

Data Center Storage and Networking

Hakim Weatherspoon

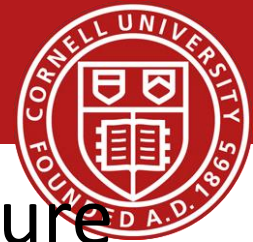
Assistant Professor, Dept of Computer Science

CS 5413: High Performance Systems and Networking

December 1, 2014

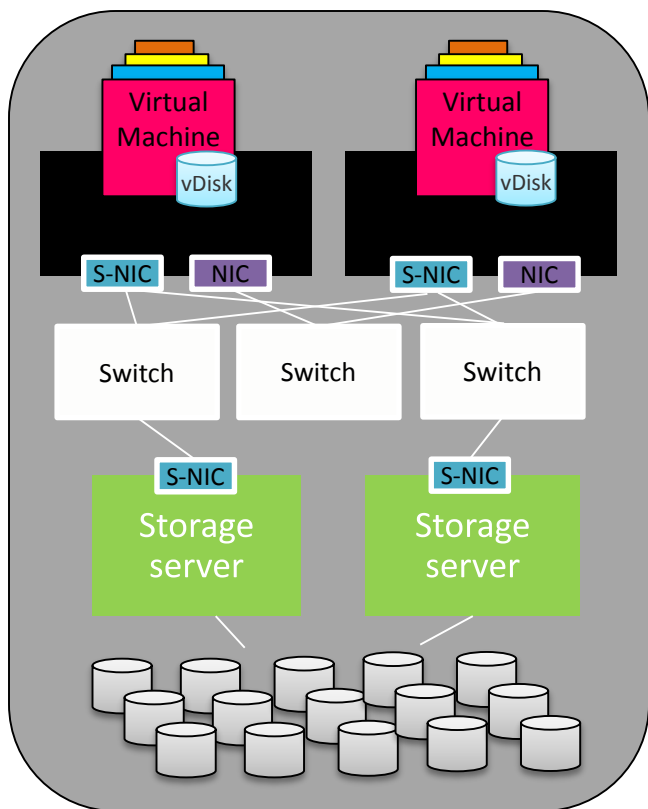
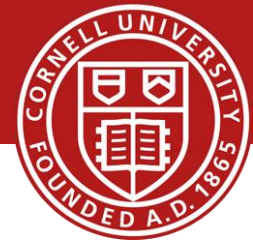
Slides from ACM SOSP 2013 presentation on “IOFlow: A Software-Defined Storage Architecture.” Eno Thereska, Hitesh Ballani, Greg O'Shea, Thomas Karagiannis, Antony Rowstron, Tom Talpey, and Timothy Zhu. In SOSP'13, Farmington, PA, USA. November 3-6, 2013. “

Goals for Today

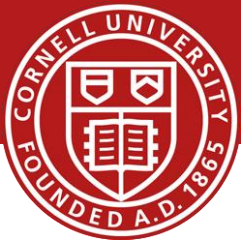


- IOFlow: a software-defined storage architecture
 - E. Thereska, H. Ballani, G. O'Shea, T. Karagiannis, A. Rowstron, T. Talpey, R. Black, T. Zhu. ACM Symposium on Operating Systems Principles (SOSP), October 2013, pages 182-196.

Background: Enterprise data centers



- General purpose applications
- Application runs on several VMs
- Separate network for VM-to-VM traffic and VM-to-Storage traffic
- Storage is virtualized
- Resources are shared

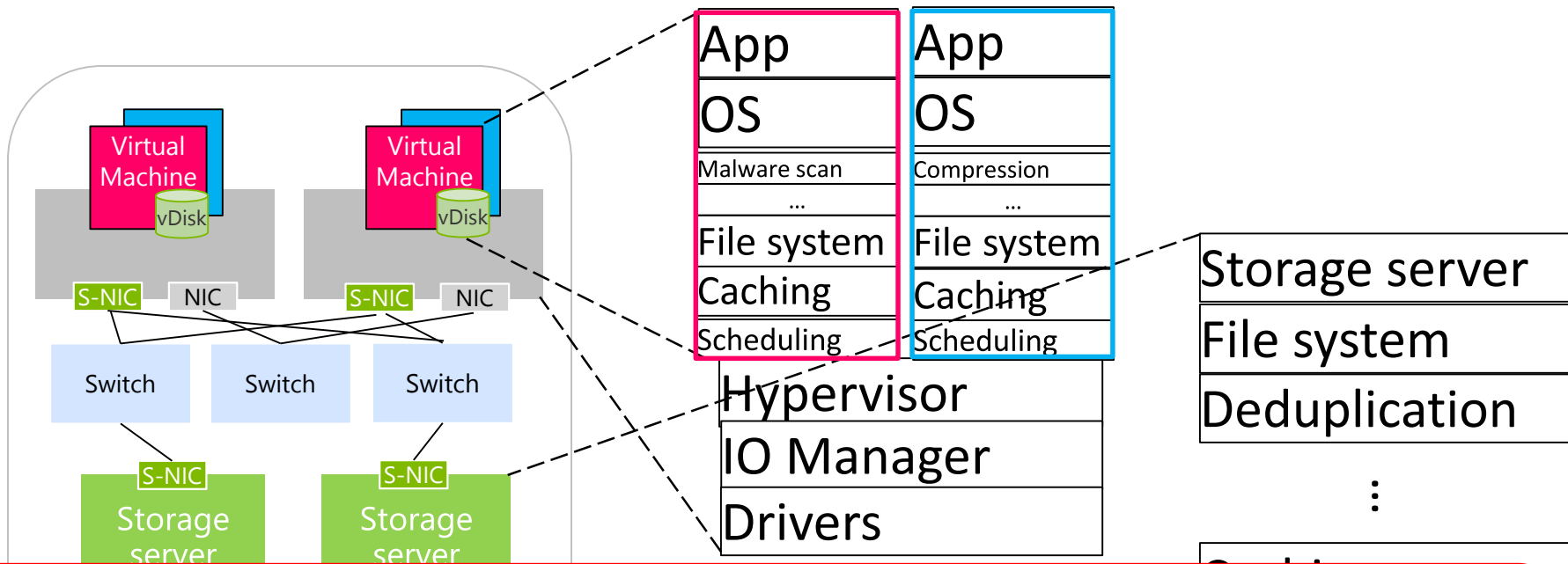
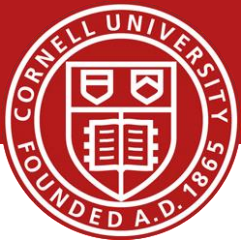


Want: predictable application behaviour and performance

Need system to provide end-to-end SLAs, e.g.,

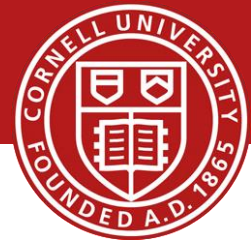
- Guaranteed storage bandwidth B
 - Guaranteed high IOPS and priority
 - Per-application control over decisions along IOs' path
-
- It is hard to provide such SLAs today

Example: guarantee aggregate bandwidth B for Red tenant



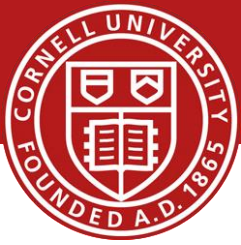
Deep IO path with 18+ different layers that are configured and operate independently and do not understand SLAs

