

# Fast and Accurate Load Balancing for Geo-Distributed Storage Systems

Kirill L. Bogdanov<sup>1</sup>

Waleed Reda<sup>1,2</sup>

Gerald Q. Maguire Jr.<sup>1</sup>

Dejan Kostic<sup>1</sup>

Marco Canini<sup>3</sup>

<sup>1</sup>KTH Royal Institute of Technology

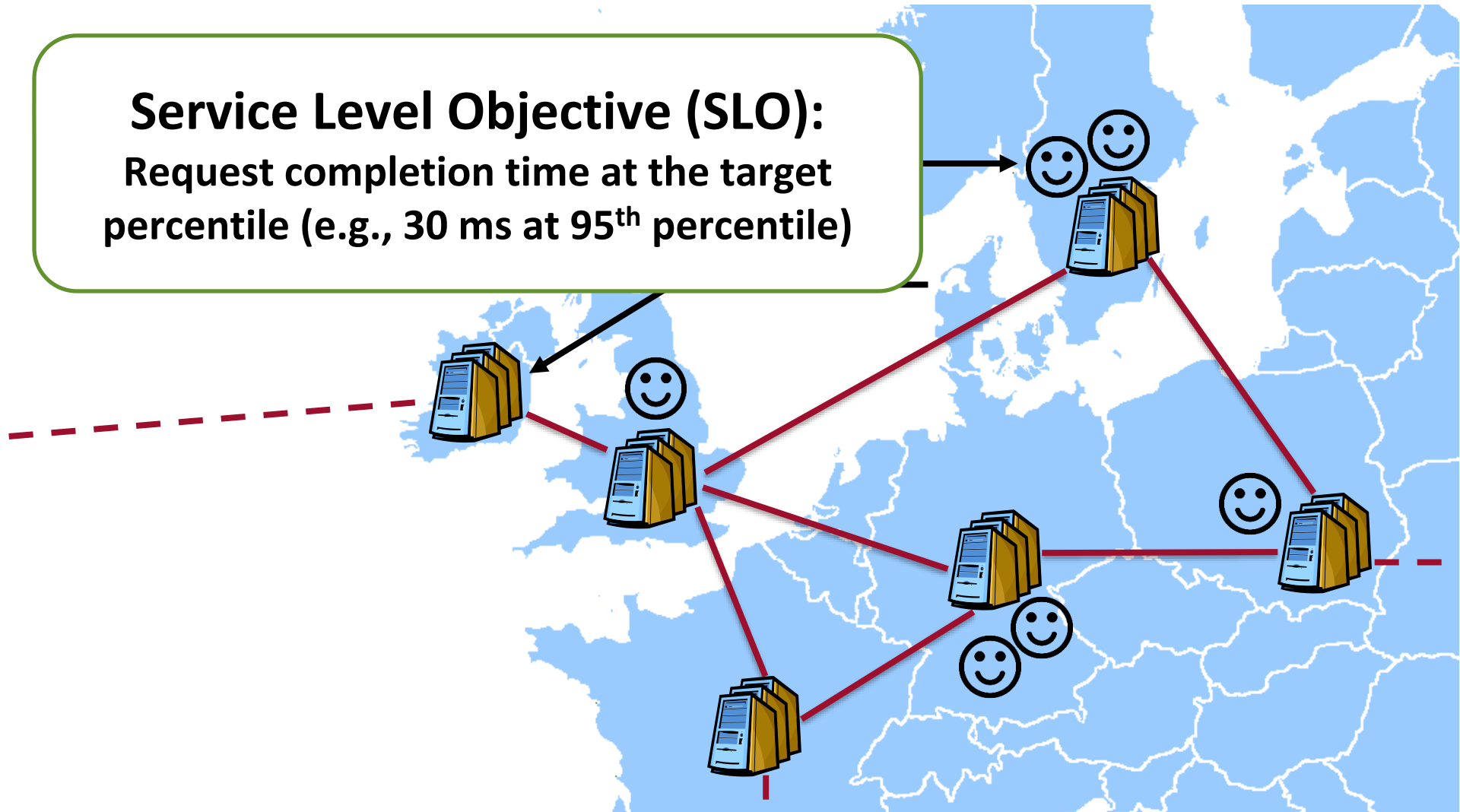
<sup>2</sup>Université Catholique de Louvain

<sup>3</sup>KAUST



# Geo-Distributed Services

**Service Level Objective (SLO):**  
Request completion time at the target  
percentile (e.g., 30 ms at 95<sup>th</sup> percentile)



