

Proving Network Latency Guarantees in Data Centers

Diana Andreea Popescu
PhD student, Computer Laboratory
Supervisor: Andrew Moore



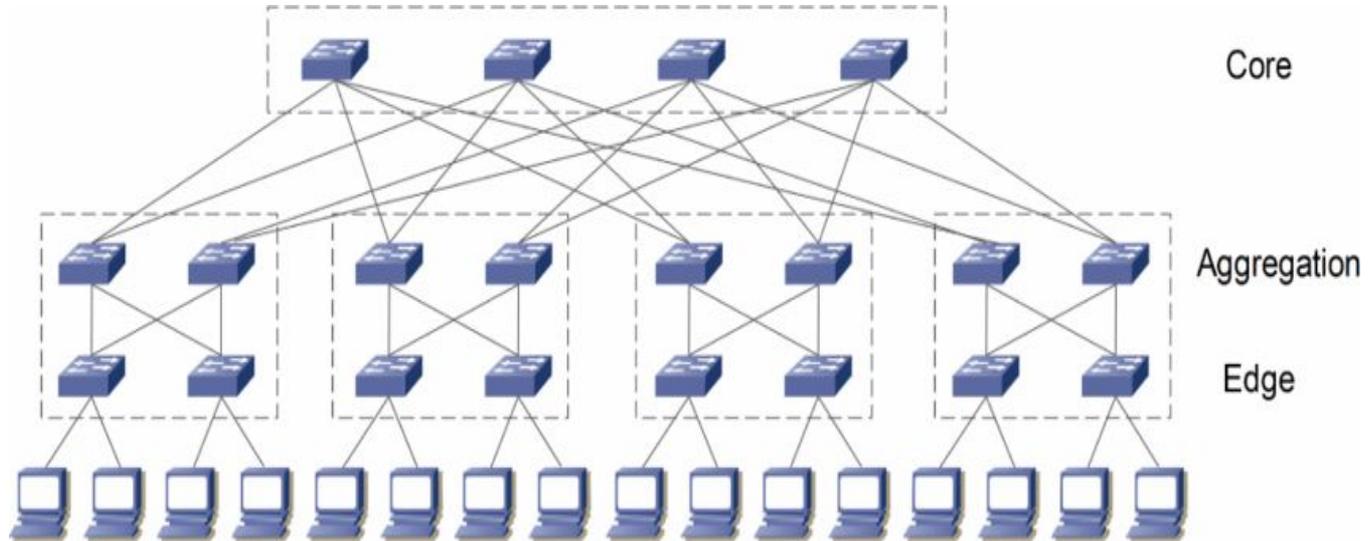
What is a data center?



Who's using them?

- Google, Facebook, Microsoft, Amazon...
- Tenants
 - Rent computing resources
 - Run their applications

How a data centre looks like



Motivation

- Data centers are **really** expensive
- Better data center use through (network) measurements
- SLA (service-level agreements) respected

Fair Application Competition *please....*



Fair Application Competition *please....*



Fair Application Competition *please*....



Fair Application Competition *please....*



**We need intelligent decisions
to solve this**

What do we usually measure?

- CPU usage
- Memory usage
- Network bandwidth
- **Network latency**

What do we usually measure?

