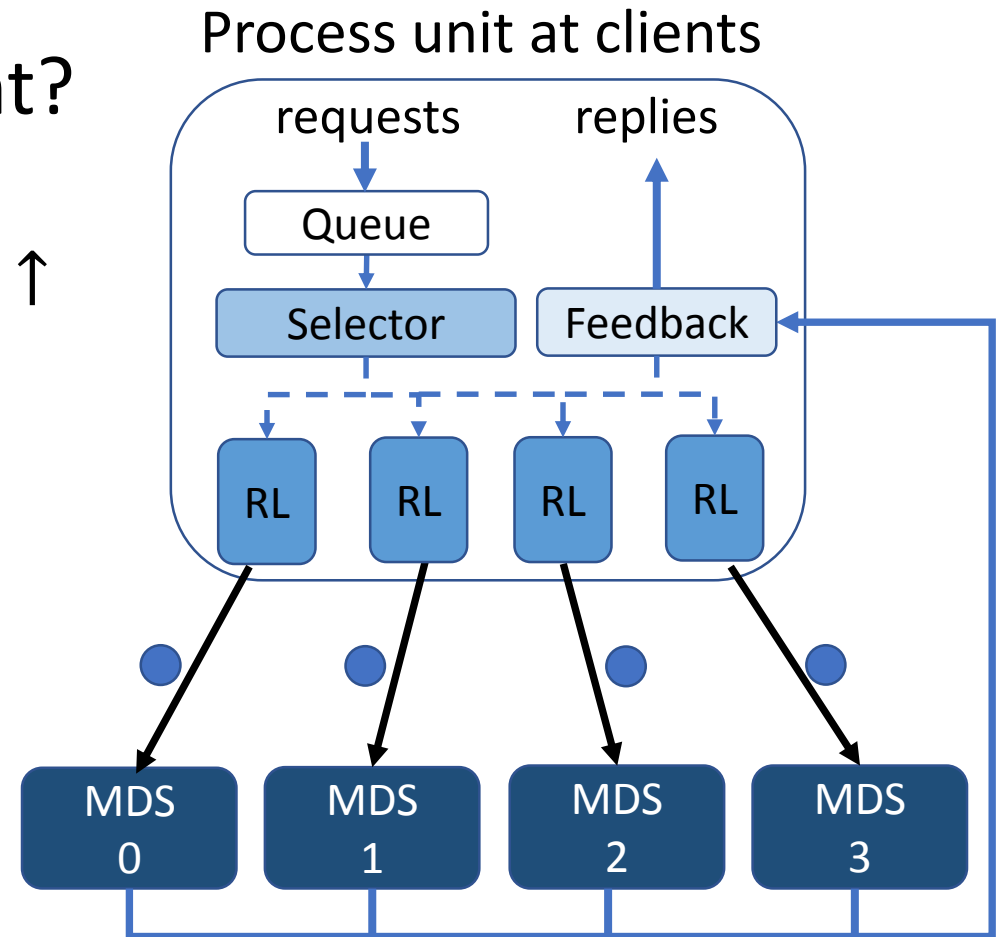
- Feedback: process feedbacks and forward replies