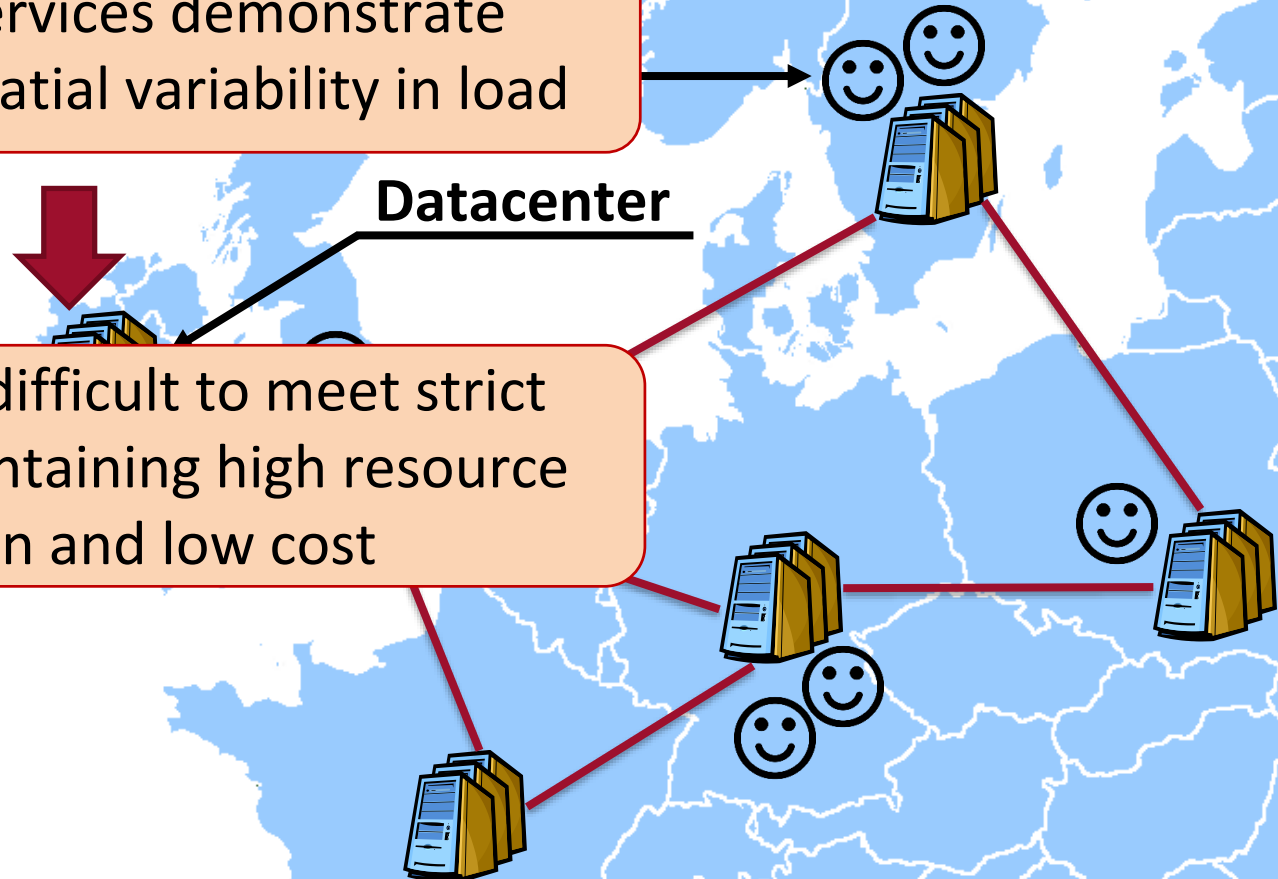
# Geo-Distributed Services

Web-based services demonstrate temporal and spatial variability in load

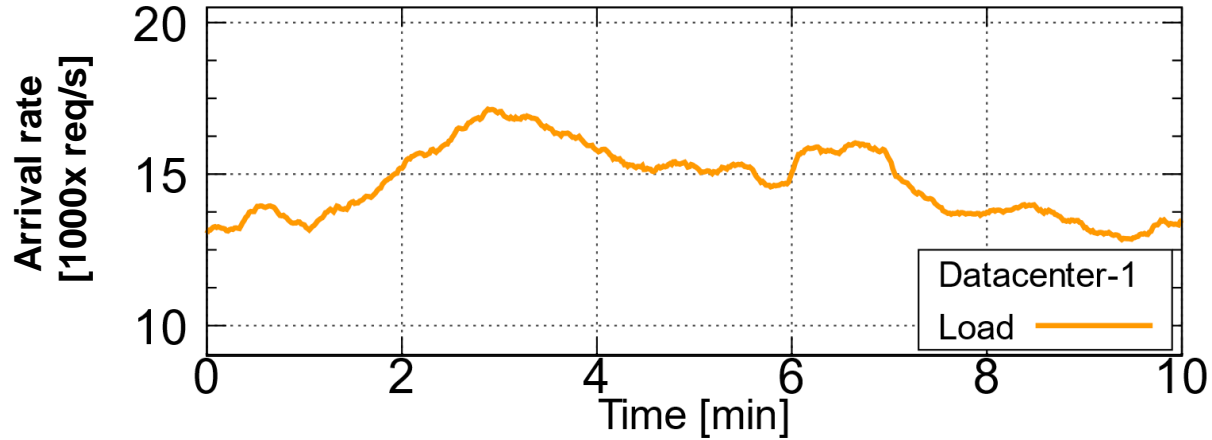


Datacenter

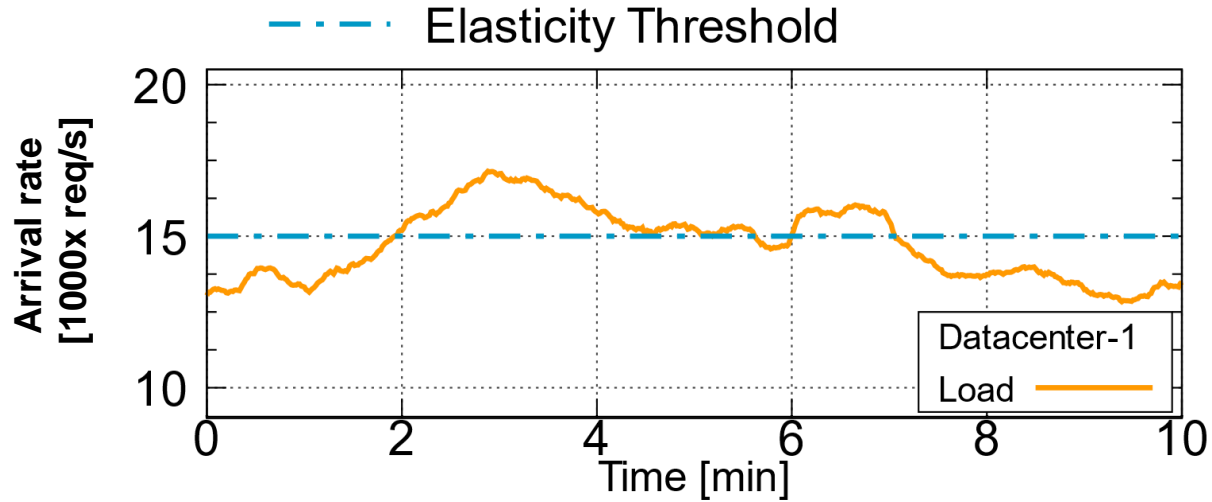
**Problem:** it is difficult to meet strict SLOs, while maintaining high resource utilization and low cost