- CPU usage
- Memory usage
- Network bandwidth
- **Network latency**

**Use this information
to make intelligent decisions**

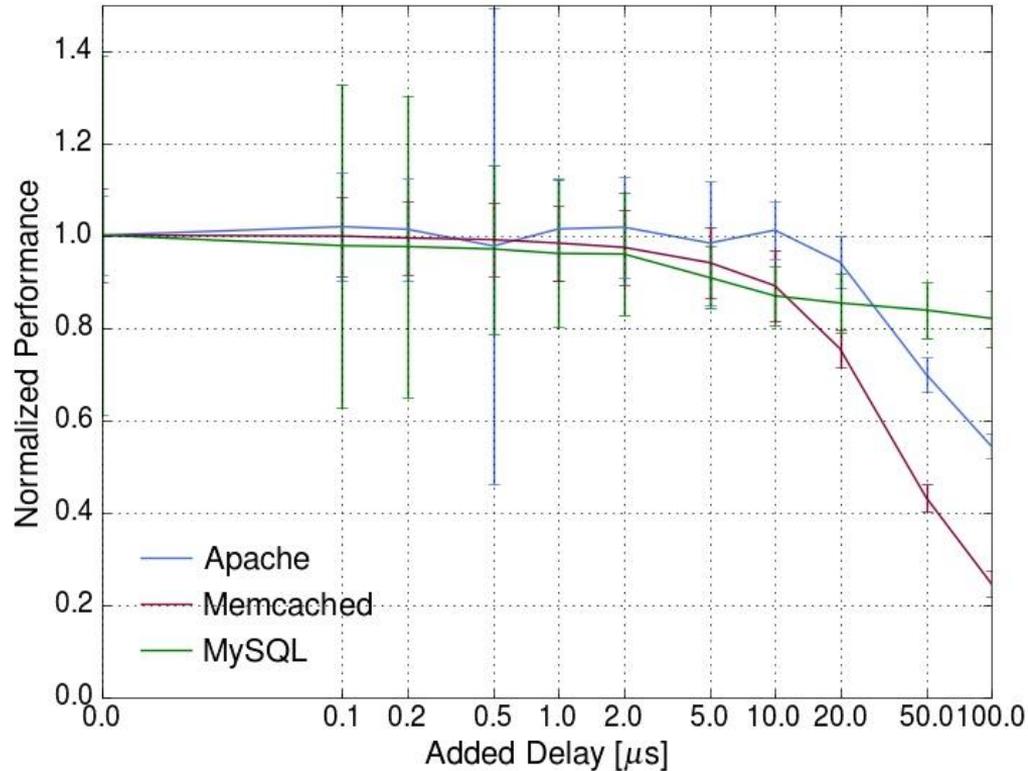
How can we measure network latency in a data centre?

- Each host has a GPS receiver
 - Clocks are synchronised, can run one-way delay measurements
 - Not scalable
- Ping
 - Measure RTT (round-trip time)
 - Overhead is high, accuracy is not good
- DTP (Datacentre Time Protocol)
 - Hardware-based solution, not immediately deployable
- PTP (Precision Time Protocol)
 - Moderate overhead, sub-us precision
- Custom solution: be able to set probe rate and not incur high overhead