Process unit at clients

requests    replies

Queue

Selector    Feedback

RL    RL    RL    RL

MDS 0    MDS 1    MDS 2    MDS 3

14

# Congestion-aware rate-control mechanism

- Restrict requests routed to each MDS per $\delta$ time window

- Gradually increase the restriction according to a cubic growth function

- Feedback module computes receiving rates after each time window and forwards to RLs
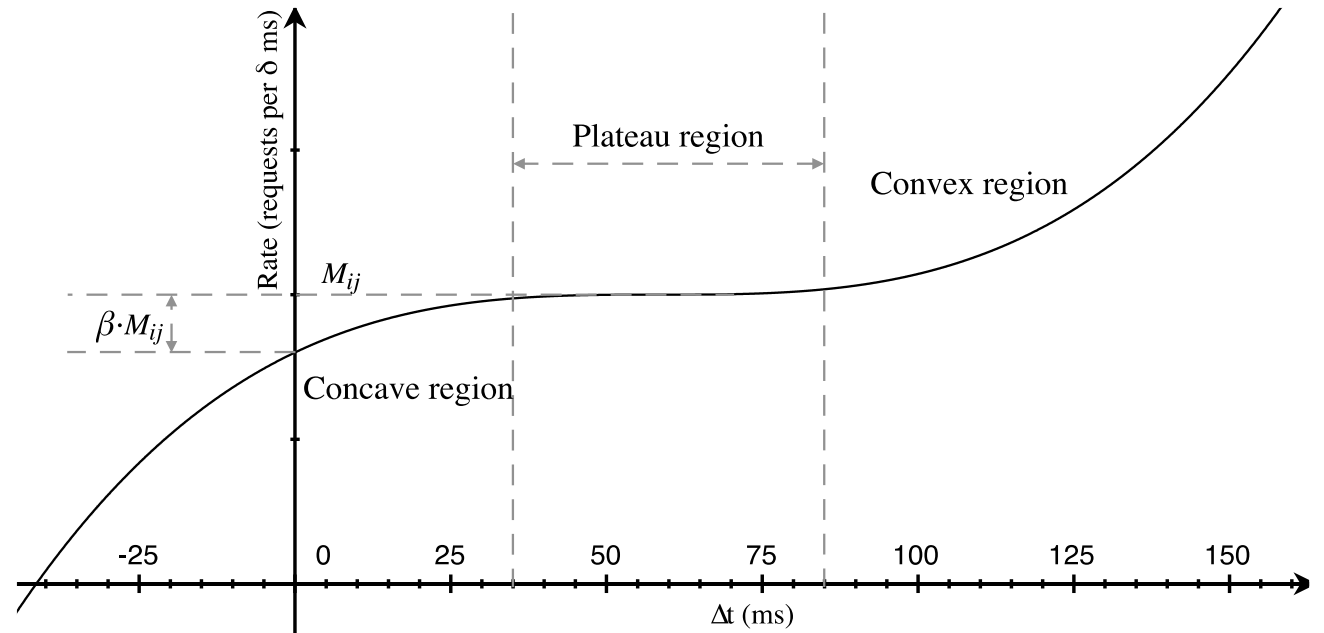
Process unit at clients

requests          replies

Queue

Selector          Feedback

RL    RL    RL    RL

MDS 0    MDS 1    MDS 2    MDS 3

# Congestion-aware rate-control mechanism

- How to identify a congestion event?
  - sending rate > receiving rate
  - elapsed time since last sending rate ↑ event > $\lambda$ (a hysteresis period )

- What to do then?
  - record current sending rate as saturated sending rate
  - reduce current sending rate

