

# pHost

## Distributed Near-Optimal Datacenter Transport Over Commodity Network Fabric

**Peter X. Gao**, Akshay Narayan, Gautam Kumar  
Rachit Agarwal, Sylvia Ratnasamy, Scott Shenker  
UC Berkeley/ICSI

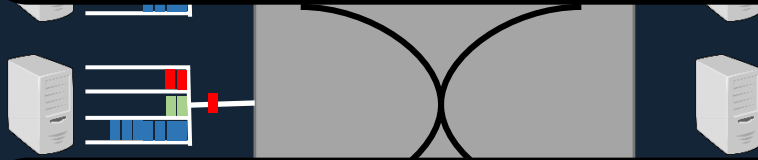
# Datacenter Transport Protocol Designs

Fastpass Silo  
D3 DeTail D2TCP PDQ  
PIAS FairCloud Hedera  
pFabric HULL  
DCTCP Timely PASE

- Low Latency (FCT)
- High Throughput
- Deadline
- Fairness

# Performance and Flexibility

Close to optimal  
performance



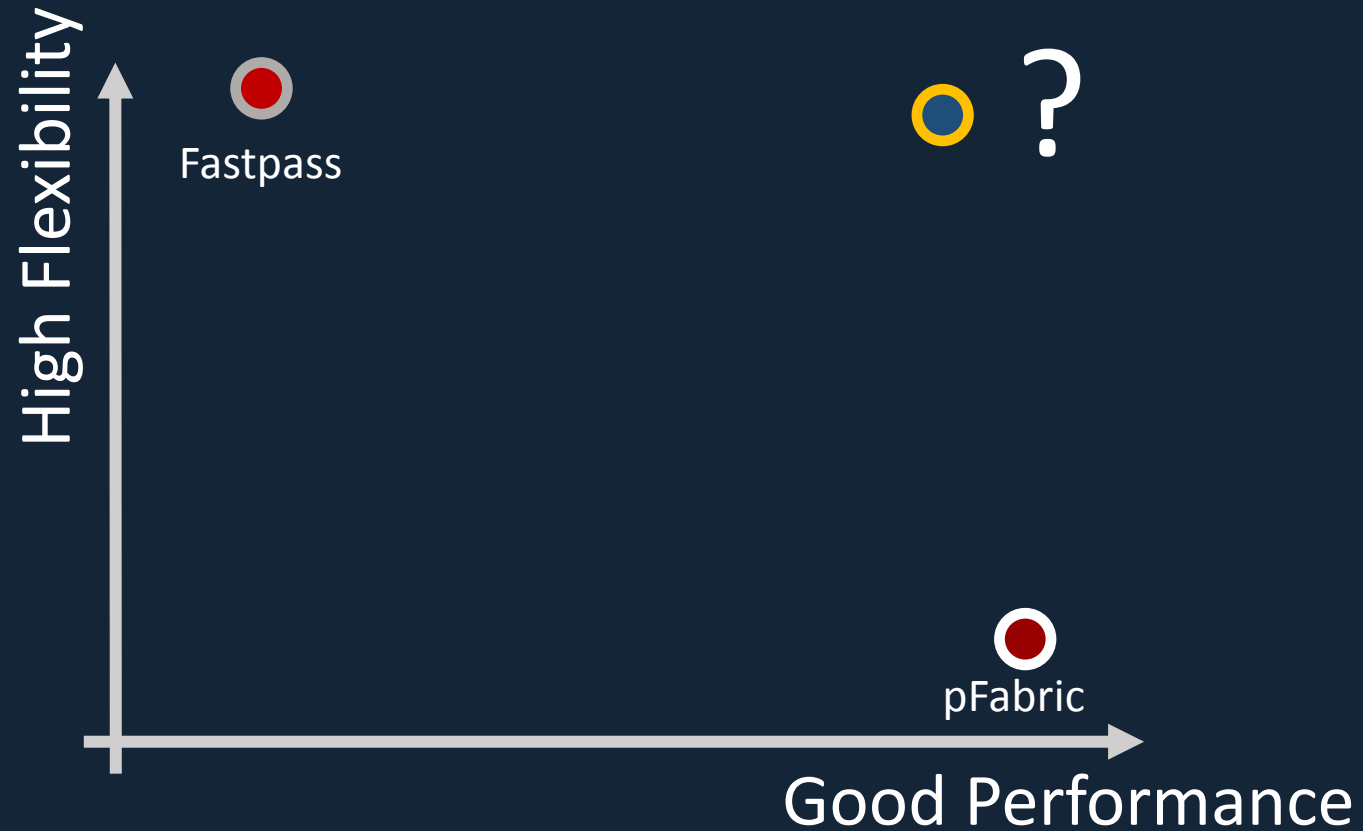
Embed policy inside  
network

Overhead on  
Scheduling Short flows

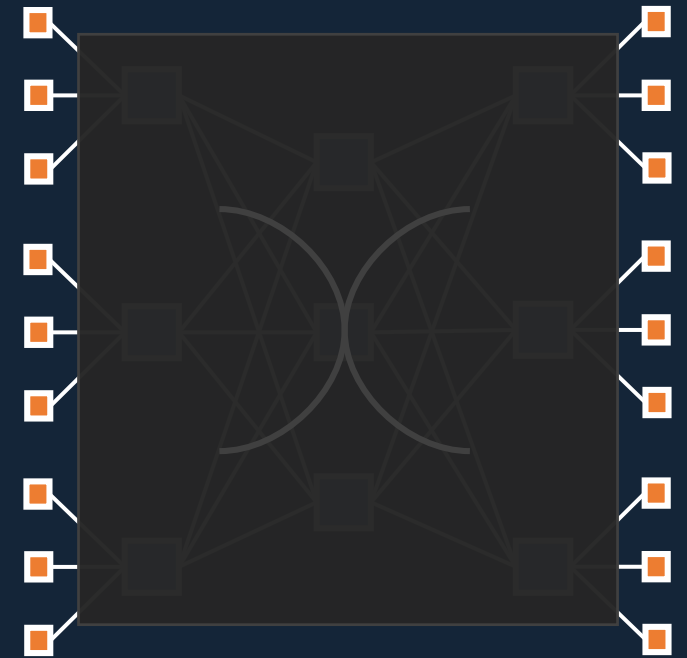
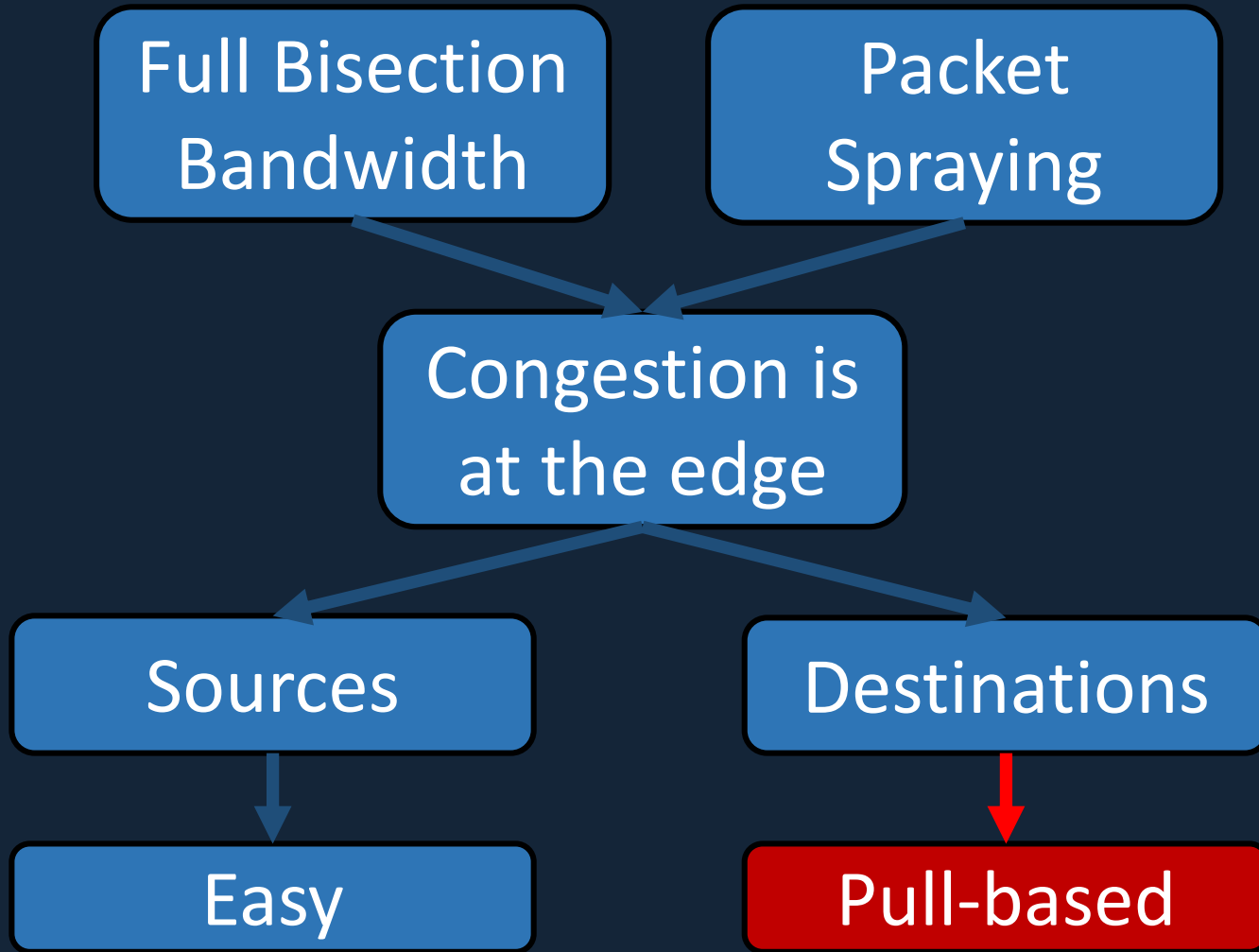


Highly flexible

# Can we get the best of both worlds?



# Design Intuition



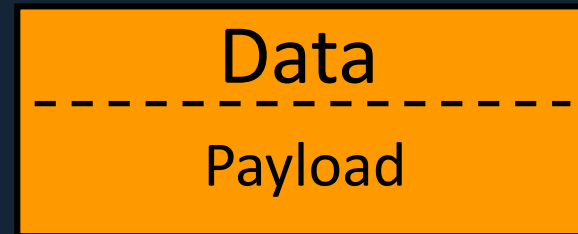
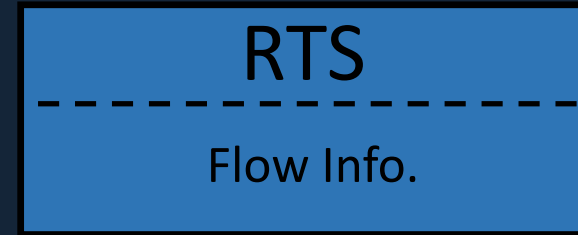
# pHost

- No specialized hardware
- No complex computation inside the fabric
- No centralized scheduler
- No explicit network feedback