Challenges in enforcing end-to-end SLAs

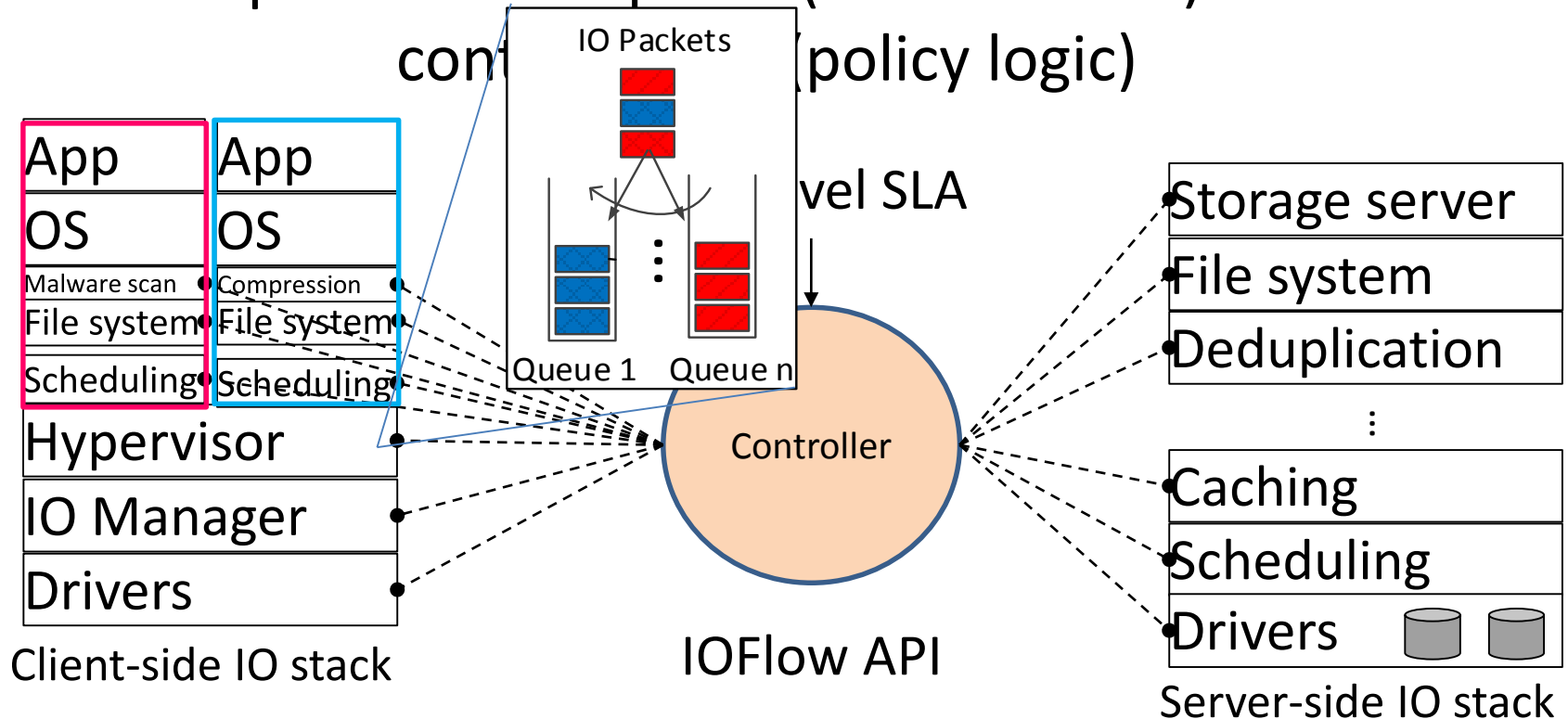


- No storage control plane
- No enforcing mechanism along storage data plane
- Aggregate performance SLAs
 - Across VMs, files and storage operations
- Want non-performance SLAs: control over IOs' path
- Want to support unmodified applications and VMs

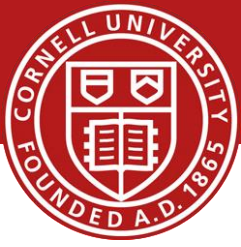
IOFlow architecture



Decouples the data plane (enforcement) from the control plane (policy logic)

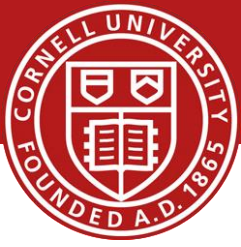


Contributions



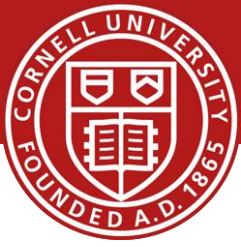
- Defined and built storage control plane
- Controllable queues in data plane
- Interface between control and data plane (IOFlow API)
- Built centralized control applications that demonstrate power of architecture

SDS: Storage-specific challenges



Low-level primitives	Old networks	SDN	Storage today	SDS
End-to-end identifier	✓ →	✓	✗	✓
Data plane queues	✓ →	✓	✗	✓
Control plane	✓ →	✓	✗	✓

Storage flows

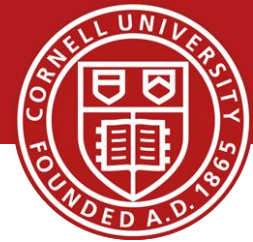


- Storage “Flow” refers to all IO requests to which an SLA applies

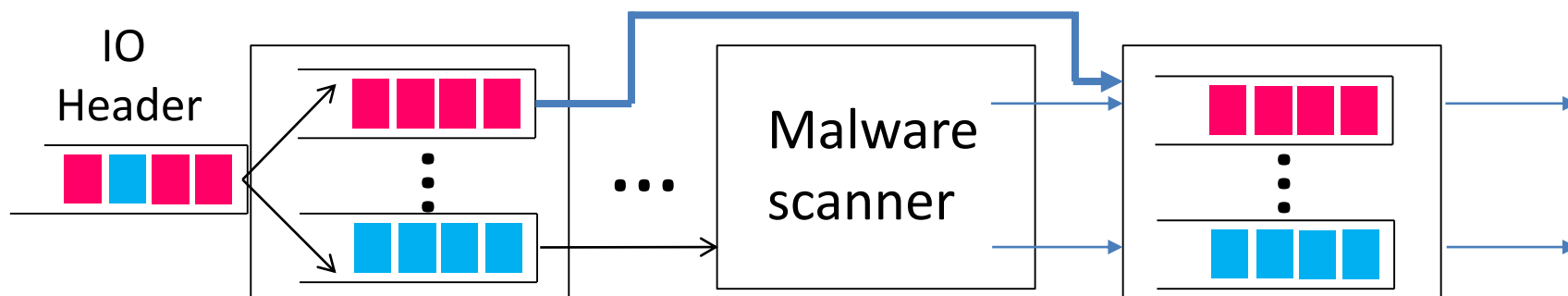
$\langle \underbrace{\{\text{VMs}\}}_{\text{source set}}, \{\text{File Operations}\}, \underbrace{\{\text{Files}\}, \{\text{Shares}\}}_{\text{destination sets}} \rangle \text{ --- } \text{SLA}$