Process unit at clients

requests    replies

Queue

Selector    Feedback

RL    RL    RL    RL

MDS 0    MDS 1    MDS 2    MDS 3

# The cubic growth function for the rate-control

- $\Delta t$: elapsed time since the last congestion event

- $M_{ij}$ : saturated sending rate

  - Changed to current sending rate adaptively whenever a congestion event happens

  - Then, current sending rate reduced to $(1 - \beta) \cdot M_{ij}$, and start to grow all over again accordingly
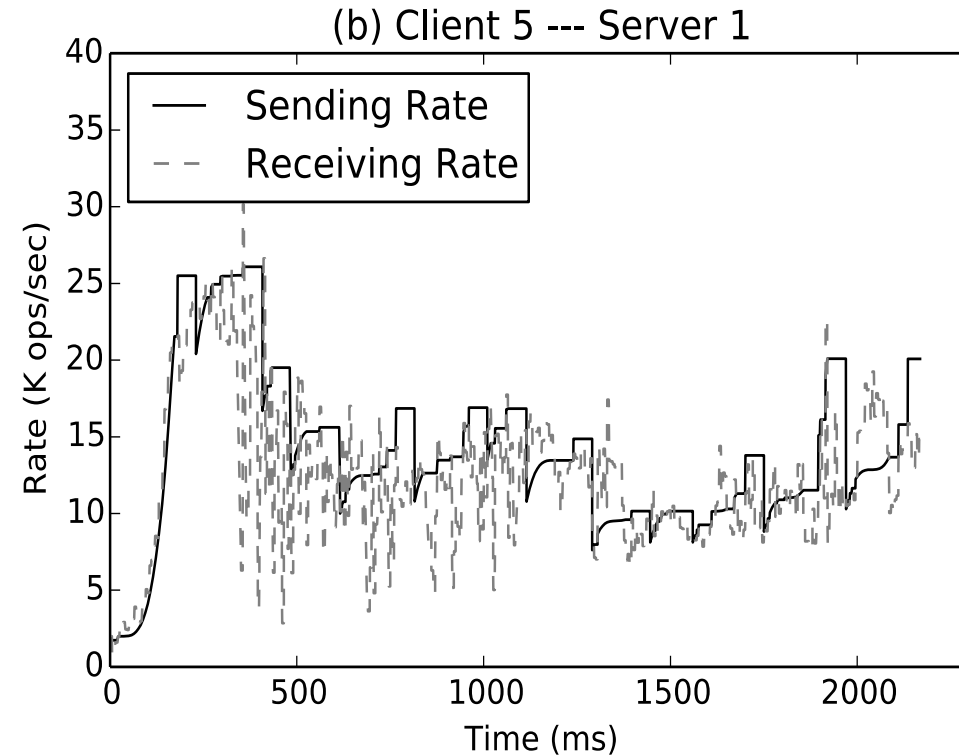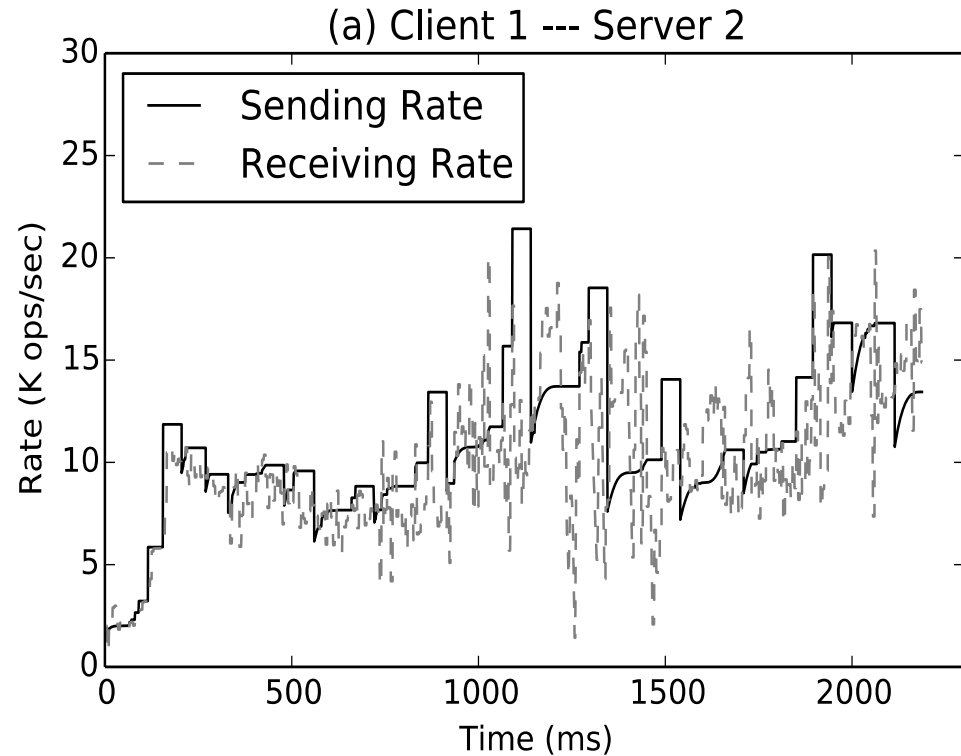
# Evaluation setup