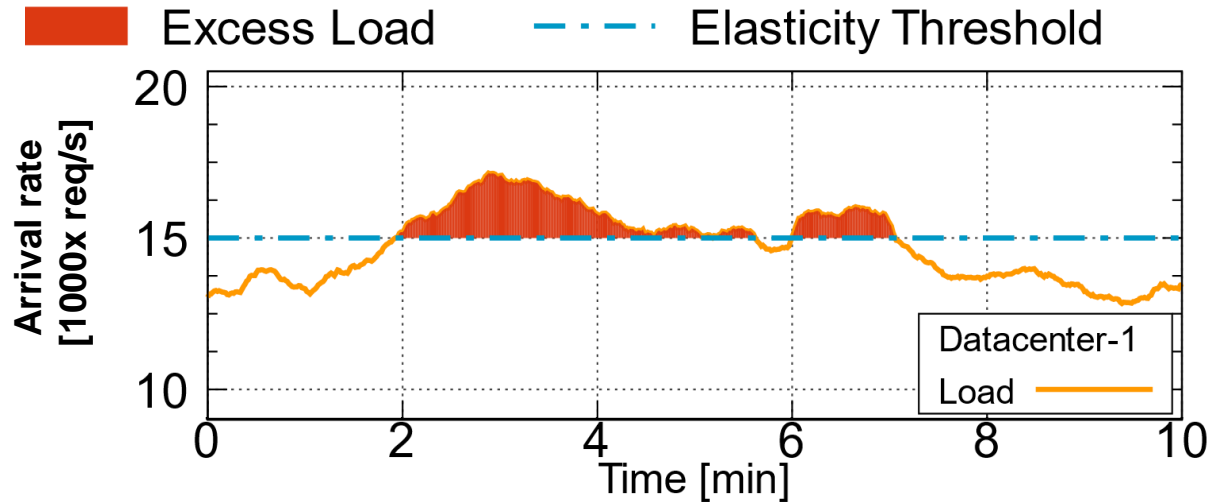
# Approach 1 - Datacenter Elasticity



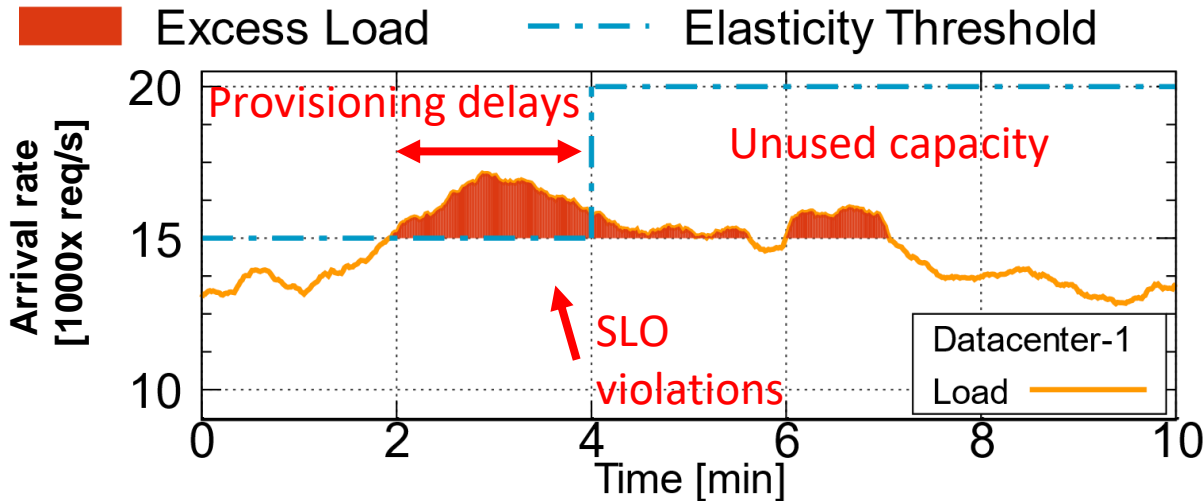
# Approach 1 - Datacenter Elasticity



# Approach 1 - Datacenter Elasticity



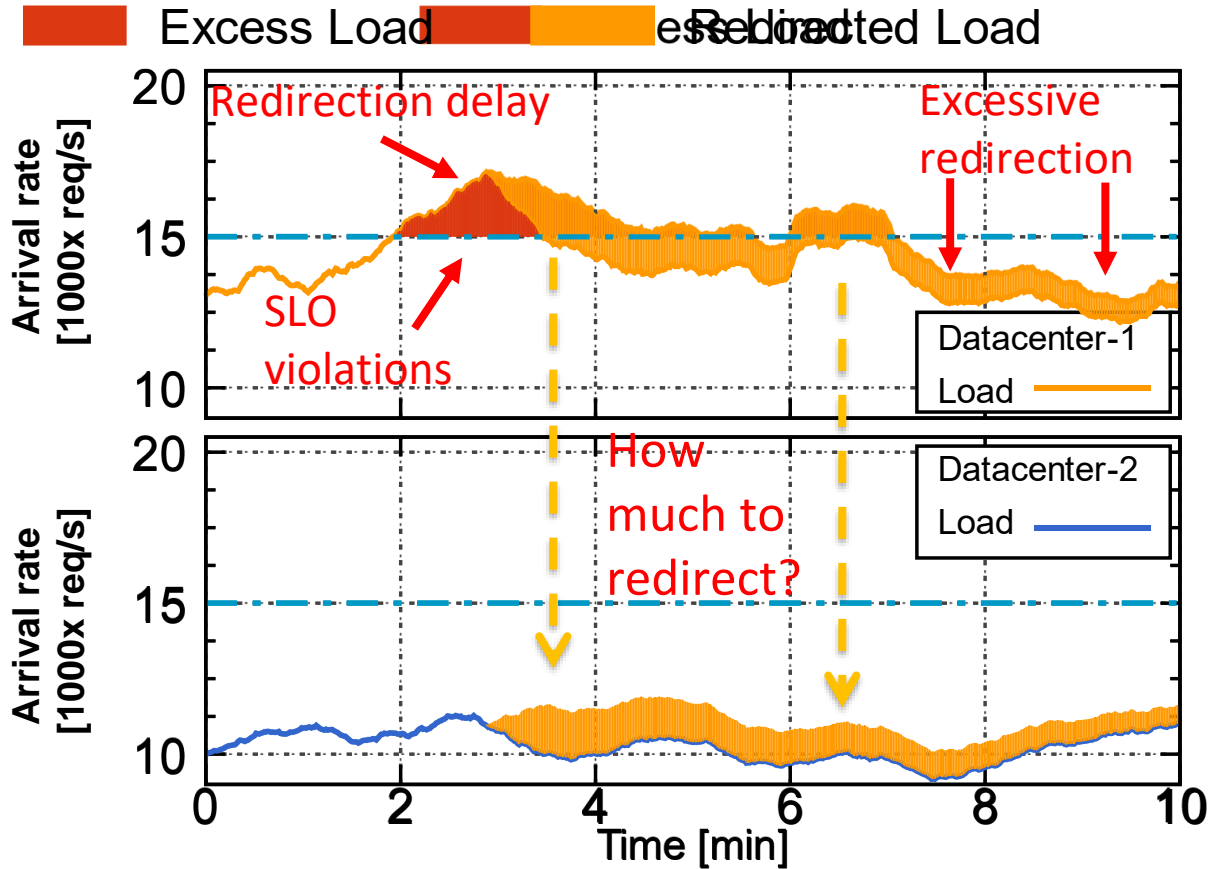
# Approach 1 - Datacenter Elasticity



Lead to overprovisioning!

- ✘ Provisioning delay (minutes) due to time needed to spawn and warm up a VM
- ✘ Hard to predict workload far into the future
- ✘ Load spikes can be short lived

# Approach 2 - Geo-Distributed Load Balancing



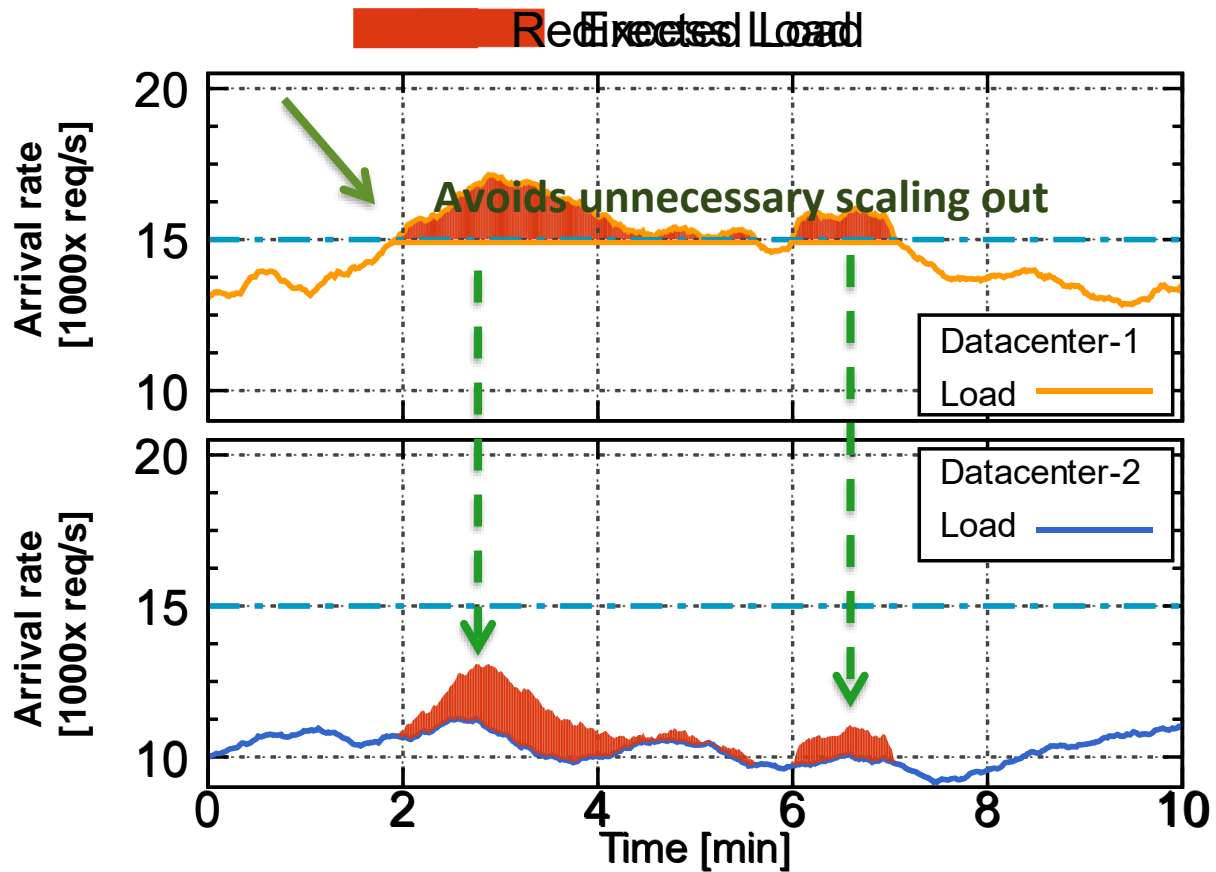
**Redirection delays**

**Inaccurate response time estimation**

**Excessive or insufficient redirection**



# Our Approach: Kurma



**Reacts to changes  
in load within  
seconds**

**Accurately  
estimates remote  
rate of SLO  
violations**

**Tames SLO  
violation at the  
target level**

# Request Completion Time

Datacenter Ireland

Server 1



Datacenter Frankfurt

Server 2



Wide Area  
Network

**Base Propagation:** Stable component associated with packet propagation along a network path

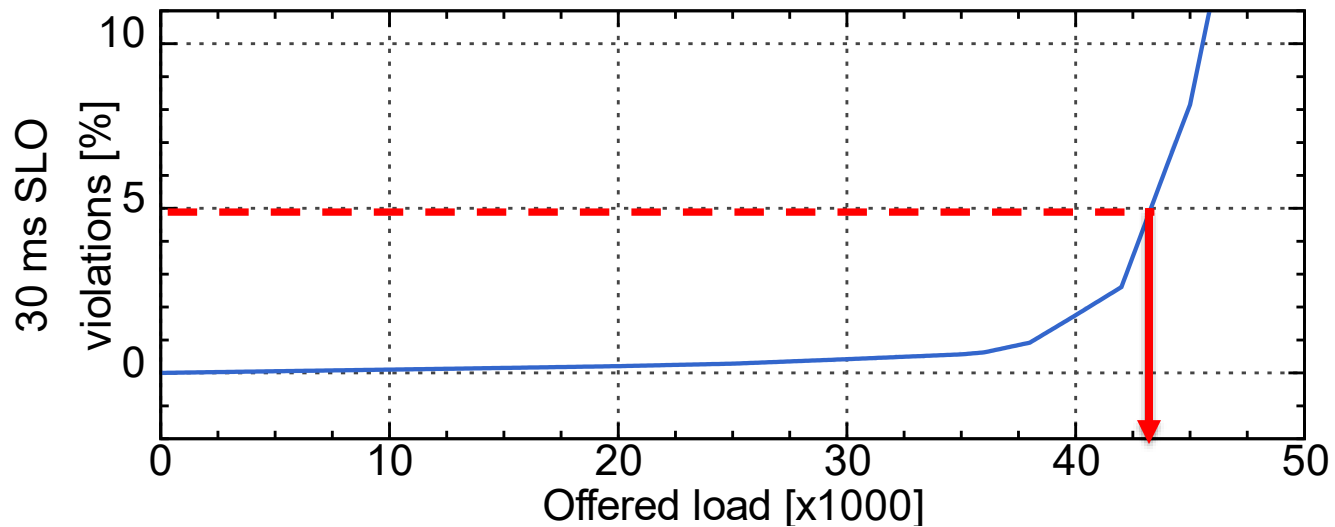
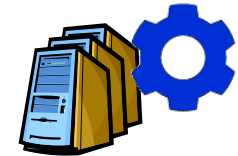
**Delay Variance:** Variable component associated with competing traffic and queuing

**Service Time:** Variable component associated with load on the server

**Kurma solves global optimization model while considering:  
Base Propagation + Delay Variance + Service Time  
at *all* datacenters**