- **Aggregate, per-operation and per-file SLAs, e.g.,**
 - $\langle \{VM\ 1-100\}, \text{write}, *, \\ \backslash \backslash \text{share} \backslash \text{db-log} \rangle \text{ --- } \text{high priority}$
 - $\langle \{VM\ 1-100\}, *, *, \\ \backslash \backslash \text{share} \backslash \text{db-data} \rangle \text{ --- } \text{min } 100,000 \text{ IOPS}$
- **Non-performance SLAs, e.g., path routing**
 - $\langle VM\ 1, *, *, \\ \backslash \backslash \text{share} \backslash \text{dataset} \rangle \text{ --- } \text{bypass malware scanner}$

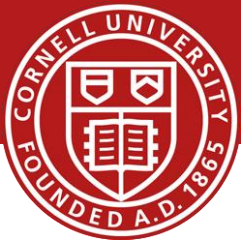
IOFlow API: programming data plane queues



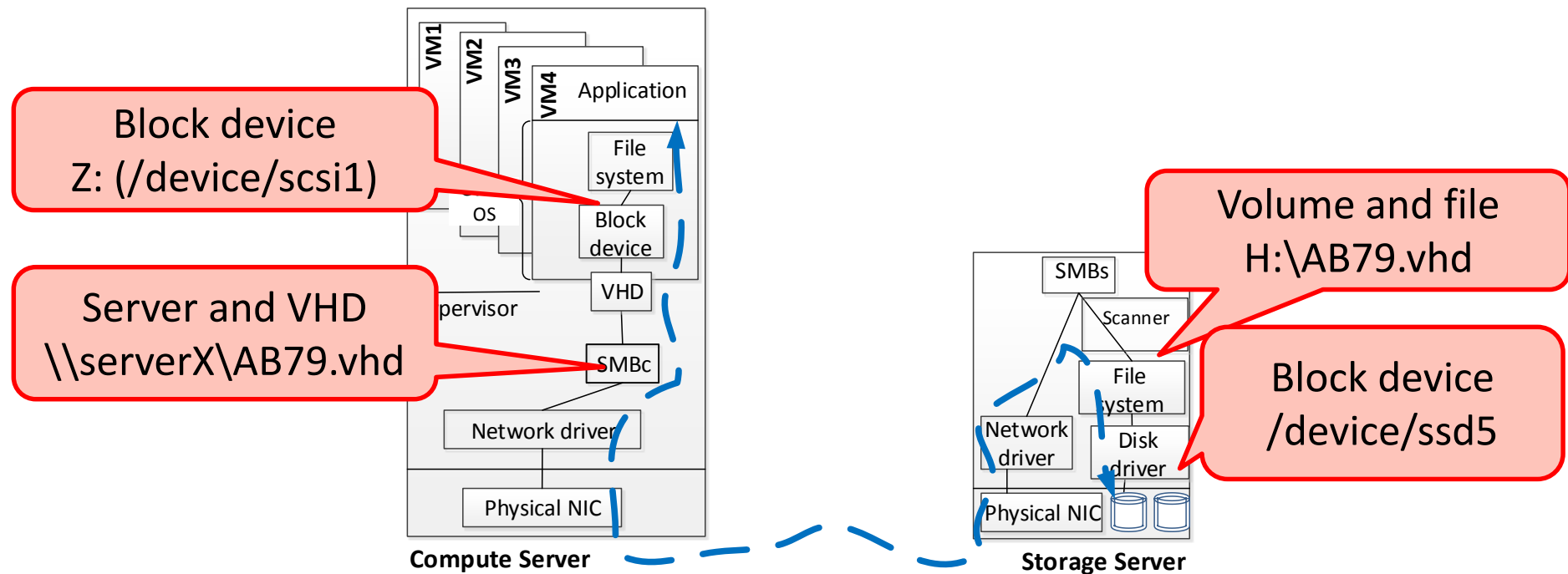
1. Classification [IO Header -> Queue]
2. Queue servicing [Queue -> *<token rate, priority, queue size>*]
3. Routing [Queue -> Next-hop]



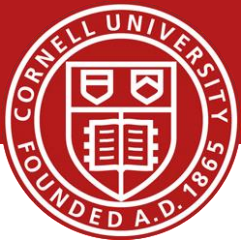
Lack of common IO Header for storage traffic



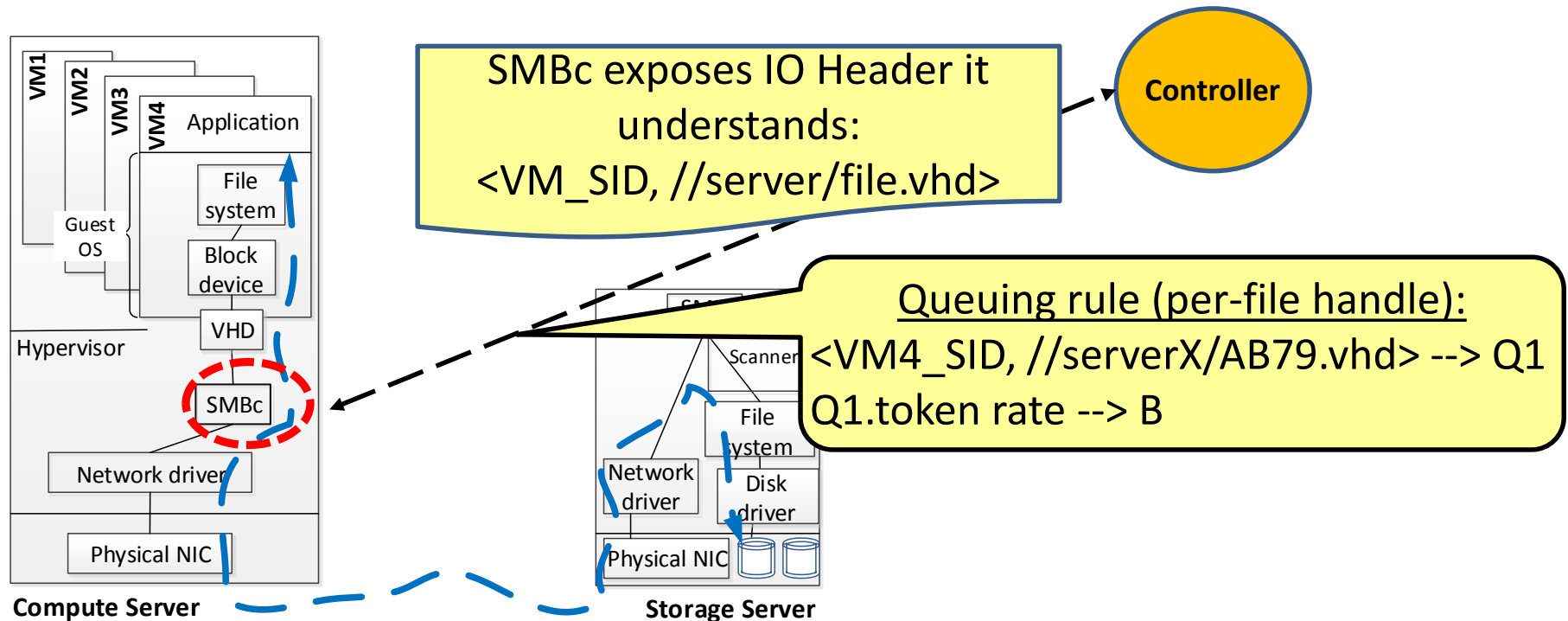
- SLA: <VM 4, *, *, \\share\dataset> --> Bandwidth B