- We implemented a prototype RMSC for simulation purposes
- Up to 8 servers to measure system scalability
- Crafted descending setup for heterogeneous experiments
- 10 clients run on separate machines launching request with Poisson arrivals
- $\delta$ = 5 ms, $\lambda$ = 10 ms, $\beta$ =0.20
- To compare against CARD, we implemented aforementioned Round-Robin, MART and LADS as well
- Refer to the paper for more setup details
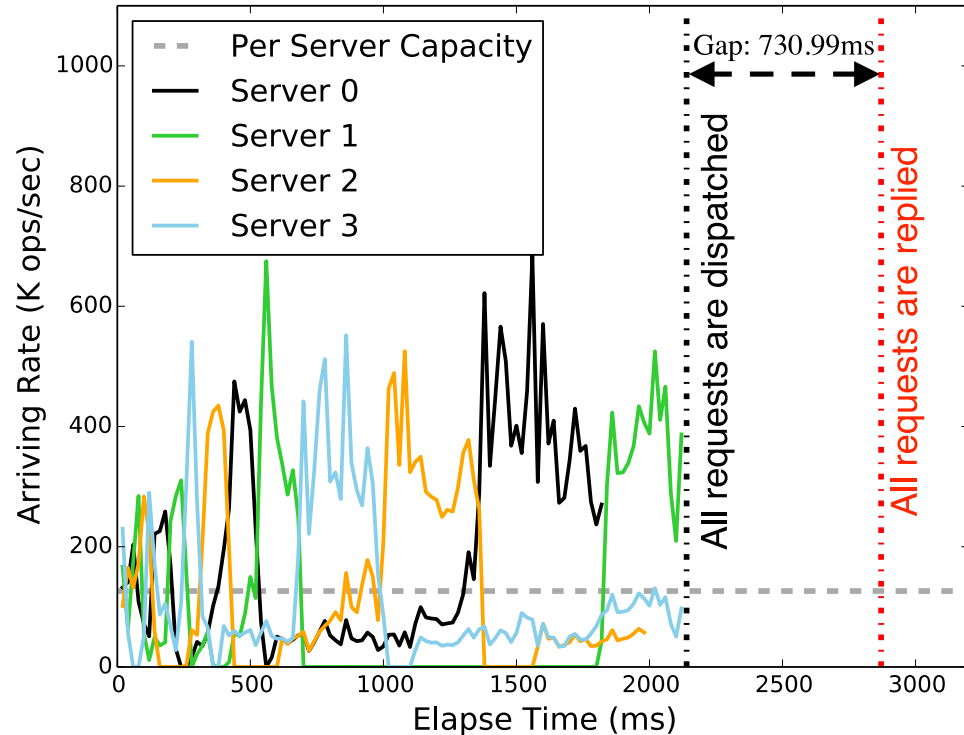
# Evaluation highlights

- Do CARD's rate-control mechanism work as expected?

    - Yes, the rate-control process is effective and adaptive

    - Loads among servers are balanced under heavy workloads

- Can CARD achieve better scalability?