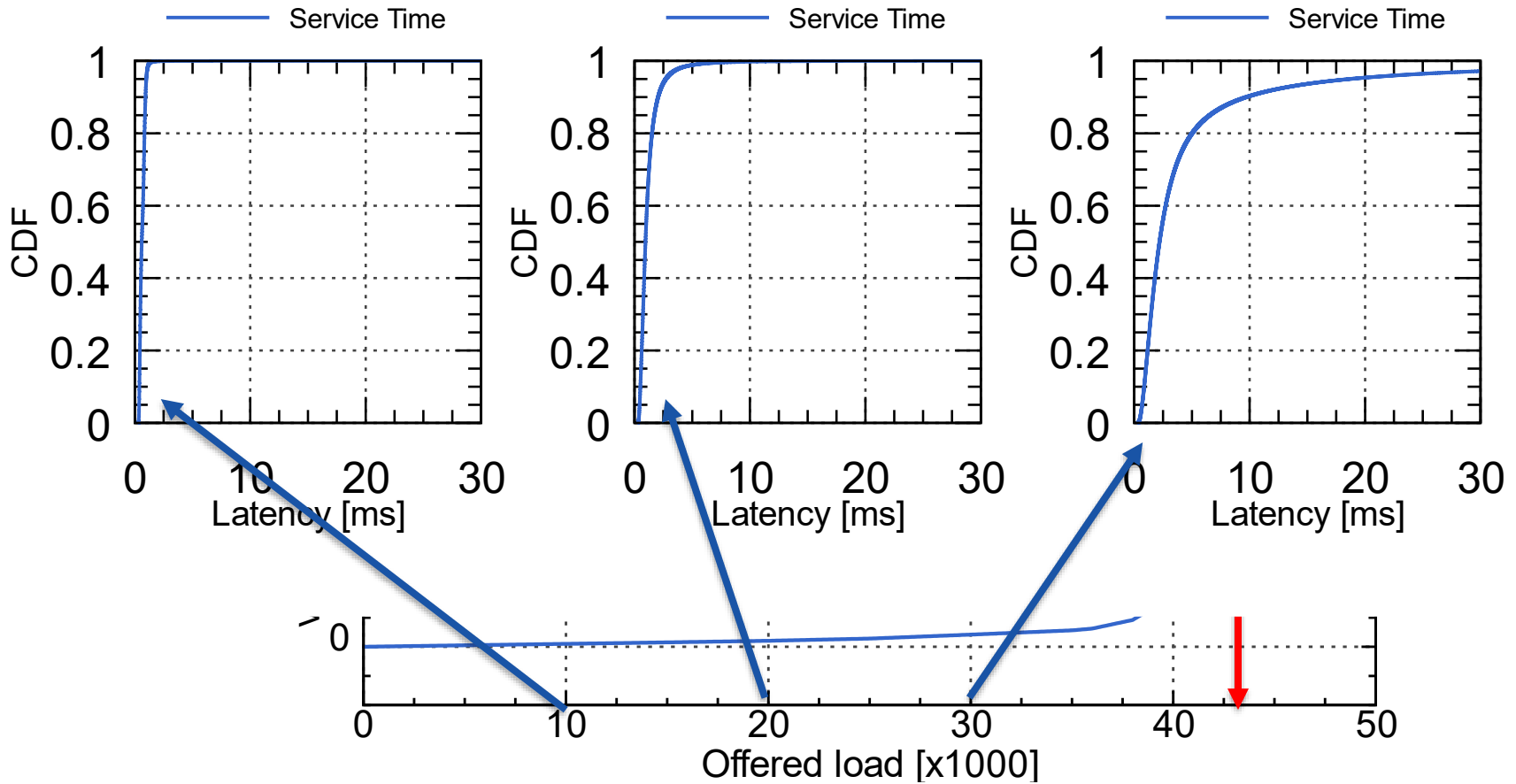
# Understanding Service Time

**Datacenter Frankfurt**  
**5 Server Cassandra cluster**





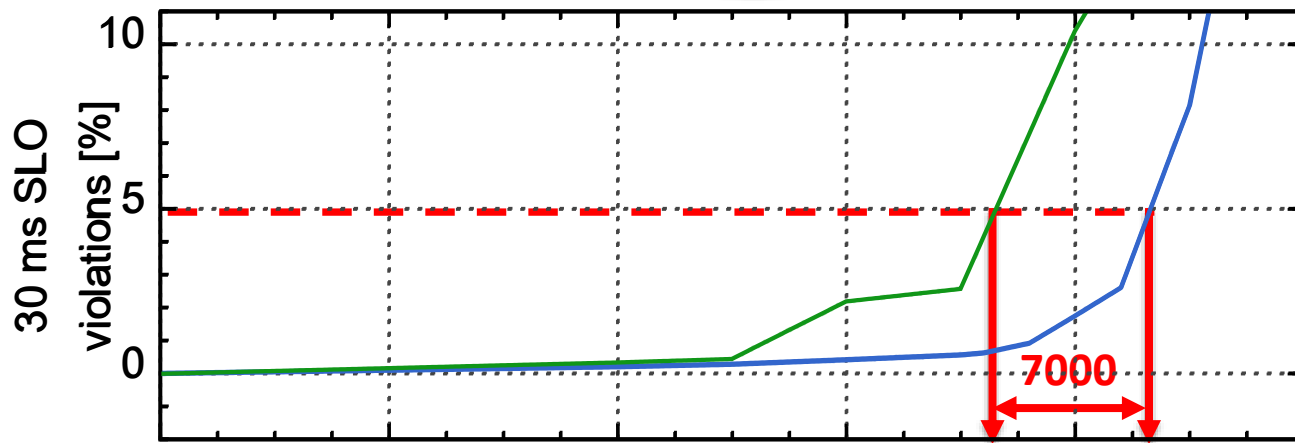
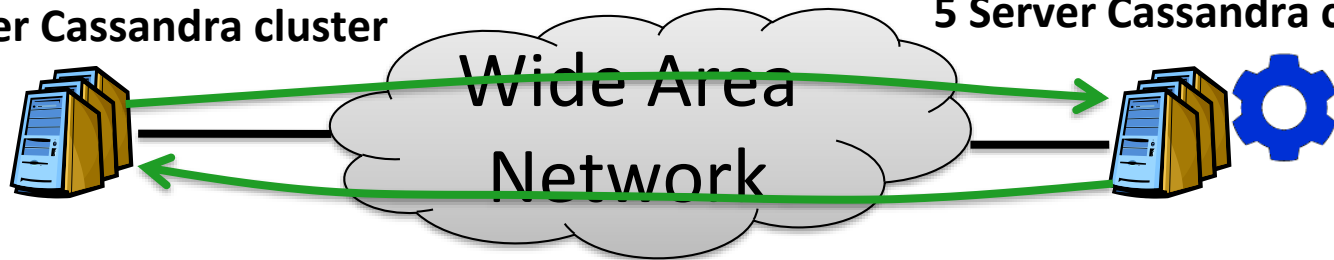
# Understanding Service Time



# Understanding Service Time

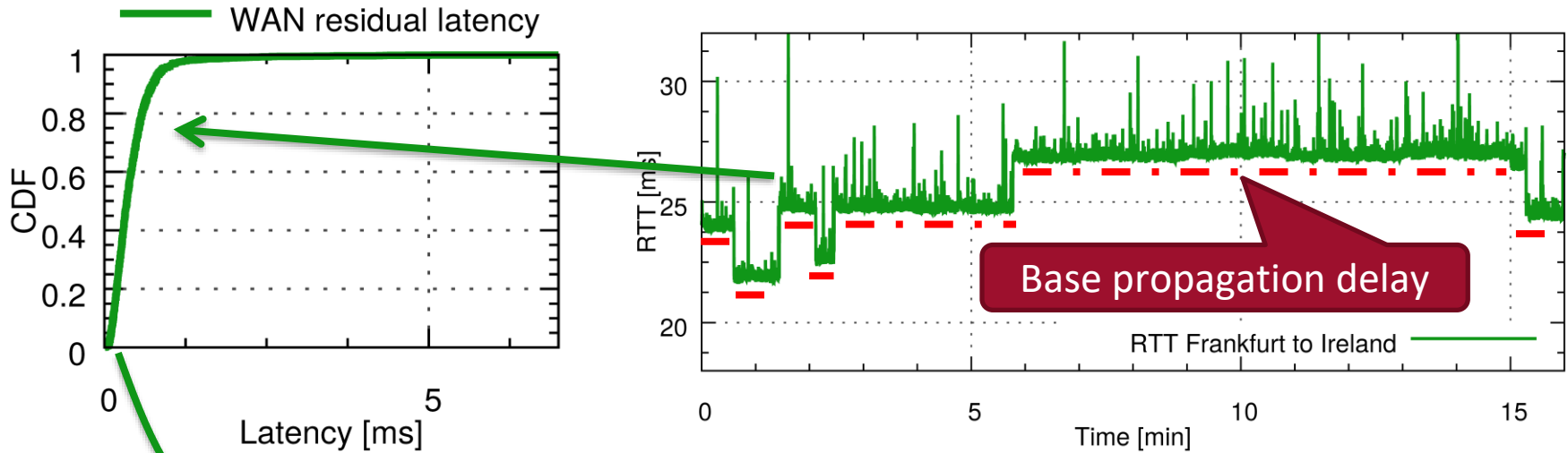
Datacenter Ireland  
5 Server Cassandra cluster