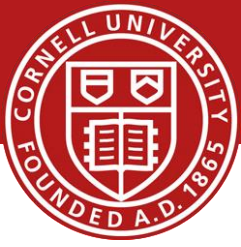
Flow name resolution through controller



- SLA: {VM 4, *, *, //share/dataset} --> Bandwidth B

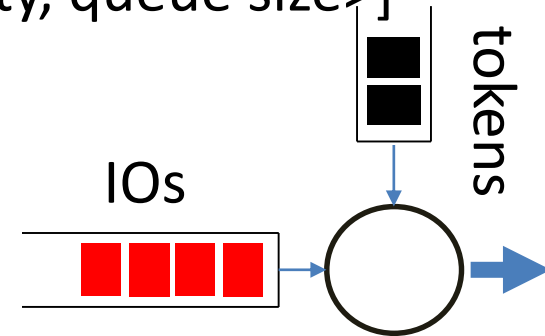


Rate limiting for congestion control



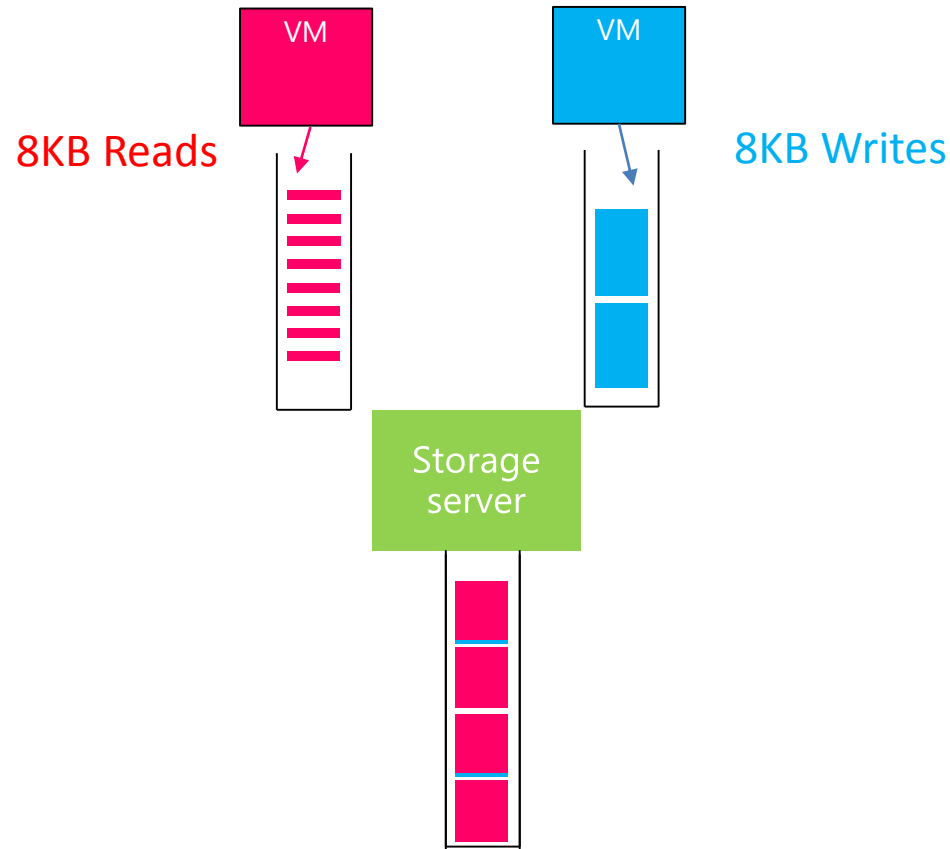
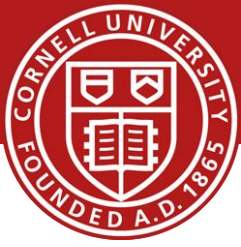
Queue servicing [Queue -> <token rate, priority, queue size>]

- Important for performance SLAs
- Today: no storage congestion control

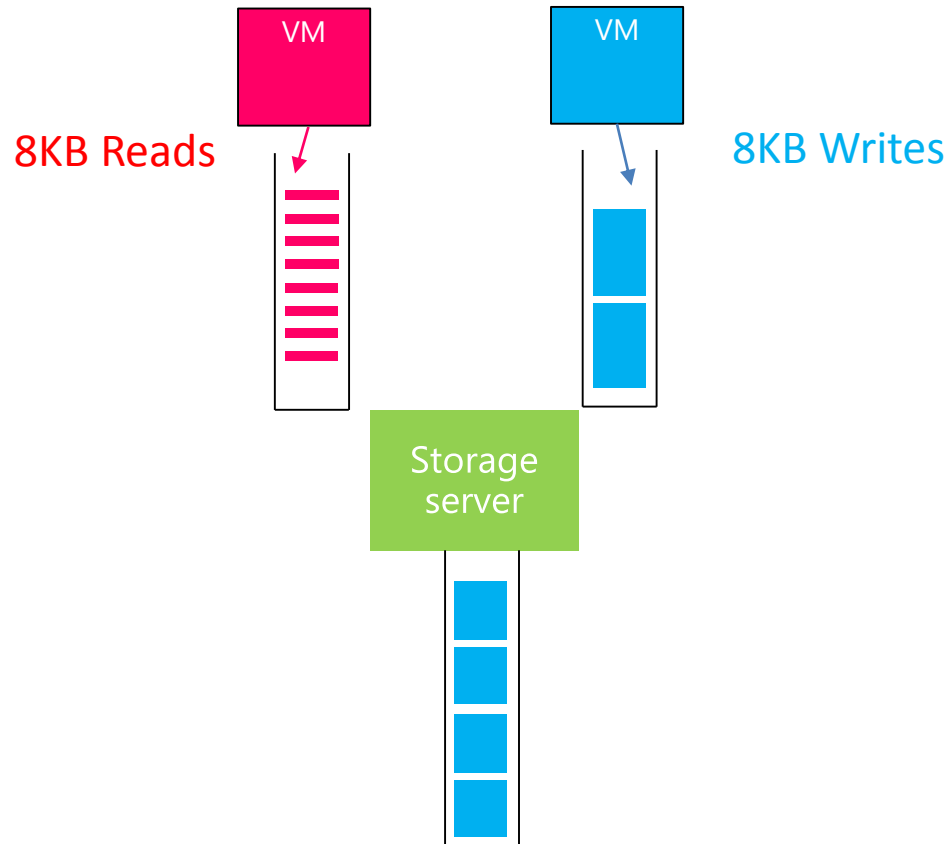
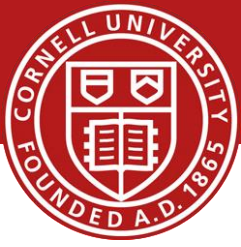


- Challenging for storage: e.g., how to rate limit two VMs, one reading, one writing to get equal storage bandwidth?