    - In homogeneous clusters: CARD $\approx$ Round-Robin > other strategies

    - In heterogeneous clusters: Yes, CARD > other strategies

# Examples of the rate-control procedure



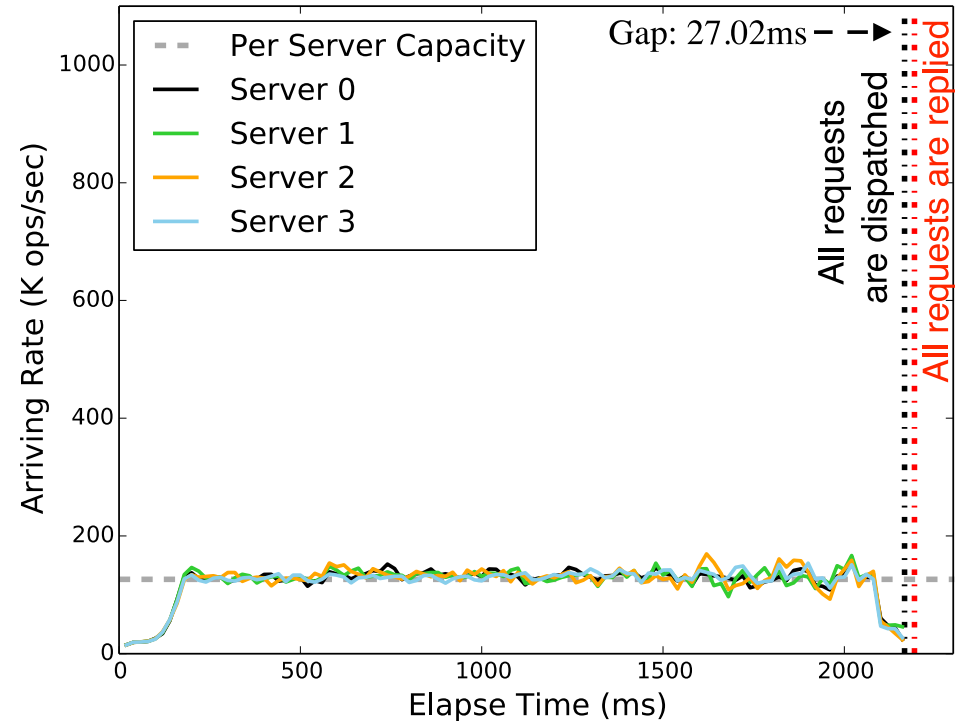(a) Client 1 --- Server 2

(b) Client 5 --- Server 1

The sending rate from each client to each server is adjusted adaptively according to the receiving rate

20
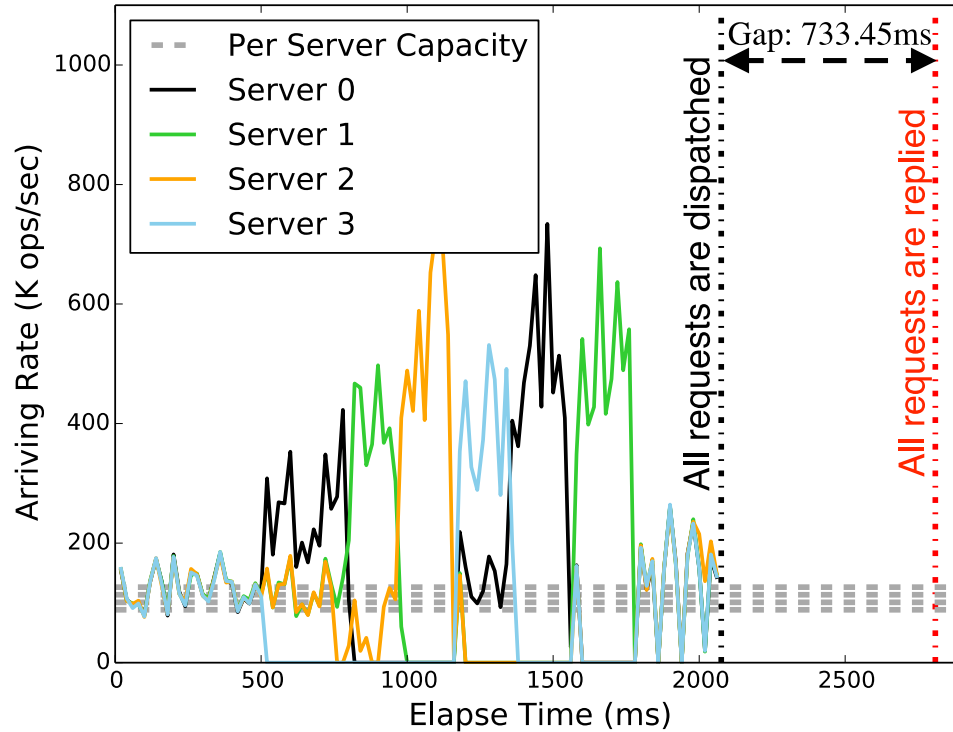
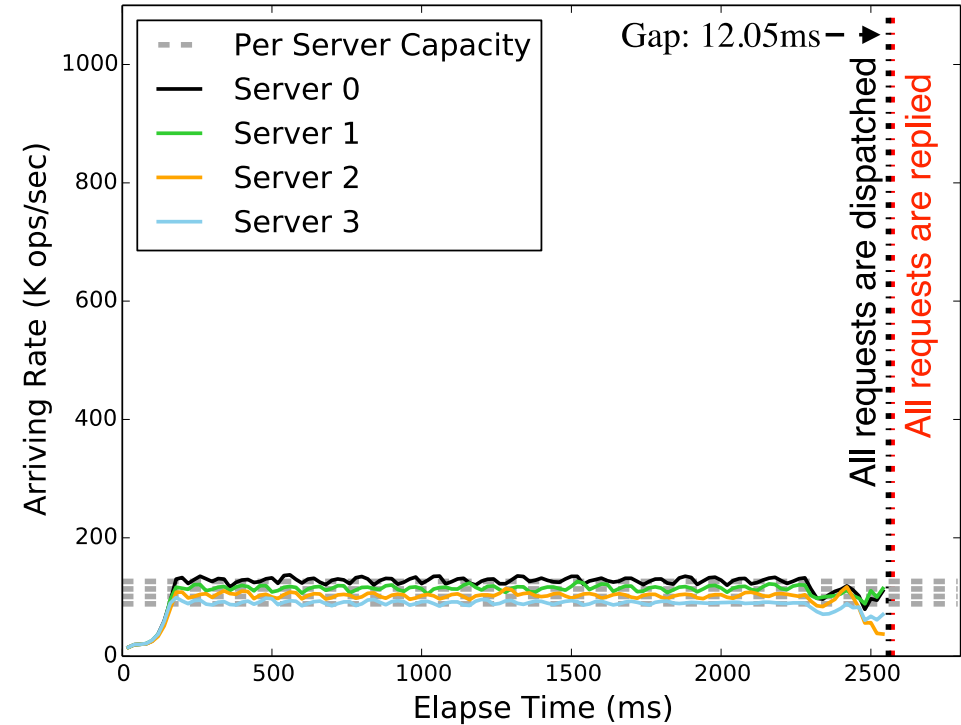# Overall arriving rates in the homogeneous cluster