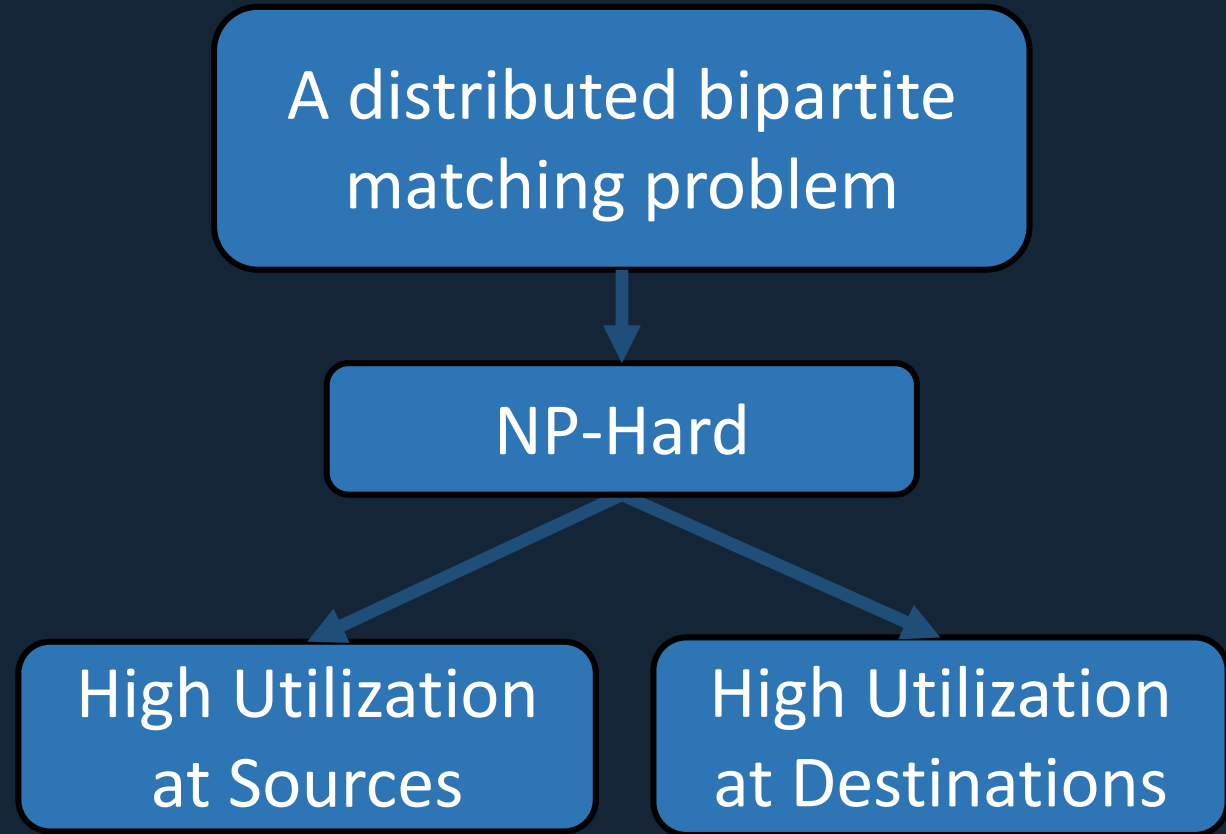
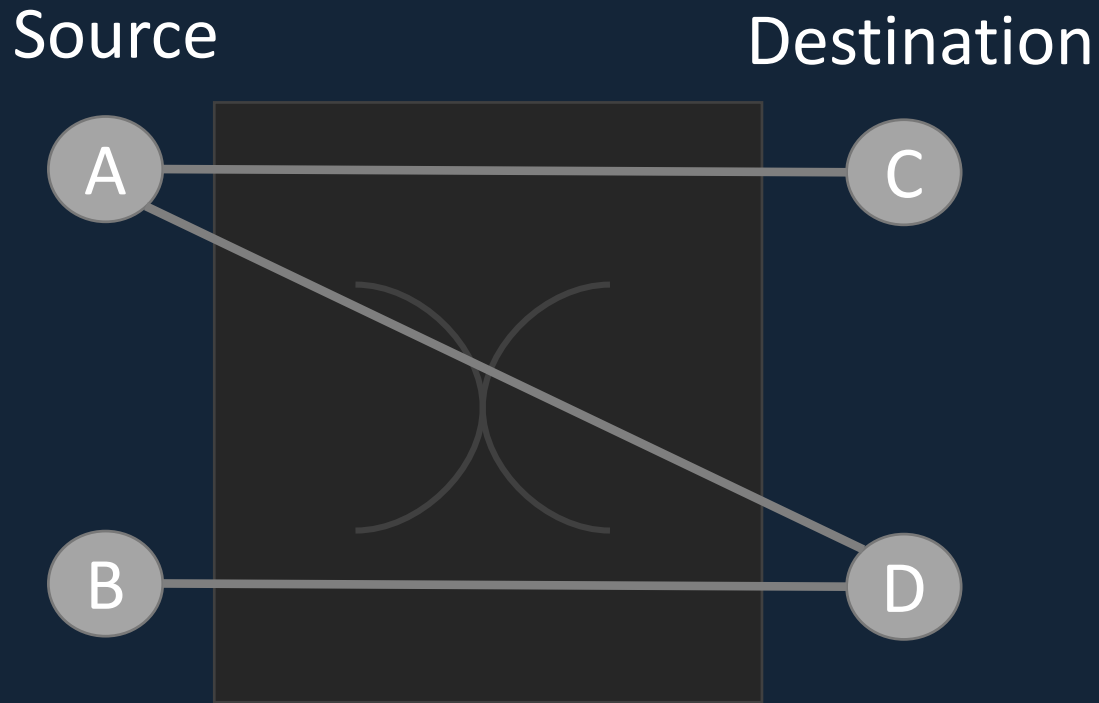
With surprisingly near optimal performance

# pHost Basics

- Source sends **RTS** on flow arrival
- Destination issues **Token** every MTU transmission time
- Source sends **Data** packet when it has a token



# Pull-based End-host Scheduling





# Ensuring high utilization at the sources

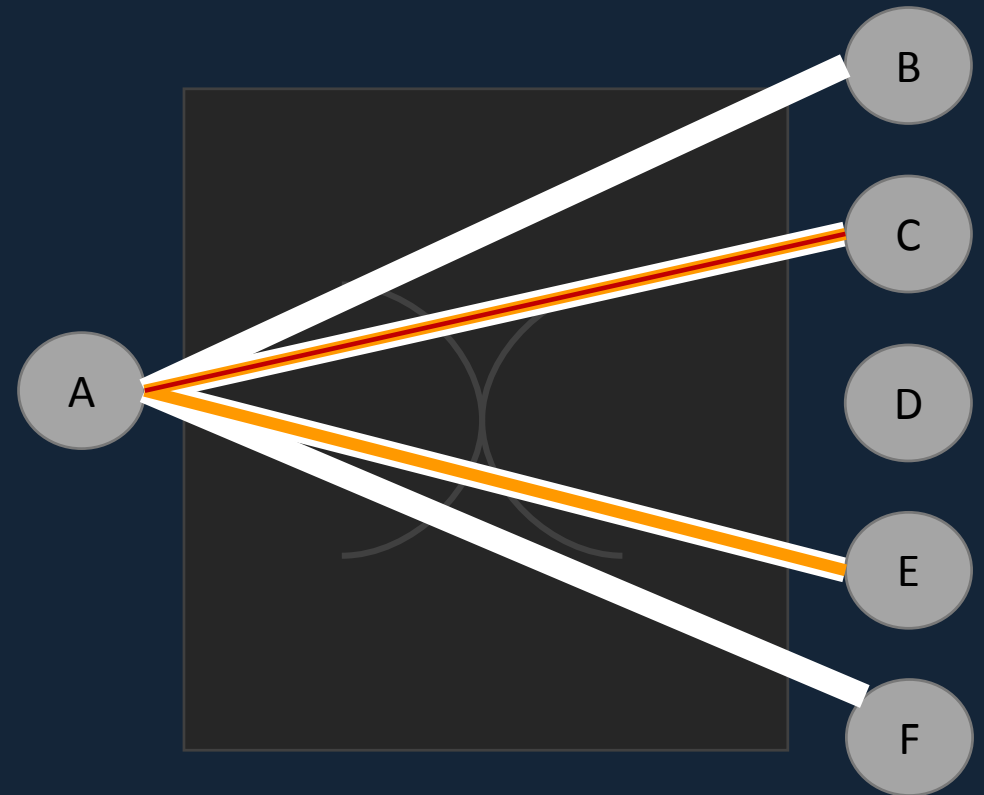
Round Trip Time between  
RTS and first Token