Datacenter Frankfurt  
5 Server Cassandra cluster



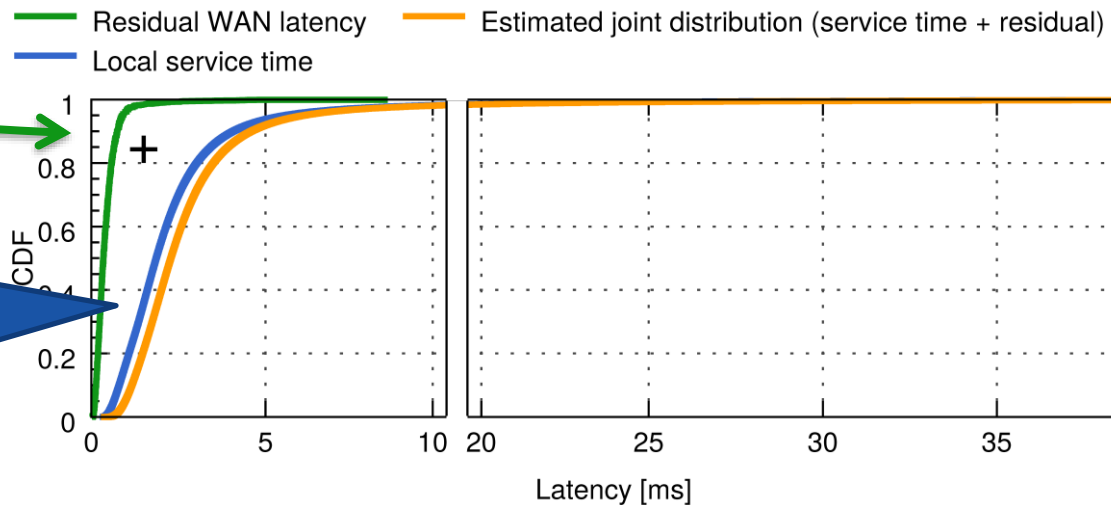
**Insight:** the farther away a remote datacenter is, the less loaded it should be to serve remote requests within a given SLO target

# Understanding WAN Latency

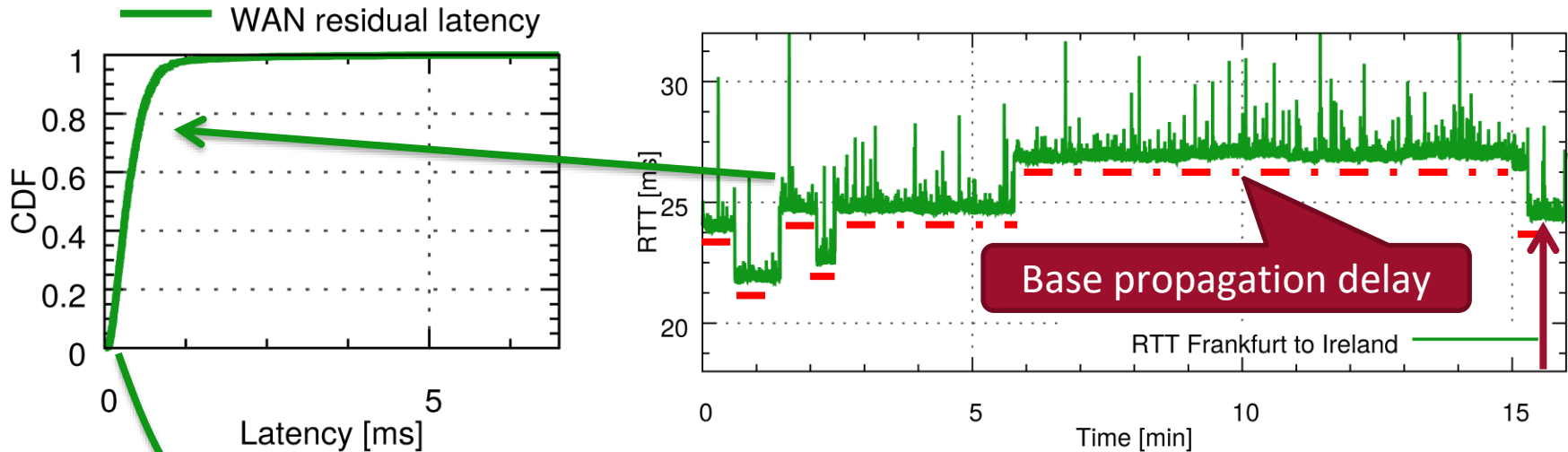


Monte Carlo Simulations

Service time distribution recorded locally at a specific load

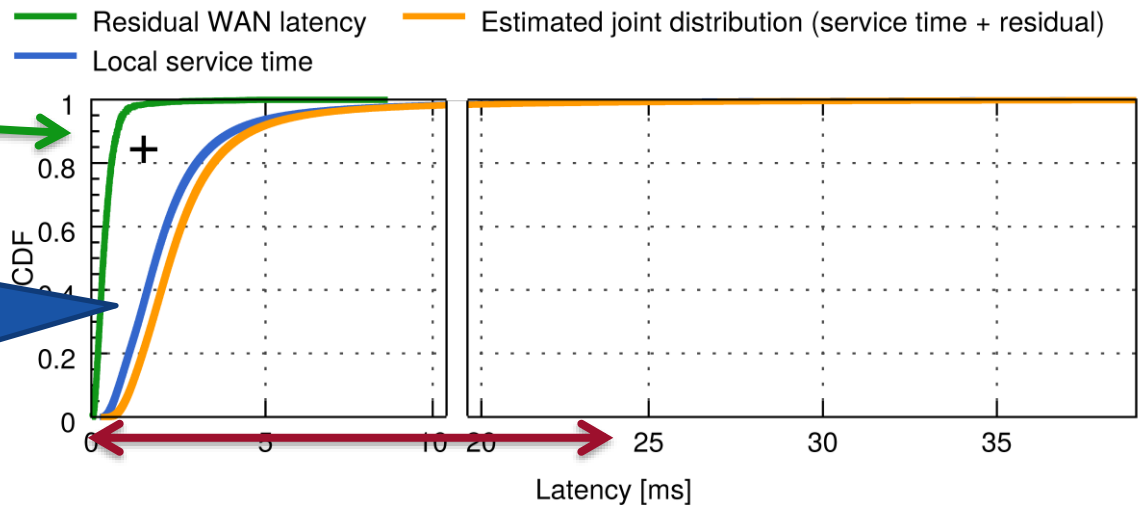


# Understanding WAN Latency



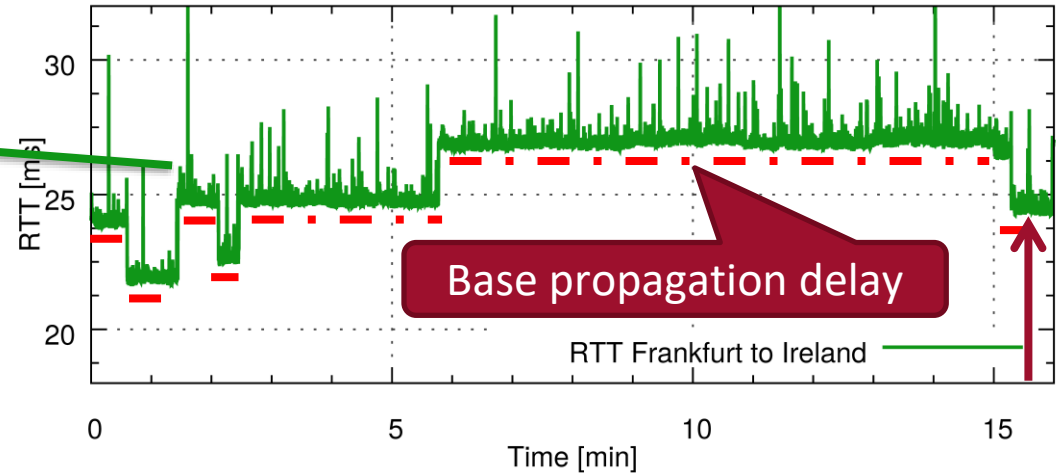
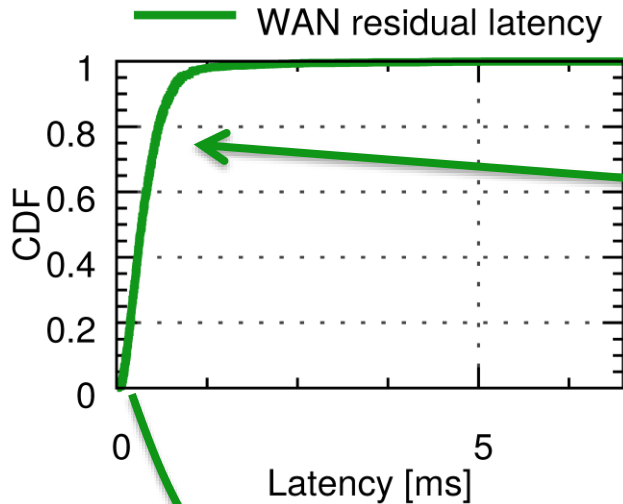
Monte Carlo Simulations

