# DETAIL: REDUCING THE FLOW COMPLETION TIME TAIL IN

### **DATACENTER NETWORKS**

David Zats, Tathagata Das, Prashanth Mohan,

Dhruba Borthakur, Randy Katz

Presented by Alexander Pokluda

February 6, 2013



.

\_

## Sophisticated Web Applications

- Rendering a page may require hundreds of requests to back-end servers
- Strict page rendering deadlines of 200-300ms must be met to ensure a positive user experience

	Courts for words along and films	
acebook 🔟 🖛 🗺	Search for people, places and things	Alexander Pokluda Home
Alexander Pokluda	Update Status 👔 Add Photos/Video	
Edit Profile	What's happening, Alexander?	Prisoners' Justice F on February 8
VORITES		2 requests from Sara Thompson
News Feed	SORT +	5 other app requests
Messages	Mike Wareing	Find More Friends
Events 6	Very excited to welcome my new niece, Aria Karen Wareing. Seven	Alexander, Try the Friend Finder
Photos	like : Comment - Yesterday at 12:16am via mobile - 18	Find friends the easy way
OUPS	r <sup>A</sup> bile Wareing James (1) early and 23 others like this	using the automatic friend
London Activist Net 20+	Lump Marsing And ure ture utilities of the the head Can't wait	ande.
Call for Toronto Polic 20+	for your announcement also!	Your Email
Ontario Clean Air Alli 20+	Yesterday at 8:06am - Like - 🖒 1	Email Password
Create Group	Lori Atkinson congratulations!! Her name is beautiful!	Find Friends
PPS	Yesterday at 9:23am 'Like ' 40 1	
D App Center	Sharon Gordon Congrats Unde Michael and Aunt Julie!	People You May Know See
Games Feed 20+	Testerday at 10:50am - Like - KO 1	Behzad Esfan (Bernard)

# **Network Complications**

- Typically a few responses arrive late giving us long tailed flow completion times
- Web applications must choose between sacrificing either quality or responsiveness
- Either option leads to financial loss



### **Network Performance Factors**

- Application workflows depend on performance of underlying network flows
- Congestion can cause round-trip-times (RTTs) to form a long-tailed distribution
- Congestion leads to
  - Packet loss and retransmissions
  - Uneven load balancing
  - Priority inversion
- Each contributes to increasing long tail of flow completion, especially for latency-sensitive short flows critical for page creation



### **Reducing the Flow Completion Time Tail**

- Flash congestion can be reduced if it can be detected early enough
- DeTail addresses this challenge by constructing a cross-layer network stack that detects congestion at lower layers to drive upper

# Contributions

 $\prec$ 



 Assessment of the causes of the long-tailed flow completion times

• A cross-layer network stack that addresses them

• Implementation-validated simulations demonstrating DeTail's significant improvement



### **Traffic Measurements**

 Intra-rack RTTs are typically low but congestion can cause them to vary by two orders of magnitude



### Impact on Workflows

### Partition-Aggregate

 At the 99.9<sup>th</sup> percentile, a 40-worker flow has 4 workers (10%) miss their 10*ms* deadlines while a 400-worker flow has 14 (3.5%) miss theirs

### Sequential

 At the 99.9<sup>th</sup> percentile, web sites must have less than 150 sequential data retrievals per page to meet 200*ms* page creation deadlines

Based on published datacenter traffic measurements for production networks

While events at the long tail occur rarely, workflows use so many flows that several will experience delays for every page creation A network that reduces the tail allows applications to render more complete pages without increasing server load



-

-

### **Cross-layer Network-based Approach**

Layer	Components	Info Exchanged
Application		Flow Priority
Transport	Reorder-Resistant Transport	
Network	Adaptive Load Balancing	Congestion Notification
Link	Lossless Fabric	
Physical		 Port Occupancy
Physical		Occupancy

SIMULATION, IMPLEMENTATION AND EXPERIMENTAL RESULTS

# Simulation and Implementation



- Simulation using CIOQ switch architecture in NS-3 Network Simulator
- NS-3 extended to include real-world processing delays
- NS-3 does not support ECN, but simulations still demonstrate impressive results



Modular

 Click software modified to have both ingress and egress queues

Router

 Rate limiters added to prevent packet buildup in driver and hardware buffers

## **Experimental Results**

To evaluate DeTail's ability to reduce the flow completion time tail, the following approaches are compared:

#### Flow Hashing (FH)

- Switches employ flow-level hashing
- Status quo and baseline

#### Lossless Packet Scatter (LPS)

- Switches employ packet scatter with PFC
- Not standard but can be deployed in current datacenters

#### DeTail

- Switches employ PFC and Adaptive Load Balancing (ALB)
- New and exciting!