MART

CARD

1) Heuristic selections cause severe herd behavior and load-oscillations

2) A data loading job is completed earlier when using CARD

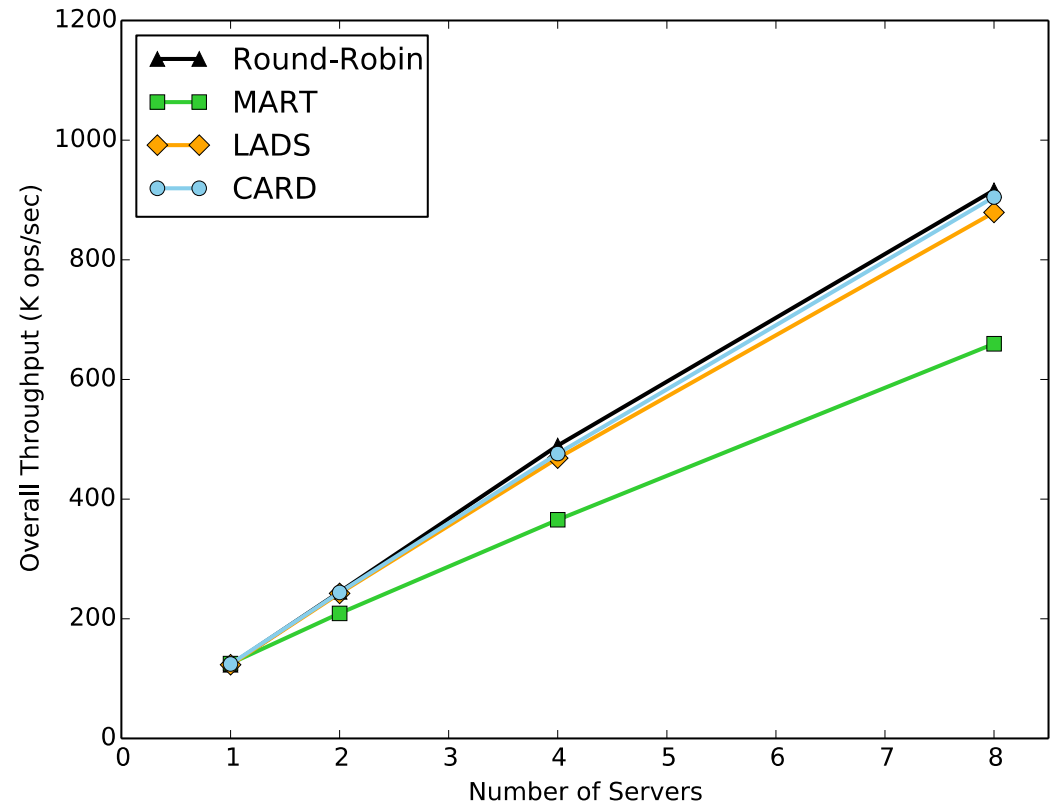# Overall arriving rates in the heterogeneous cluster