Service time distribution recorded locally at a specific load



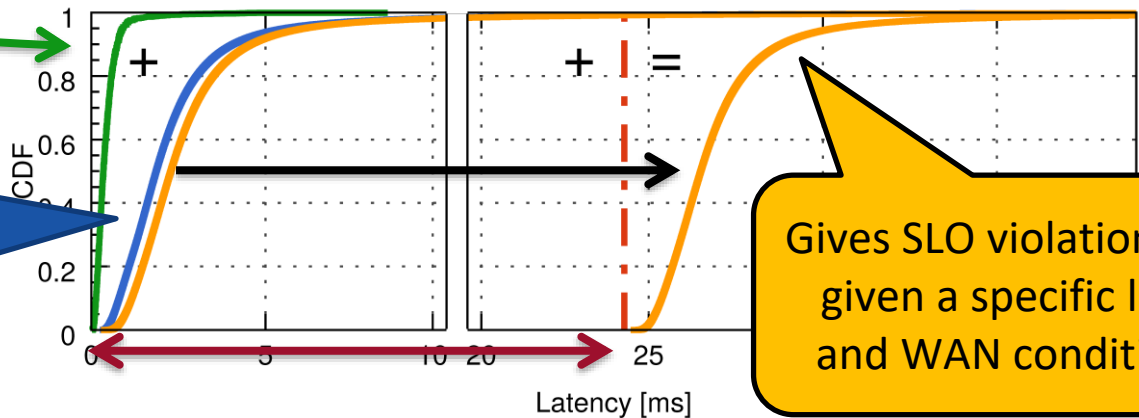


# Understanding WAN Latency



Monte Carlo Simulations

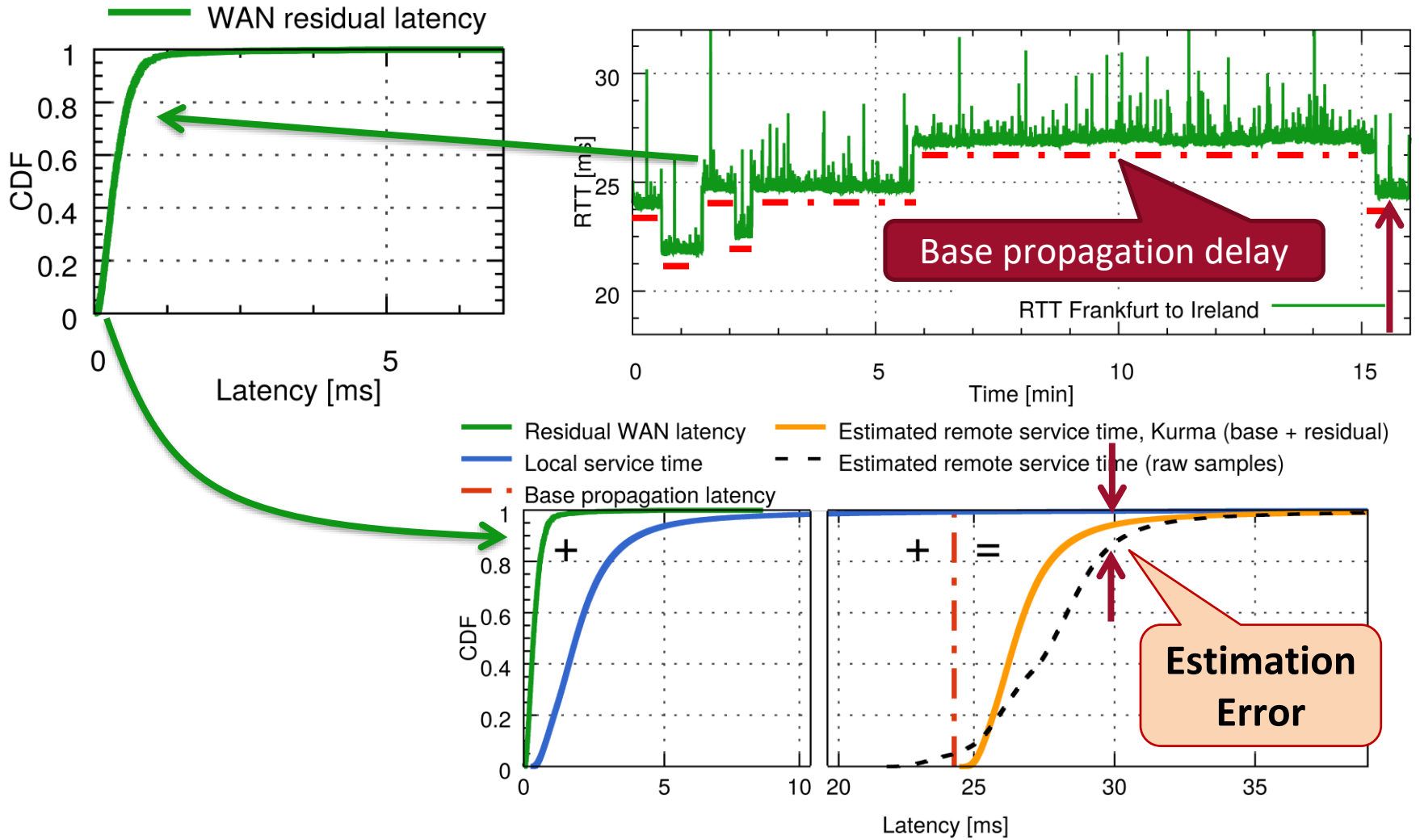
- Residual WAN latency
- Local service time
- Base propagation latency
- Estimated remote service time, Kurma (base + residual)