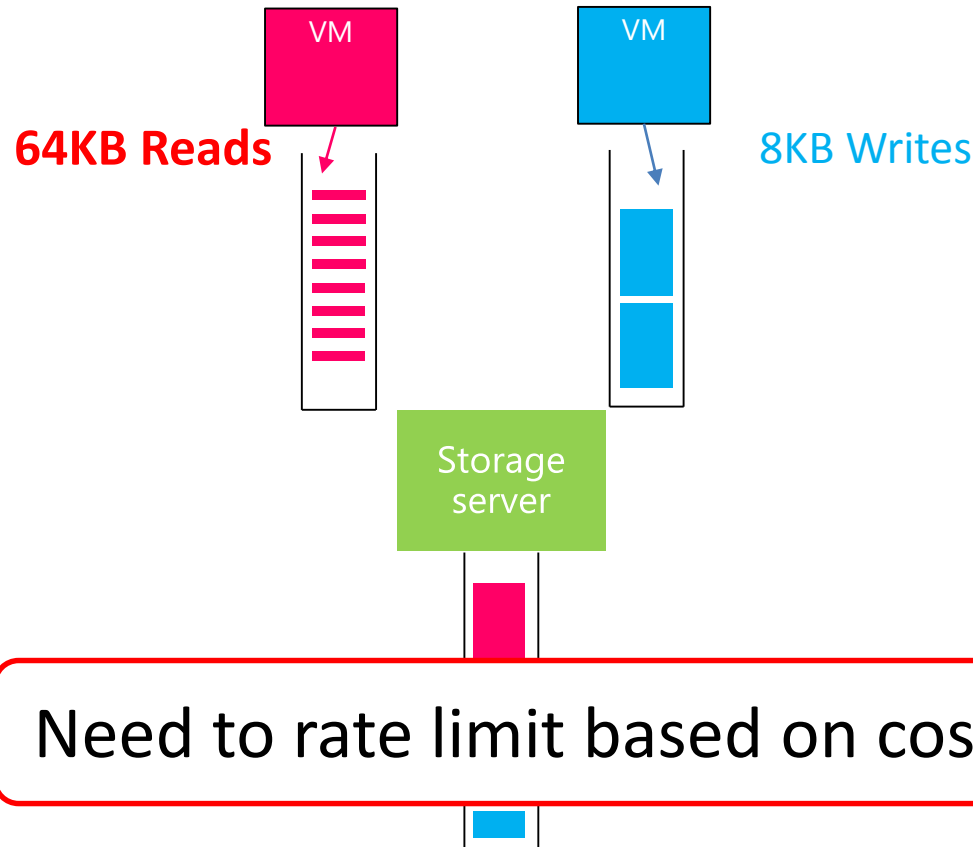
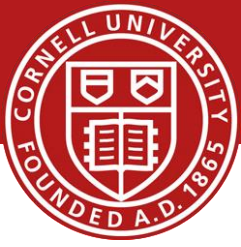
Rate limiting on payload bytes does not work



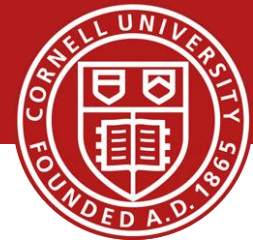
Rate limiting on bytes does not work



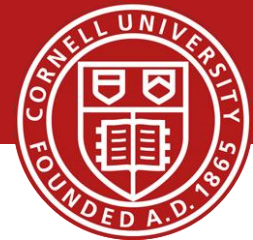
Rate limiting on IOPS does not work



Rate limiting based on cost



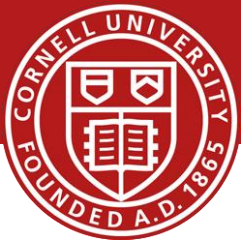
- Controller constructs empirical cost models based on device type and workload characteristics
 - RAM, SSDs, disks: read/write ratio, request size
- Cost models assigned to each queue
 - *ConfigureTokenBucket [Queue -> cost model]*
- Large request sizes split for pre-emption



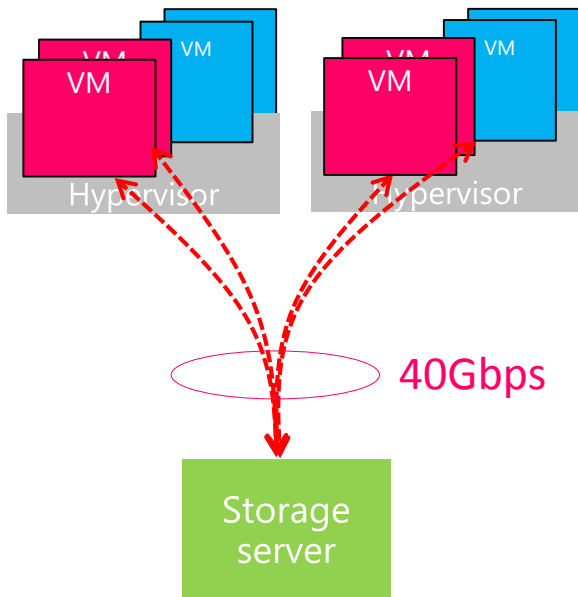
- Classification [IO Header -> *Queue*]
 - Per-layer metadata exposed to controller
 - Controller out of critical path
- Queue servicing [Queue -> *<token rate, priority, queue size>*]
 - Congestion control based on operation cost
- Routing [Queue -> *Next-hop*]

How does controller enforce SLA?

Distributed, dynamic enforcement

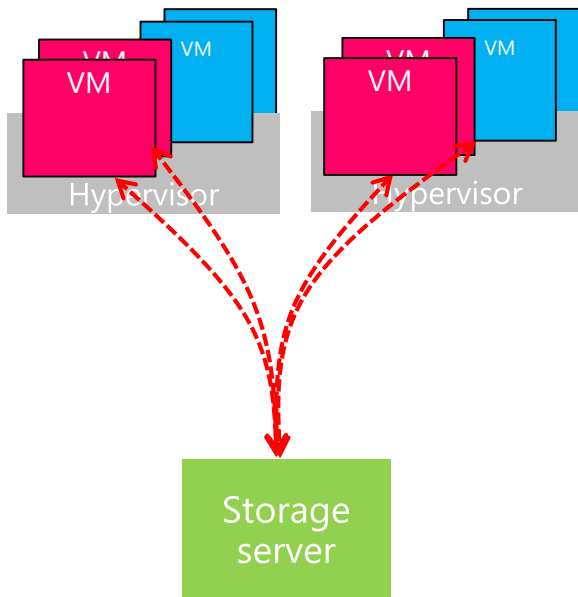
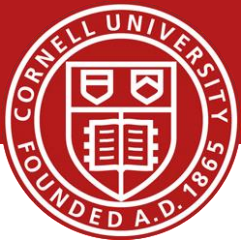


- $\langle \{\text{Red VMs 1-4}\}, *, * // \text{share/dataset} \rangle \rightarrow \text{Bandwidth 40 Gbps}$



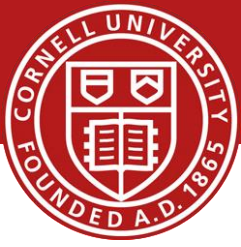
- SLA needs per-VM enforcement
- Need to control the aggregate rate of VMs 1-4 that reside on different physical machines
- Static partitioning of bandwidth is sub-optimal