Simulator predictions are closely matched by implementation measurements! The simulator is used to evaluate larger topologies and wider range of workflows

### Microbenchmarks: All-to-All Workload

- FatTree topology with 128 servers in 4 pods with 4 ToR and 4 aggregate switches each
- Each server randomly retrieves data from another

99%

Completion Time

% Reduction in

100

80

60

40

20

Servers also engaged in low-priority background flows



CDF of completion times of 8*KB* data retrievals at 2000 retrievals/second



99%



Reduction by DeTail over FH in 99<sup>th</sup> and 99.9<sup>th</sup> percentile completion times of 2*KB*, 8*KB* and 32*KB* retrievals. DeTail provides up to 70% reduction at the 99<sup>th</sup> percentile.

# Microbenchmarks: Front-end/Back-end Workload

- Same FatTree topology as before
- Servers in first three pods retrieve data from randomly chosen servers in fourth pod
- Servers also engaged in low-priority background flows



Reduction by DeTail over FH in 99<sup>th</sup> and 99.9<sup>th</sup> percentile completion times of 2*KB*, 8*KB* and 32*KB* retrievals. DeTail achieves 30% - 65% reduction in completion times at the 99.9<sup>th</sup> percentile.

# **Topological Asymmetries**

#### **Disconnected Link**

 Same as all-to-all workload but with one disconnected aggregate to core link



DeTail provides 10% - 89% reduction—almost an order of magnitude improvement—compared to FH for 8*KB* retrievals

#### **Degraded Link**

 Same as all-to-all workload but with one 1*Gpbs* downgraded to 100*Mbps*



DeTail provides 91% reduction compared to FH for 8KB retrievals

# Web Workloads: Sequential

- Servers randomly assigned to be frond-end or back-end
- Front-end servers retrieve data from randomly chosen back-end servers
- Each sequential workflow consists of 10 sequential data retrievals of 2KB, 4KB, 8KB, 16KB or 32KB



DeTail provides 71% - 76% reduction in 99.9<sup>th</sup> percentile completion times of individual data retrievals and 54% reduction overall

### Web Workloads: Partition-Aggregate

- Servers randomly assigned to be frond-end or back-end
- Front-end servers retrieve data in parallel from randomly chosen back-end servers
- Each partition-aggregate workflow consists of 10, 20, or 40 data retrievals 2KB in size



DeTail provides 78 - 88% reduction in 99.9<sup>th</sup> percentile completion times



### **Related Work**

#### **Internet Protocols**

- TCP Modifications: NewReno, Vegas, SACK
- Buffer Management: RED and Fair Queuing
- Operate at coarse-grained timescales inappropriate for datacenter workloads

#### Datacenter Networks

- Topologies: FatTrees, VL2, BCube, DCell
- Traffic Management: DCTCP, Hull, D<sup>3</sup>, Datacenter Bridging
- Bound by performance of flow hashing

#### **HPC Interconnects**

- Credit-based flow control
- Adaptive Load Balancing: UGAL, PAR
- These mechanisms have not been evaluated for web-facing datacenter networks

### Summary

DeTail is an approach for *reducing the tail completion times* of short, latency sensitive flows critical for page creation

DeTail employs cross-layer, in-network mechanisms to reduce packet losses, prioritize flows, and balance traffic

By making its flow completion statistics robust to congestion, DeTail can reduce 99.9<sup>th</sup> percentile flow completion times by over 50% for many workloads



.



.

-

## **Photo Credits**

- Railroad crossing: Toledo Blade
  - http://www.toledoblade.com/frontpage/2008/03/04/Railroad-crossing-barriers-tested-in-Michigan.html
- Pin-the-Tail-on-the-Donkey: The City Patch
  - <u>http://thecitypatch.com/2012/04/02/pin-the-tail-on-the-donkey-will-never-quite-be-the-same/</u>
- Clip-art from Office.com