LADS



CARD

1) A basic threshold throttling strategy is not sufficient enough

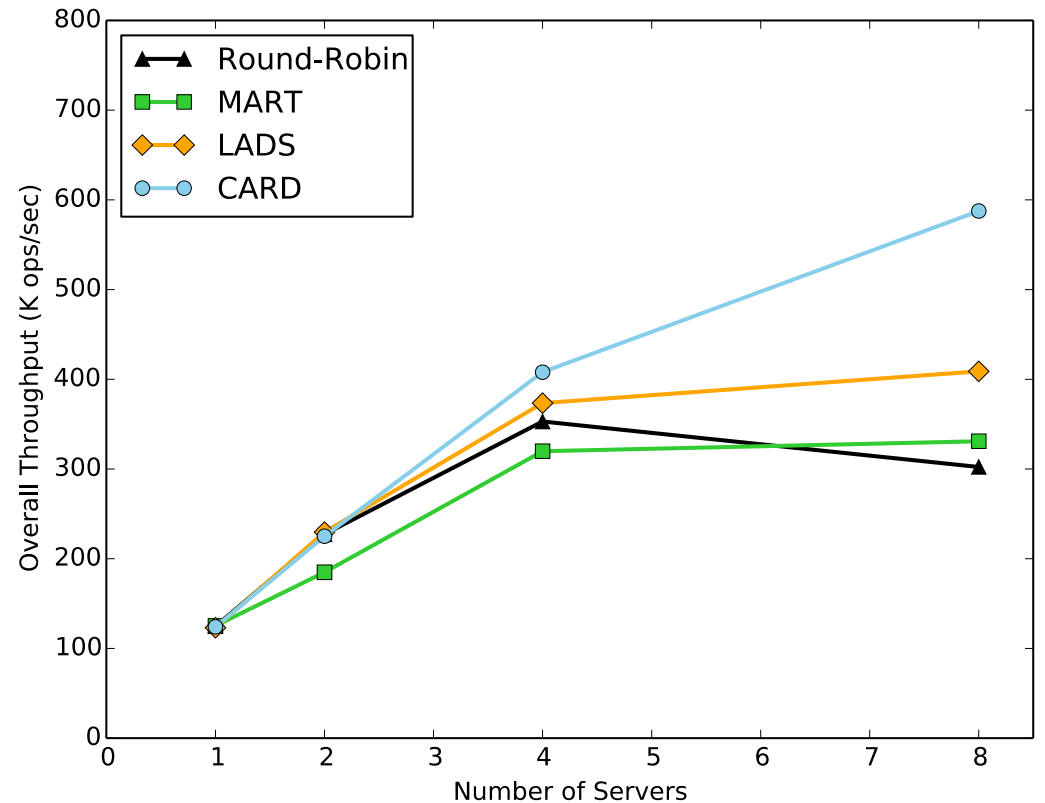2) Arriving rates are stabilized around servers' capacity when using CARD

# Overall throughput in the homogeneous cluster

- Heuristic selection is a bad choice under heavy workloads

- In ideal homogenous environments, Round-Robin and CARD achieve great scalability

# Overall throughput in the heterogeneous cluster

- Round-Robin is ineligible when facing heterogenous setups

- CARD outperforms other strategies and achieves excellent scalability



24

# Summary: CARD

- Adaptive client-side throttling method: easy and efficient

- Redirect requests from the overloaded server to the underloaded server adaptively under heavy workloads

- Degrade into pure Round-Robin when facing light-weight workloads

- Boosts throughput significantly over competing strategies in heterogeneous environments