Assign a BDP worth of  
Free Tokens initially

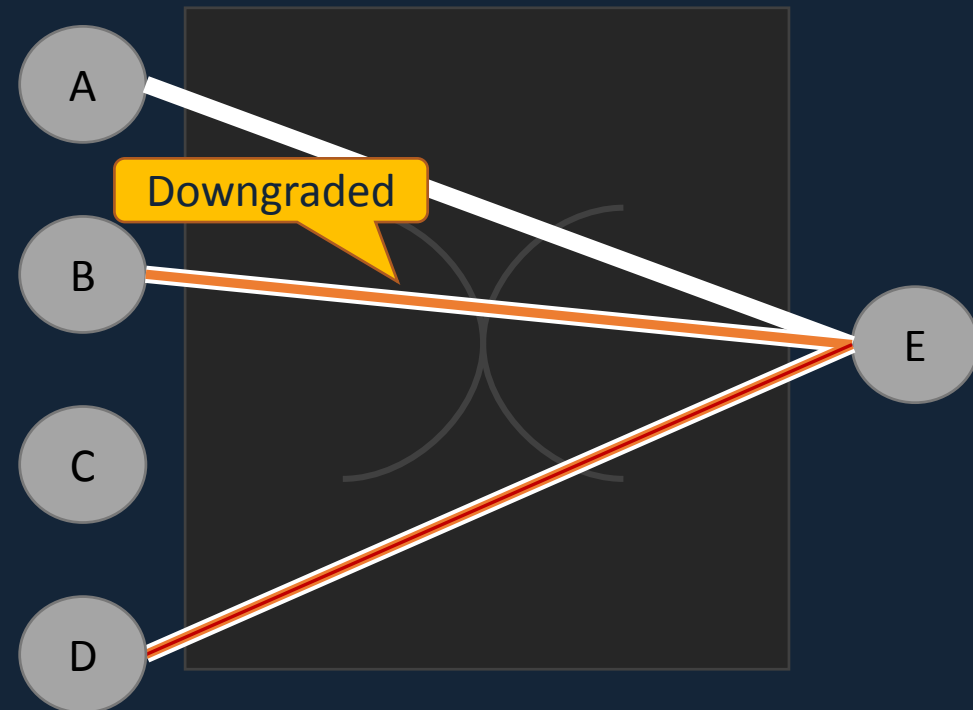
# Ensuring high utilization at the sources

- Allows sources to send multiple RTSes in parallel
- Sources may receive tokens from a **subset of the flows**
- Sources choose **one of the flows** depending on the scheduling policy



# Ensuring high utilization at the destination

- Tokens are **expired** if the source does not consume them
- If tokens issued to a flow expire, the flows are **downgraded** for a while



# Express Policy at Endhosts

- Destination side
  - Using token assignment
- Source side
  - Using token selection

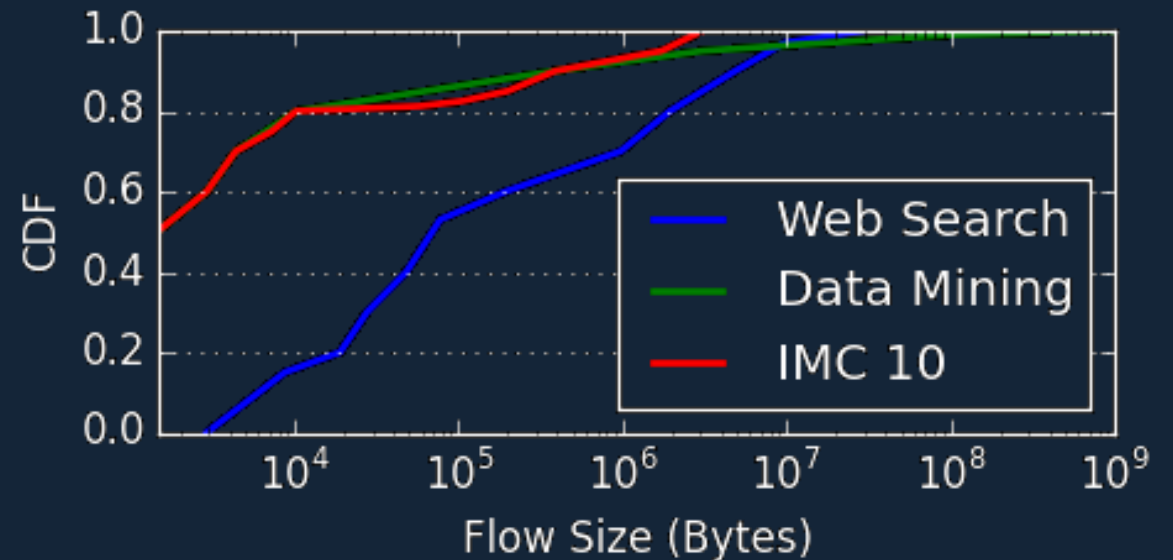
## Examples:

- SRPT
  - Shortest Remaining Flow Size
- Deadline
  - Earliest Deadline First
- Multi-tenant sharing
  - Round Robin / LRU