Service time distribution recorded locally at a specific load

Gives SLO violation rate given a specific load and WAN conditions

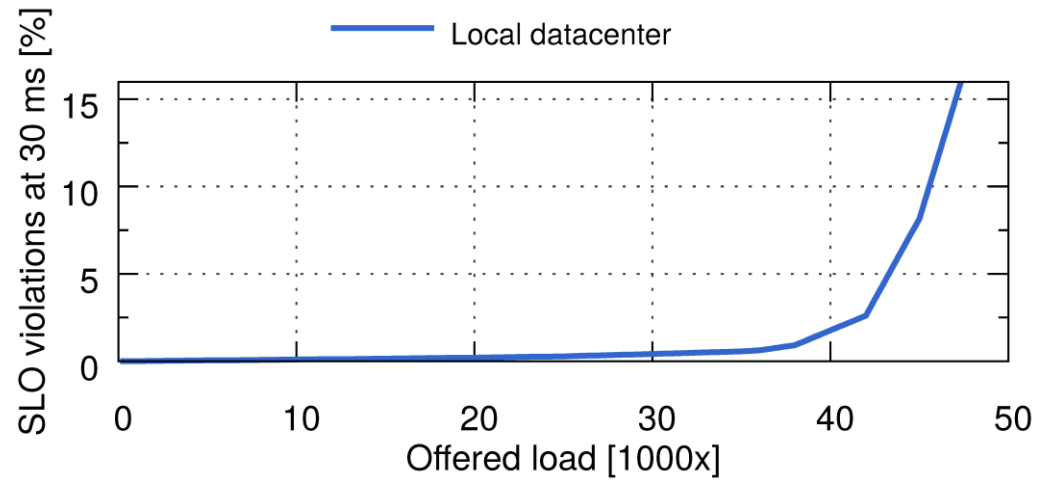
# Understanding WAN Latency





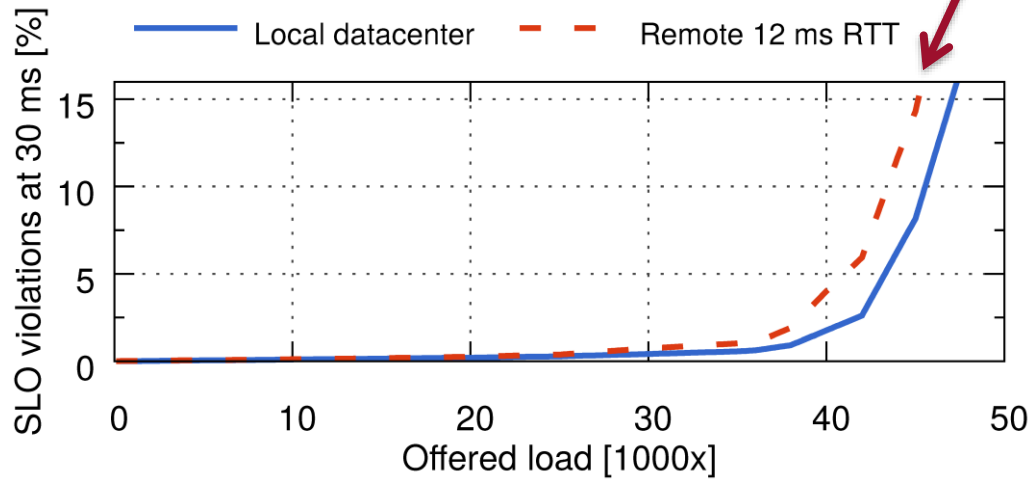
# Incorporating WAN and Load

5 VM Cassandra cluster



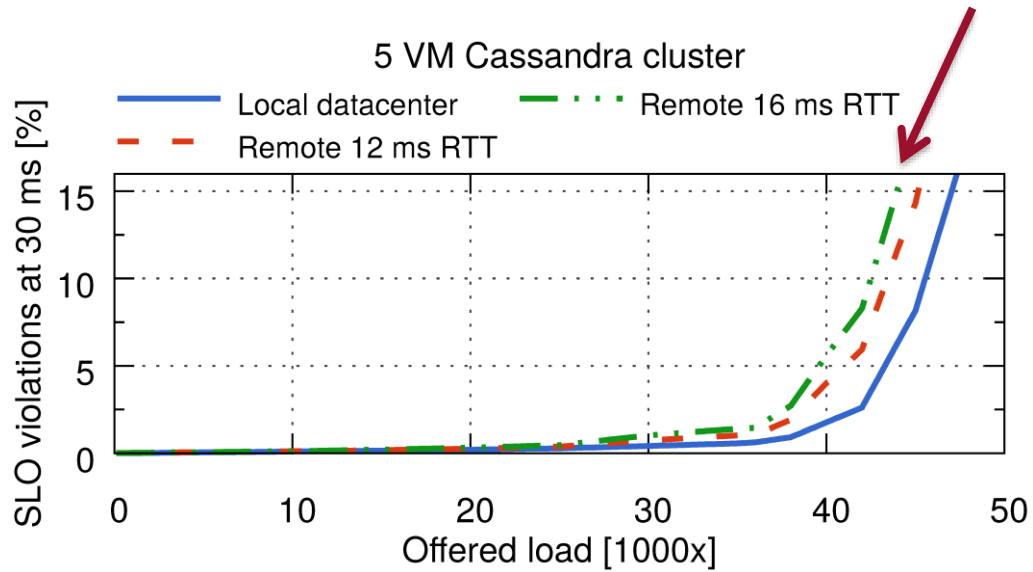
# Incorporating WAN and Load

5 VM Cassandra cluster



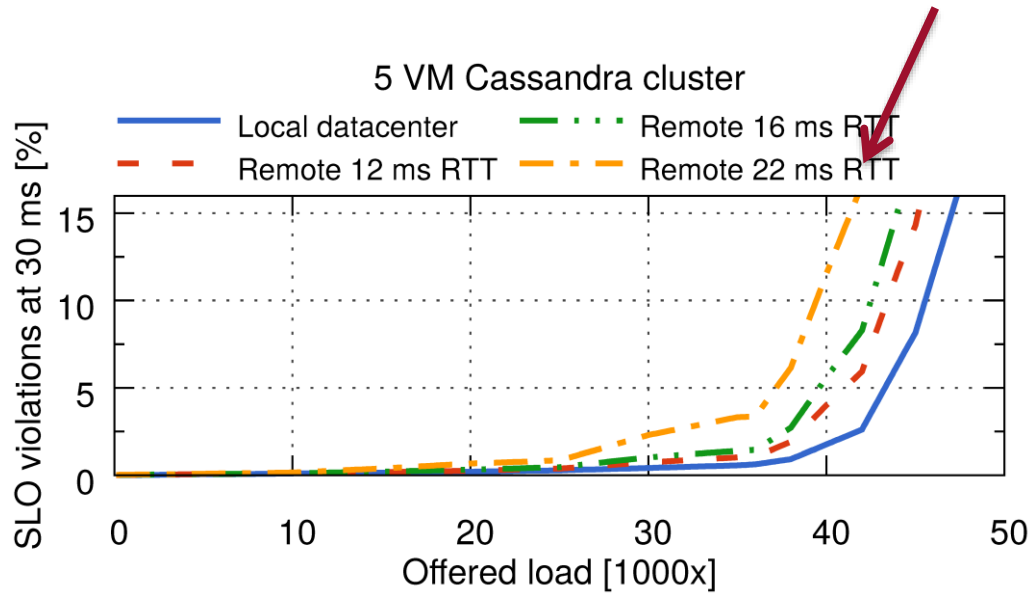


# Incorporating WAN and Load



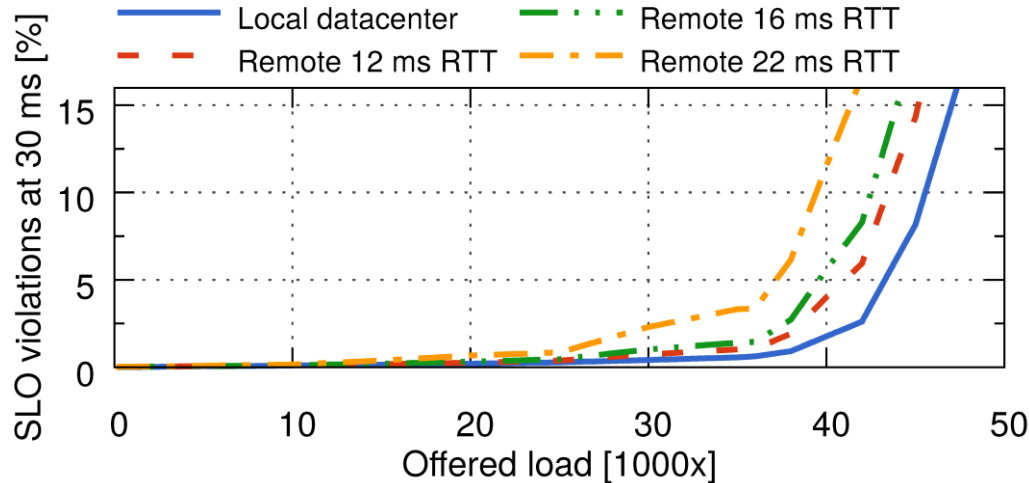


# Incorporating WAN and Load



# Optimisation Model

5 VM Cassandra cluster



+

Runtime load in each datacenter  
 $\{\lambda_1, \lambda_2, \lambda_3\}$

- Optimisation Problem**
- ✓ Minimize global SLO violations (KurmaPerf)
  - ✓ Minimize the cost of running a service (KurmaCost)

# Implementation

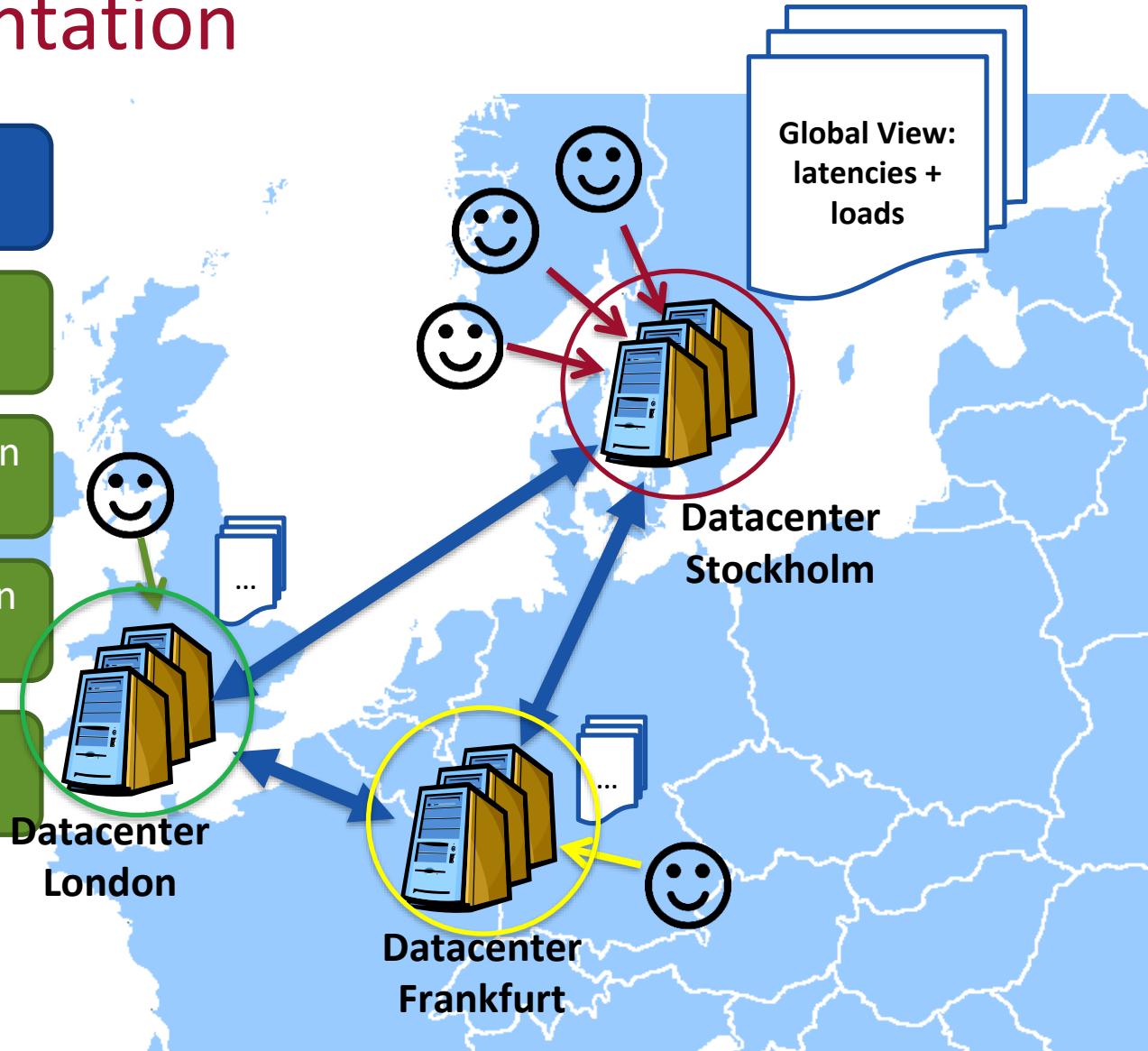
Each Epoch  
2.5 sec  $\rightarrow$  0.4Hz