Work-conserving solution



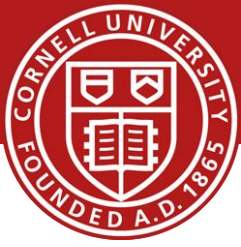
- VMs with traffic demand should be able to send it as long as the aggregate rate does not exceed 40 Gbps
- **Solution:** *Max-min fair sharing*

Max-min fair sharing



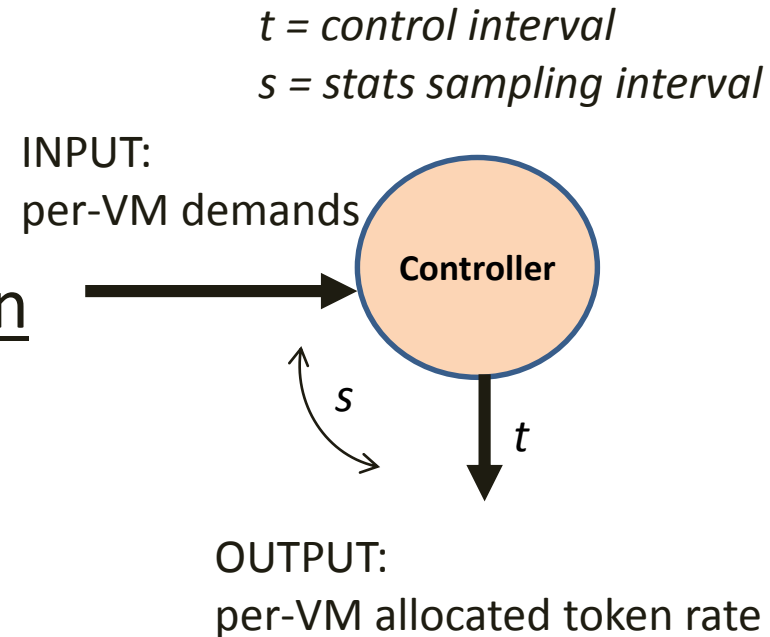
- Well studied problem in networks
 - Existing solutions are distributed
 - Each VM varies its rate based on congestion
 - Converge to max-min sharing
 - *Drawbacks*: complex and requires congestion signal
- But we have a centralized controller
 - Converts to simple algorithm at controller

Controller-based max-min fair sharing

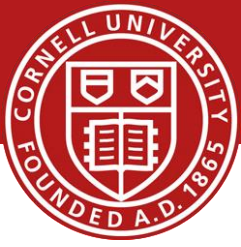


What does controller do?

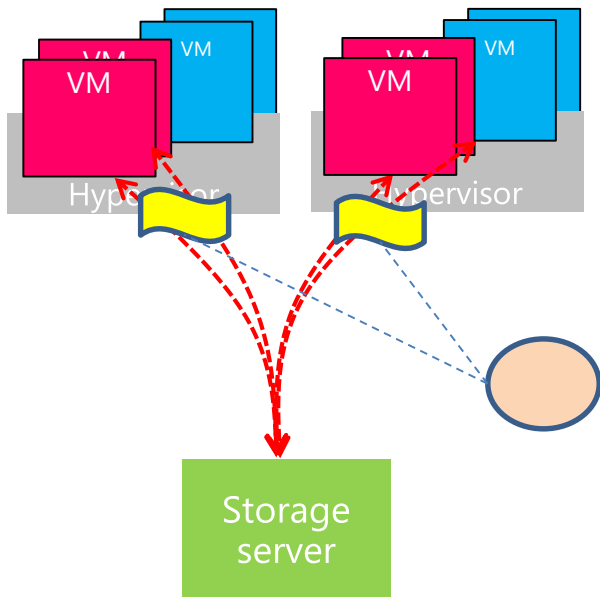
- Infers VM demands
- Uses centralized max-min within a tenant and across tenants
- Sets VM token rates
- Chooses best place to enforce



Controller decides *where* to enforce



Minimize # times IO is queued and distribute rate limiting load



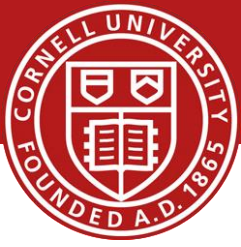
SLA constraints

- Queues where resources shared
- Bandwidth enforced close to source
- Priority enforced end-to-end

Efficiency considerations

- Overhead in data plane \sim # queues
- Important at 40+ Gbps

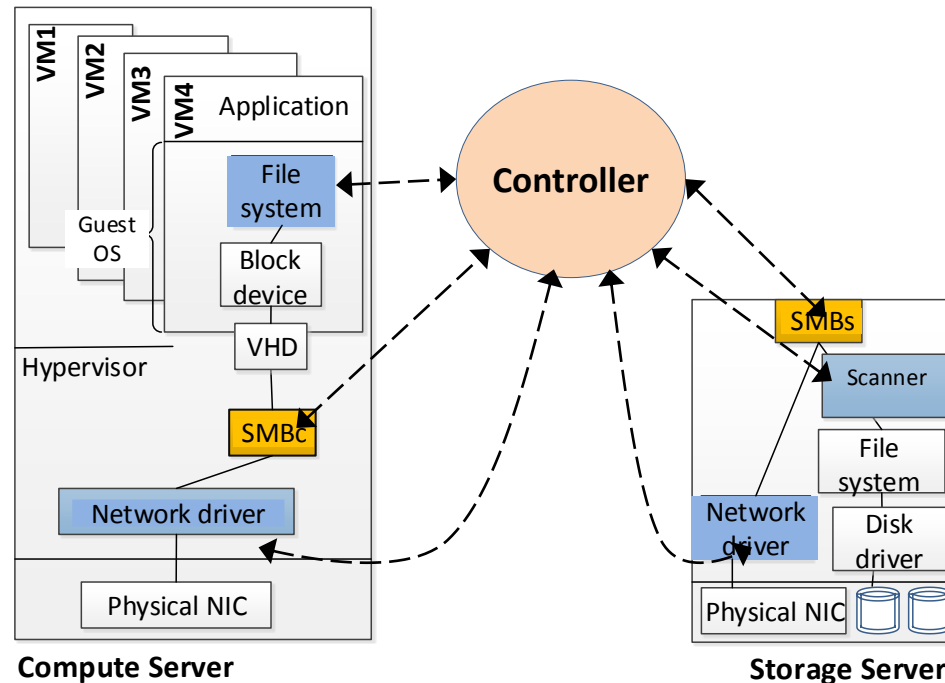
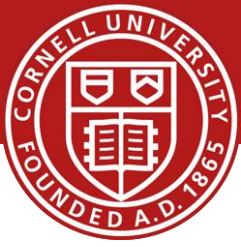
Centralized vs. decentralized control



Centralized controller in SDS allows for simple algorithms that focus on SLA enforcement and ***not*** on distributed system challenges

Analogous to benefits of centralized control in software-defined networking (SDN)

IOFlow implementation



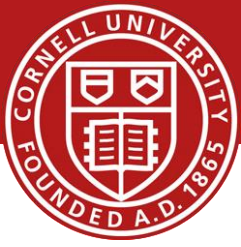
2 key layers for VM-to-Storage performance SLAs