# Evaluation Setup

- Packet-level Simulator
- Topology
  - Two-tier multi-rooted tree
  - 9 Racks, 144 nodes
  - Cut-through switching
- Mean Slowdown  
 $= \text{mean}(FCT/OPT)$

- Workloads:



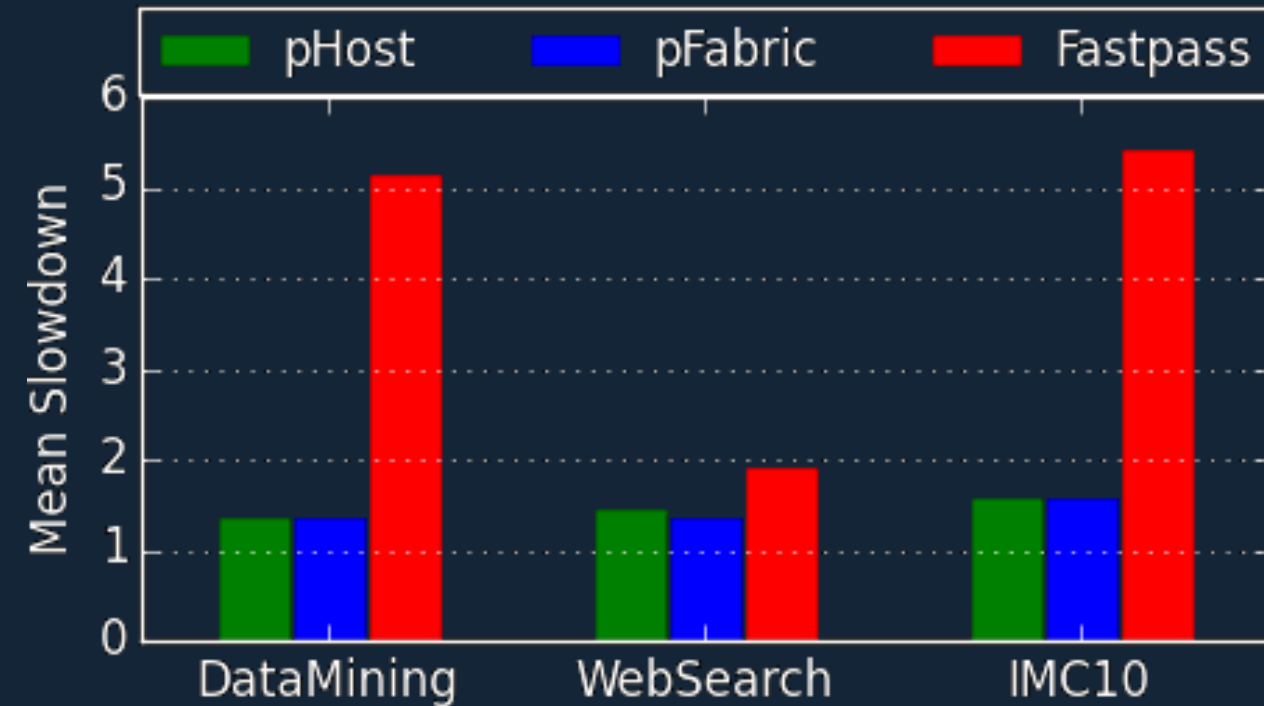
# Evaluation Outline

Can pHost match pFabric's performance?

Is pHost robust?

Is pHost flexible?