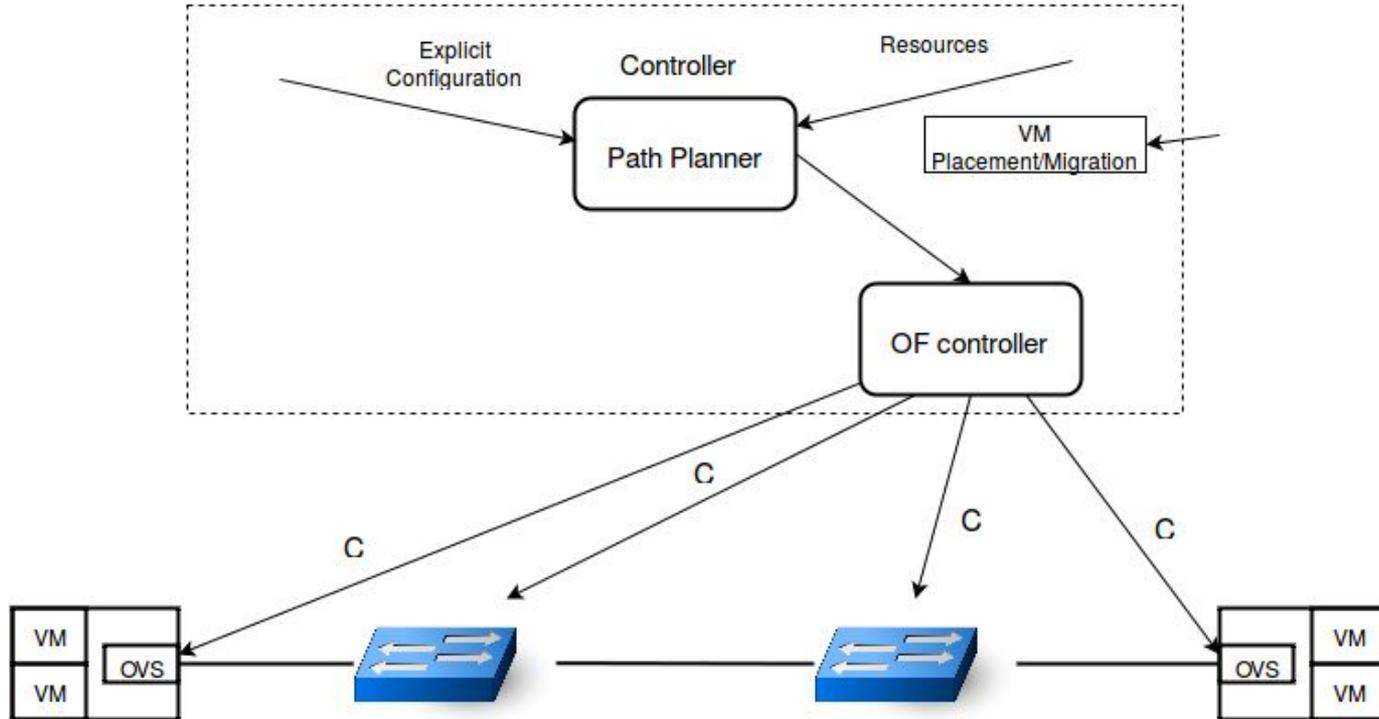
Current Work

- Representative communication patterns
- Model the relationship between application performance and network latency
- Applications:
 - Key-value Stores (Memcached)
 - Relational Databases (MySQL)
 - Web Servers (Apache)
 - Distributed Machine Learning
 - Graph processing

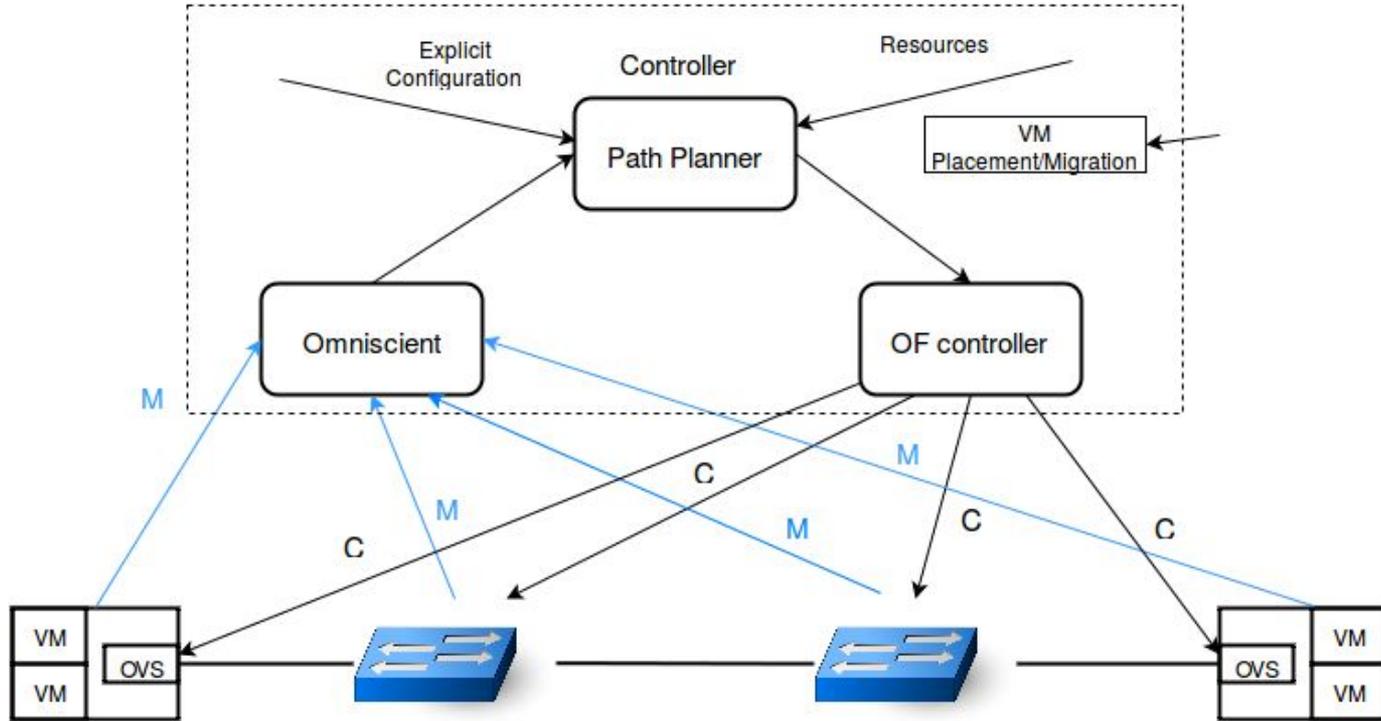
Delay Effect on Application Performance



Omniscient



Omniscient



Conclusion

Our objective: a higher degree of automation



diana.popescu@cl.cam.ac.uk

www.cl.cam.ac.uk/~dap53