4 other layers

- . Scanner driver (routing)
- . User-level (routing)
- . Network driver
- . Guest OS file system

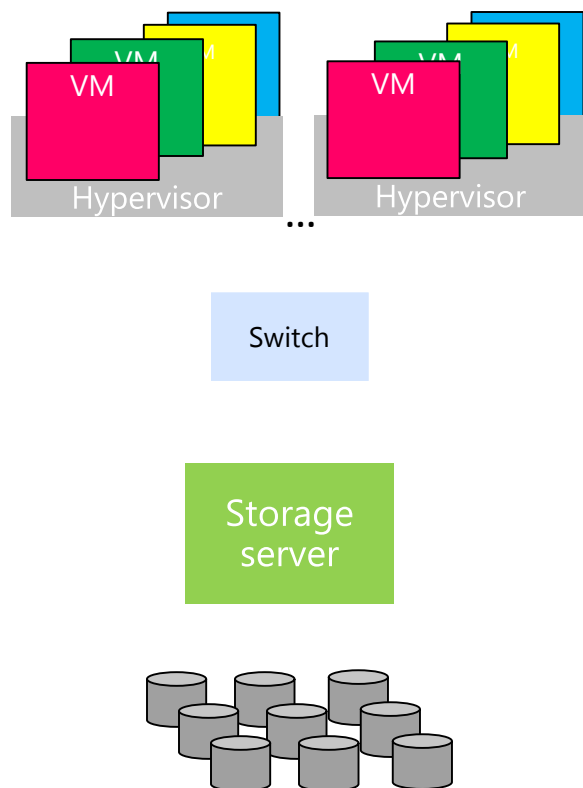
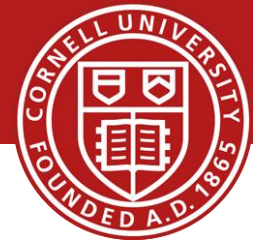
Implemented as filter drivers on top of layers

Evaluation map



- IOFlow's ability to enforce end-to-end SLAs
- Aggregate bandwidth SLAs
- Priority SLAs and routing application in paper
- Performance of data and control planes

Evaluation setup



Clients: 10 hypervisor servers, 12 VMs each
4 tenants (Red, Green, Yellow, Blue)
30 VMs/tenant, 3 VMs/tenant/server

Storage network:

Mellanox 40Gbps RDMA RoCE full-duplex

1 storage server:

16 CPUs, 2.4GHz (Dell R720)

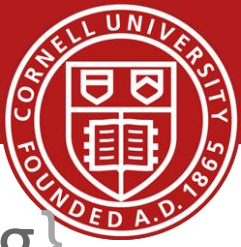
SMB 3.0 file server protocol

3 types of backend: RAM, SSDs, Disks

Controller: 1 separate server

1 sec control interval (configurable)

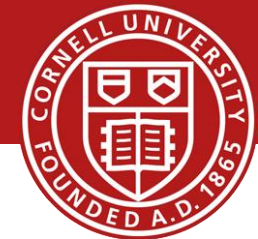
Workloads



- 4 Hotmail tenants {Index, Data, Message, Log}
- Used for trace replay on SSDs (see paper)
- IoMeter is parametrized with Hotmail tenant characteristics (read/write ratio, request size)

	Index	Data	Message	Log
Read %	75%	61%	56%	1%
IO Sizes	4/64 KB	8 KB	4/64 KB	0.5/64 KB
Seq/rand	Mixed	Rand	Rand	Seq
# IOs	32M	158M	36M	54M

Enforcing bandwidth SLAs

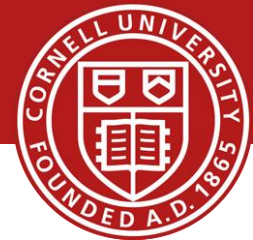


4 tenants with different storage bandwidth SLAs

Tenants have		
Tenant		SLA
▪ Red	Red	{VM1 – 30} -> Min 800 MB/s
	Green	{VM31 – 60} -> Min 800 MB/s
	Yellow	{VM61 – 90} -> Min 2500 MB/s
	Blue	{VM91 – 120} -> Min 1500 MB/s

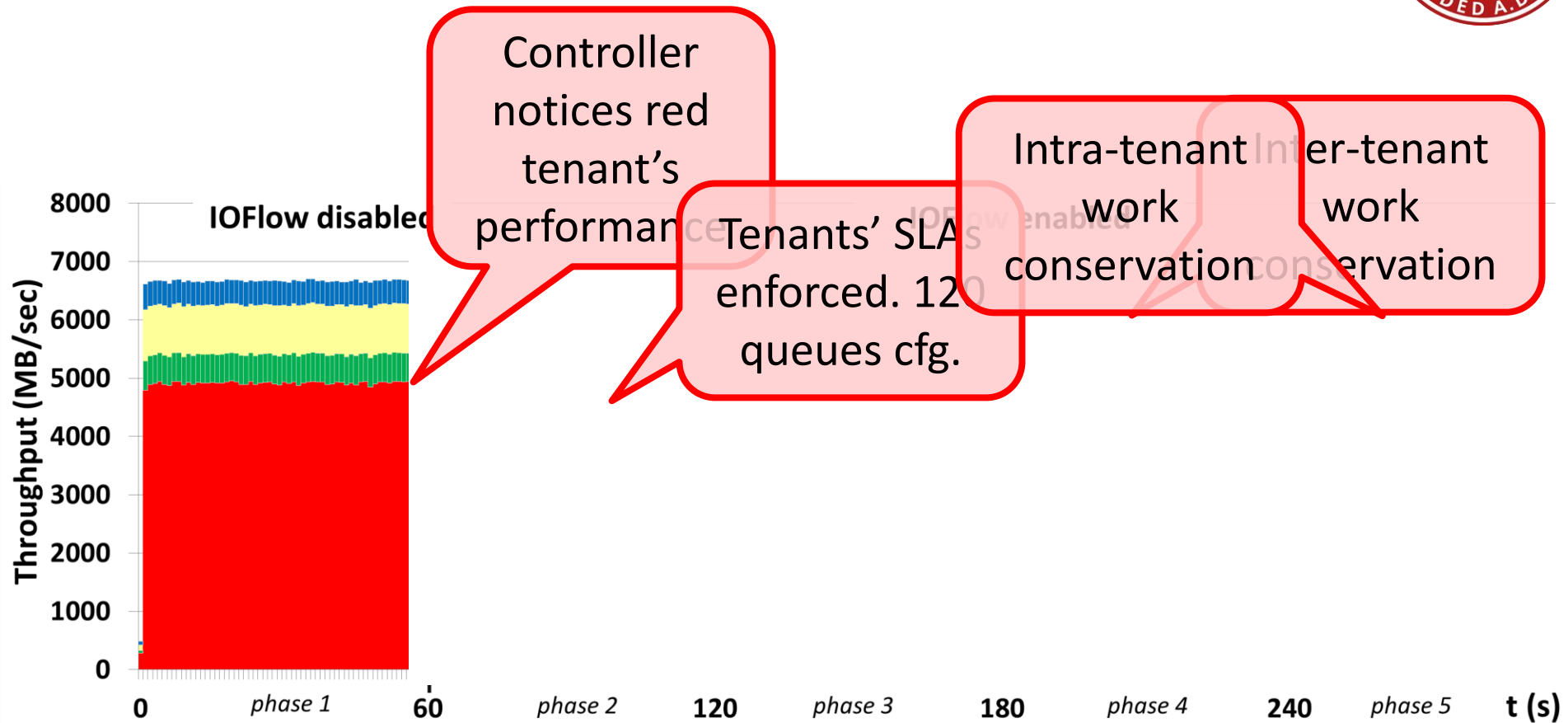
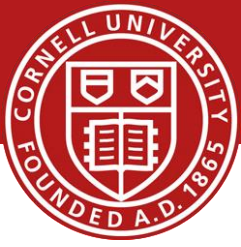
/second