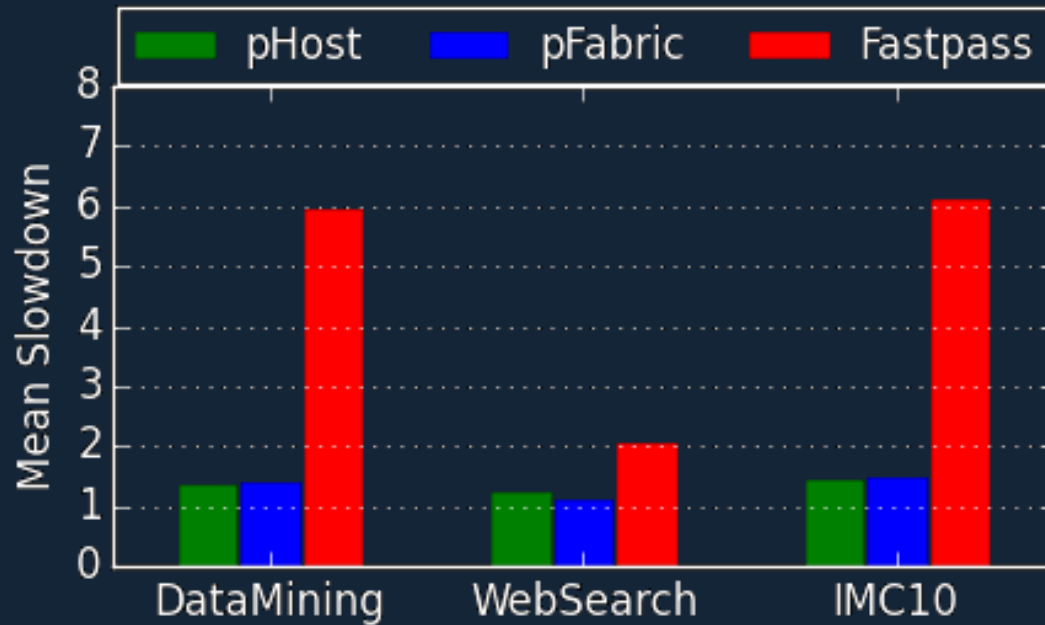
# Overall Mean Slowdown



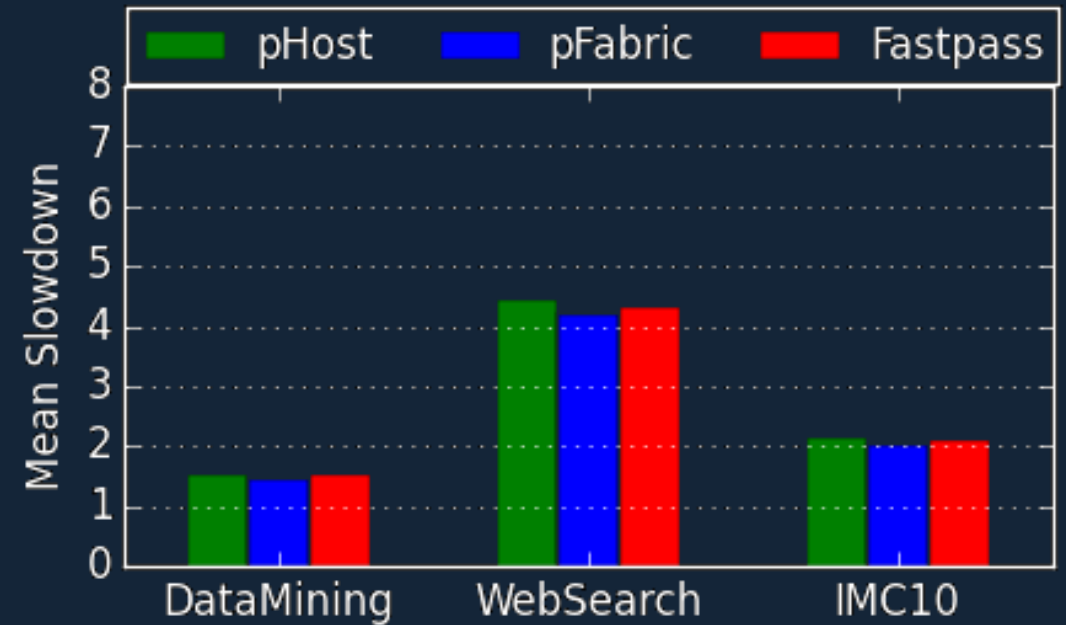
pHost matches pFabric  
and 4x better than Fastpass