Perform run-time WAN  
latency measurements

Aggregate load information  
(rates of requests)

Exchange metrics to obtain  
global view

Solve decentralized  
performance model









# Evaluation Setup

## Geo-distributed Cassandra cluster

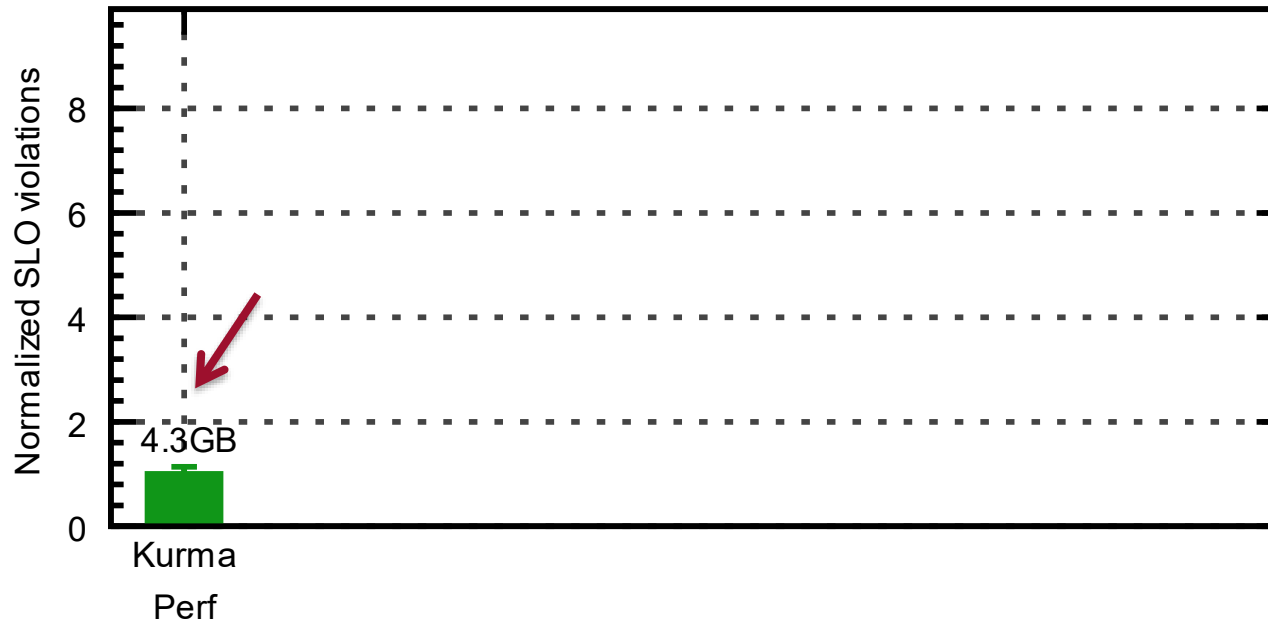
- 3 Amazon EC2 datacenter (Ireland, Frankfurt, London)
- 5 x r5.large VMs per datacenter
- SLO: 30 ms at the 95<sup>th</sup> percentile
- Modified YCSB to replay workload traces  
(World Cup <http://ita.ee.lbl.gov/html/contrib/WorldCup.html>)

## Experiments:

- **Minimizing SLO violations for reads**
- **Maintaining Target SLO (accuracy)**
- **Cost Savings for 1 min billing intervals (simulations)**
- Reads and writes, scalability, etc. [link here](#).



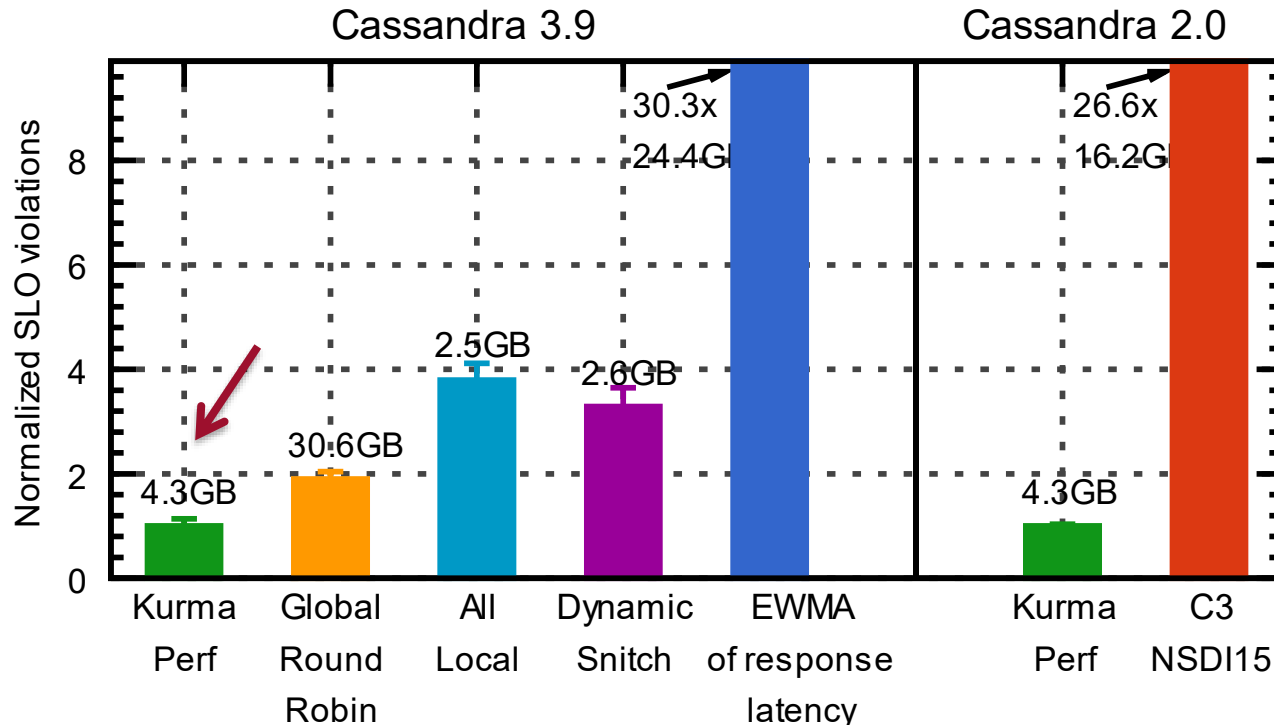
# Cumulative Normalized SLO Violations



Kurma's SLO violations are at 2.4%

The numbers shown above the bars indicate the amount of inter-datacentre traffic transferred, whiskers → 75<sup>th</sup> percentile

# Cumulative Normalized SLO Violations



Kurma's SLO violations are at 2.4%

The numbers shown above the bars indicate the amount of inter-datacentre traffic transferred, whiskers → 75<sup>th</sup> percentile











# Conclusion

# Q&A

Kurma – **fast and accurate** load balancer for geo-distributed systems that takes advantage of spatial variability in load

**Decouples** end-to-end response time into components of base propagation latency, network congestion, and service time distribution

By operating at the granularity of a few seconds, Kurma **reduces SLO violations** or **lowers the costs** of running services by avoiding excessive global service overprovisioning