Things to look for

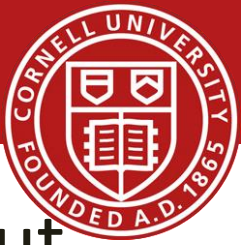


- Distributed enforcement across 4 competing tenants
 - Aggressive tenant(s) under control
- Dynamic inter-tenant work conservation
 - Bandwidth released by idle tenant given to active tenants
- Dynamic intra-tenant work conservation
 - Bandwidth of tenant's idle VMs given to its active VMs

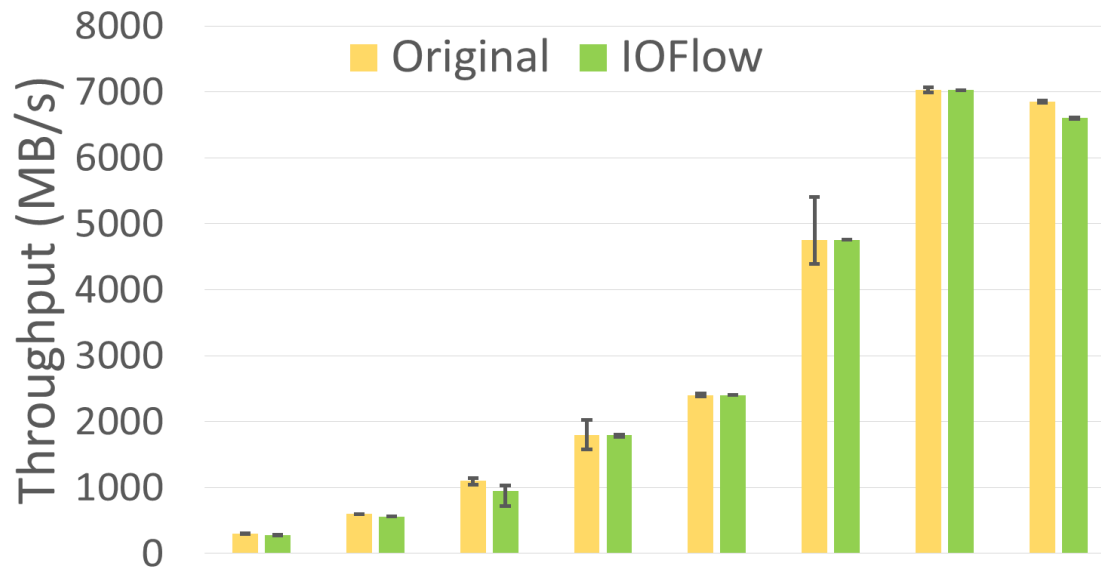
Results



Data plane overheads at 40Gbps RDMA

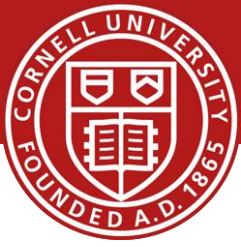


- Negligible in previous experiment. To bring out worst case varied IO sizes from 512Bytes to 64KB

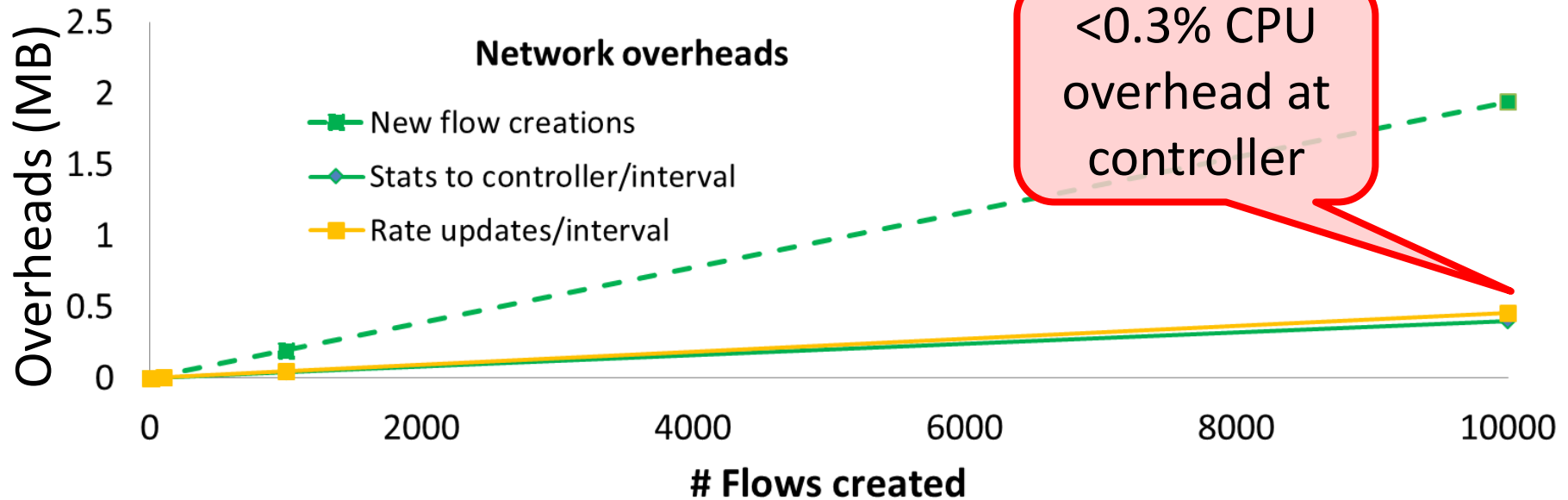


Reasonable overheads for enforcing SLAs

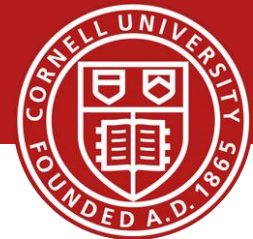
Control plane overheads: network and CPU



- Controller configures queue rules, receives statistics and updates token rates every interval



Before Next time



- Final Project Presentation/Demo
 - Due **Friday, December 12.**
 - Presentation and Demo
 - Written submission required:
 - Report
 - Website: index.html that points to report, presentation, and project (e.g. code)
- ***Required review and reading for Wednesday, December 3***
 - Plug into the Supercloud, D. Williams, H. Jamjoom, H. Weatherspoon. *IEEE Internet Computing*, Vol. 17, No 2, March/April 2013, pp 28-34.
 - <http://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=6365162>
- Check piazza: <http://piazza.com/cornell/fall2014/cs5413>
- Check website for updated schedule