pHost uses commodity  
switches

# Digging Deeper



Short Flow Slowdown



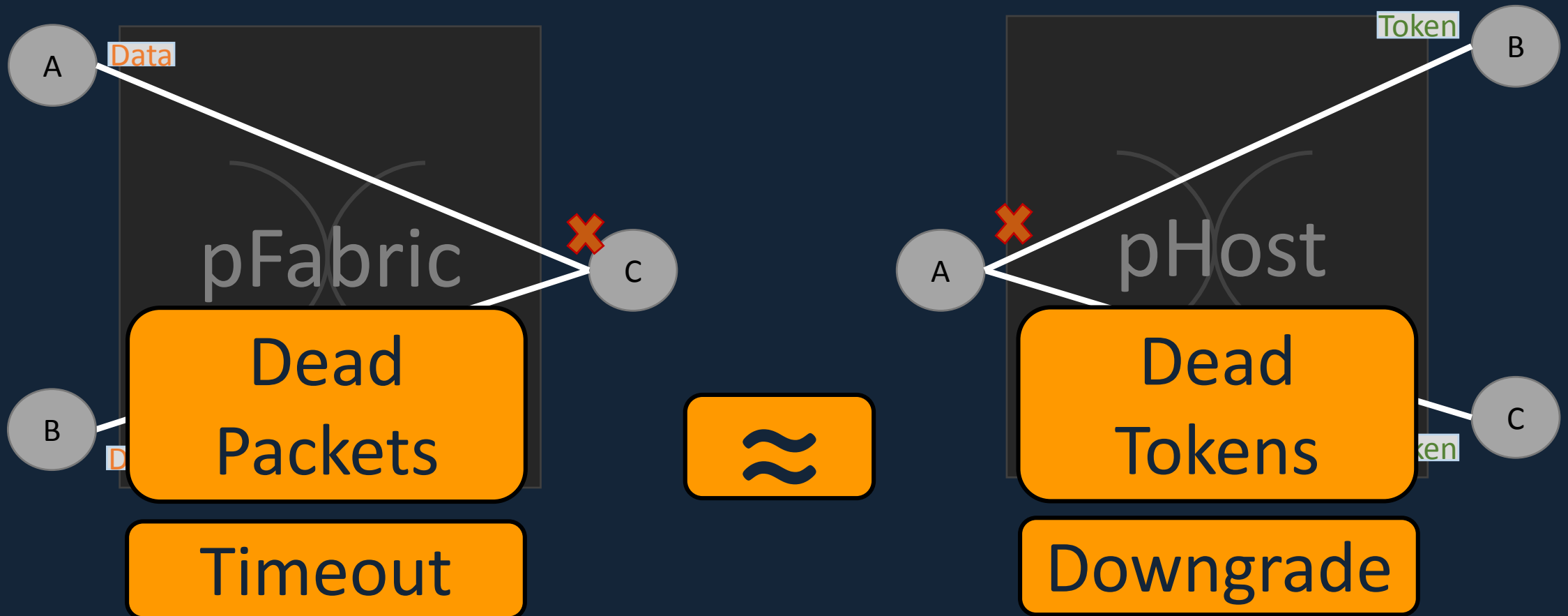
Long Flow Slowdown

Overall performance dominated by  
short flow performance

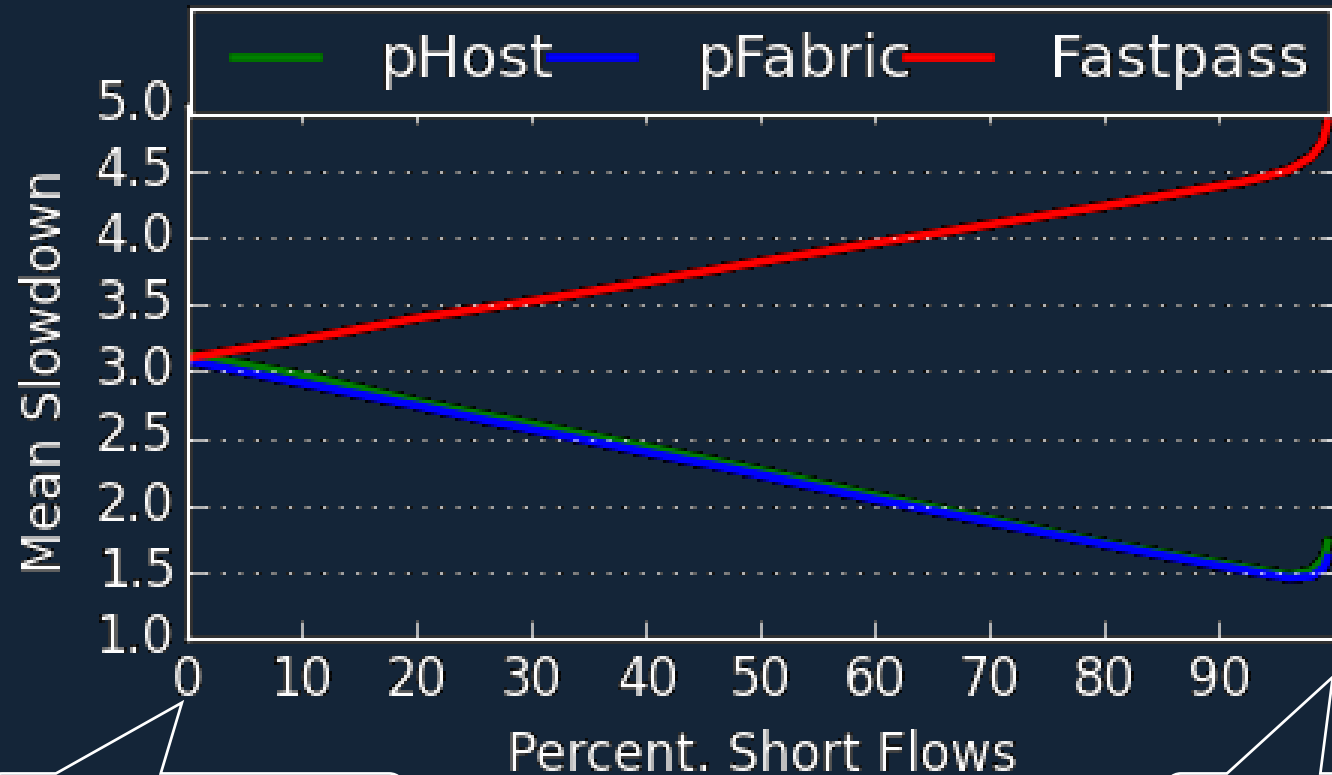


# Digging Deeper

How does pHost match pFabric's performance?



# Is pHost robust?

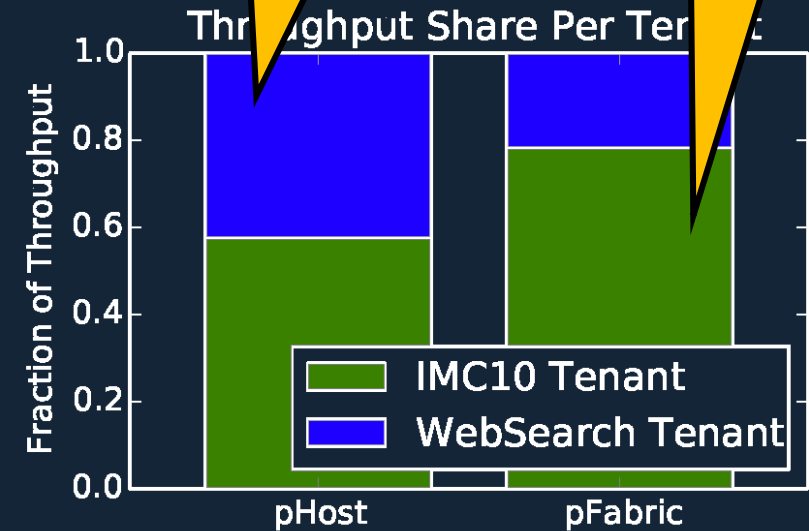
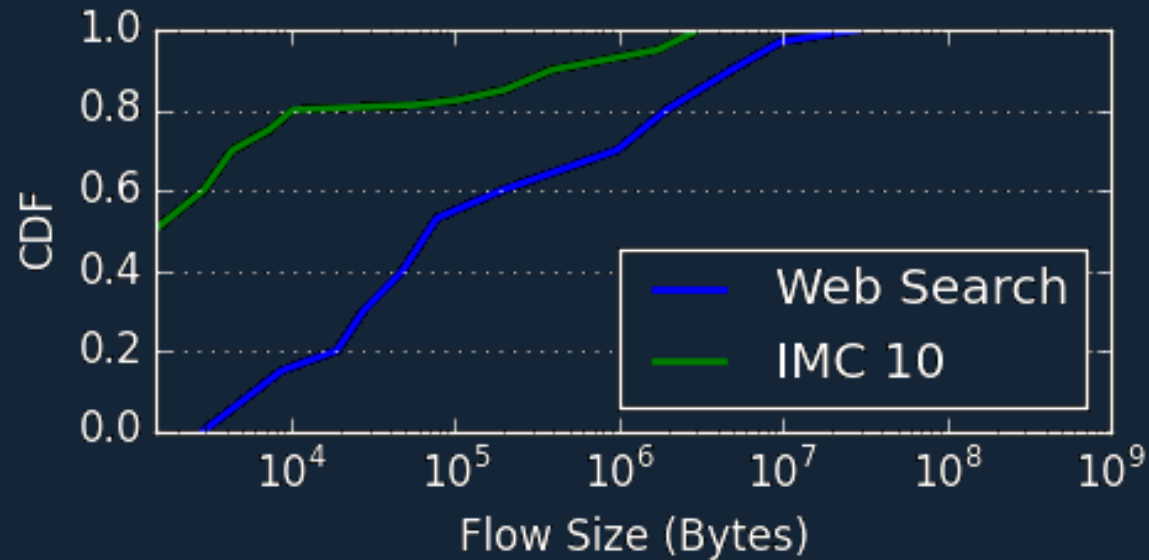


All Long (700 pkts)

All Short (3 pkts)

pHost is robust against a wide range of workloads

# Is pHost flexible?



pHost is more flexible

# More Results

Metrics	Traffic Patterns	Sensitivity
<ul style="list-style-type: none"><li>• Normalized FCT</li><li>• Tail Slowdown</li><li>• Deadline</li><li>• Throughput</li></ul>	<ul style="list-style-type: none"><li>• Random</li><li>• Incast</li><li>• Permutation</li></ul>	<ul style="list-style-type: none"><li>• Stability Analysis</li><li>• Varying Traffic Load</li><li>• Varying Switch Buffer Size</li><li>• Parameter Sensitivity</li></ul>

Please read our paper for details

# pHost

- No specialized hardware
- No complex computation inside the fabric
- No centralized scheduler
- No explicit network feedback

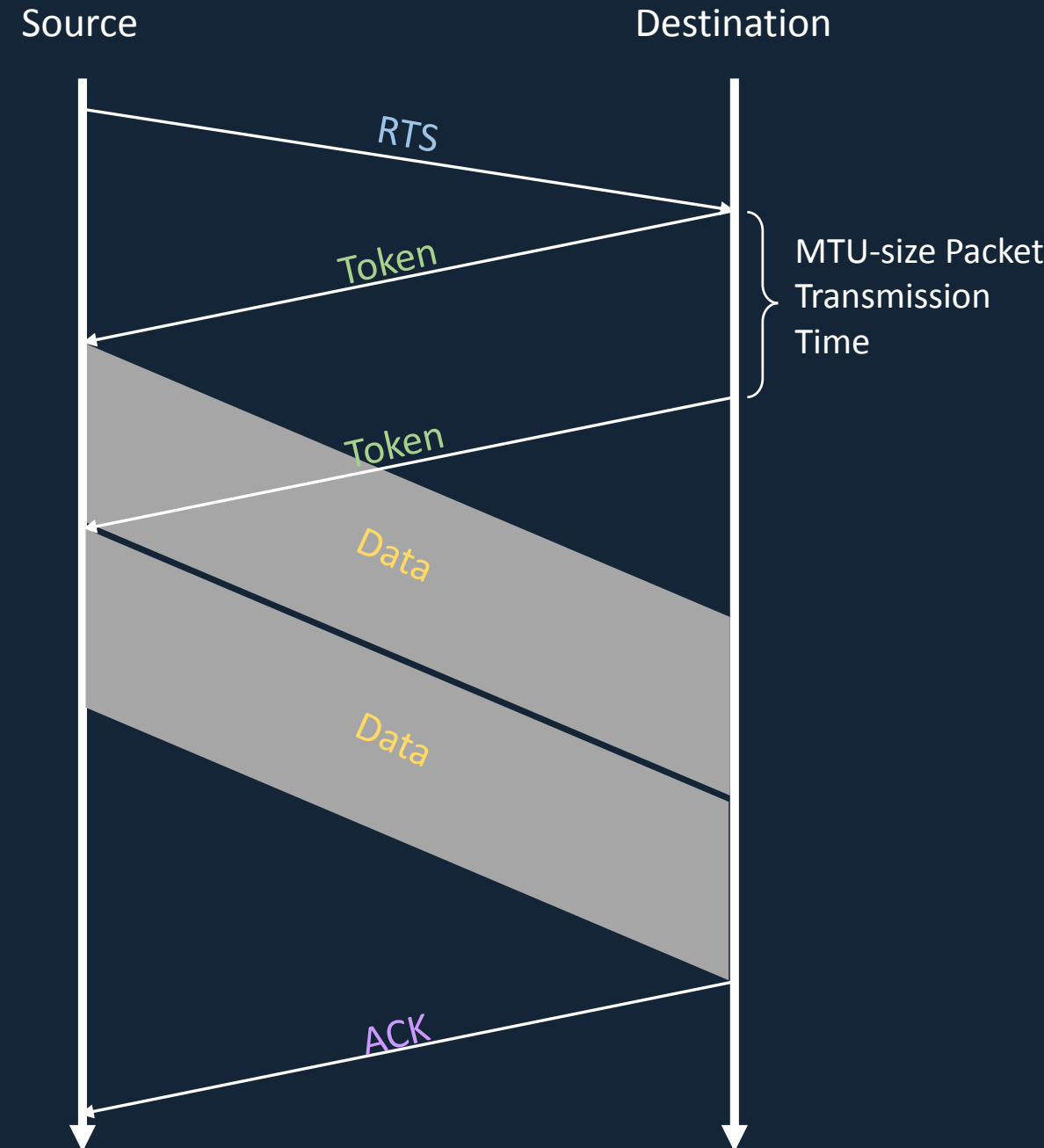
With ~~surprisingly~~ near optimal performance





# pHost Basics

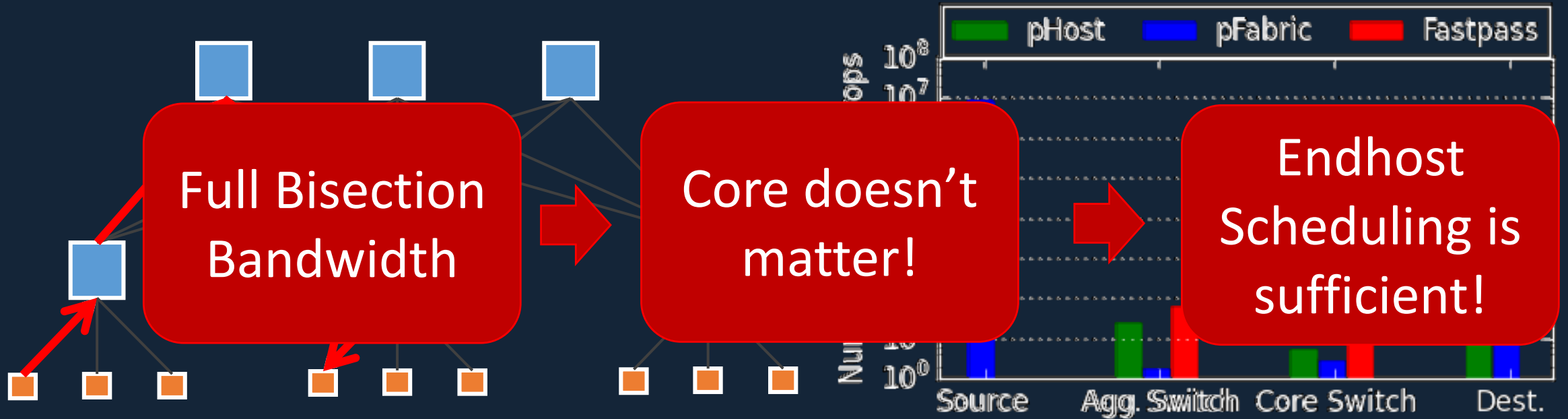
- Source sends an *RTS* on flow arrival
- The destination issues a *Token* to the source every *MTU-size packet transmission time*
- The source can consume that token by sending a *Data* packet to the destination
- A token expires if the source does not use it quickly
- Finally, when the destination receives all the data packets, it sends an *ACK* to the source



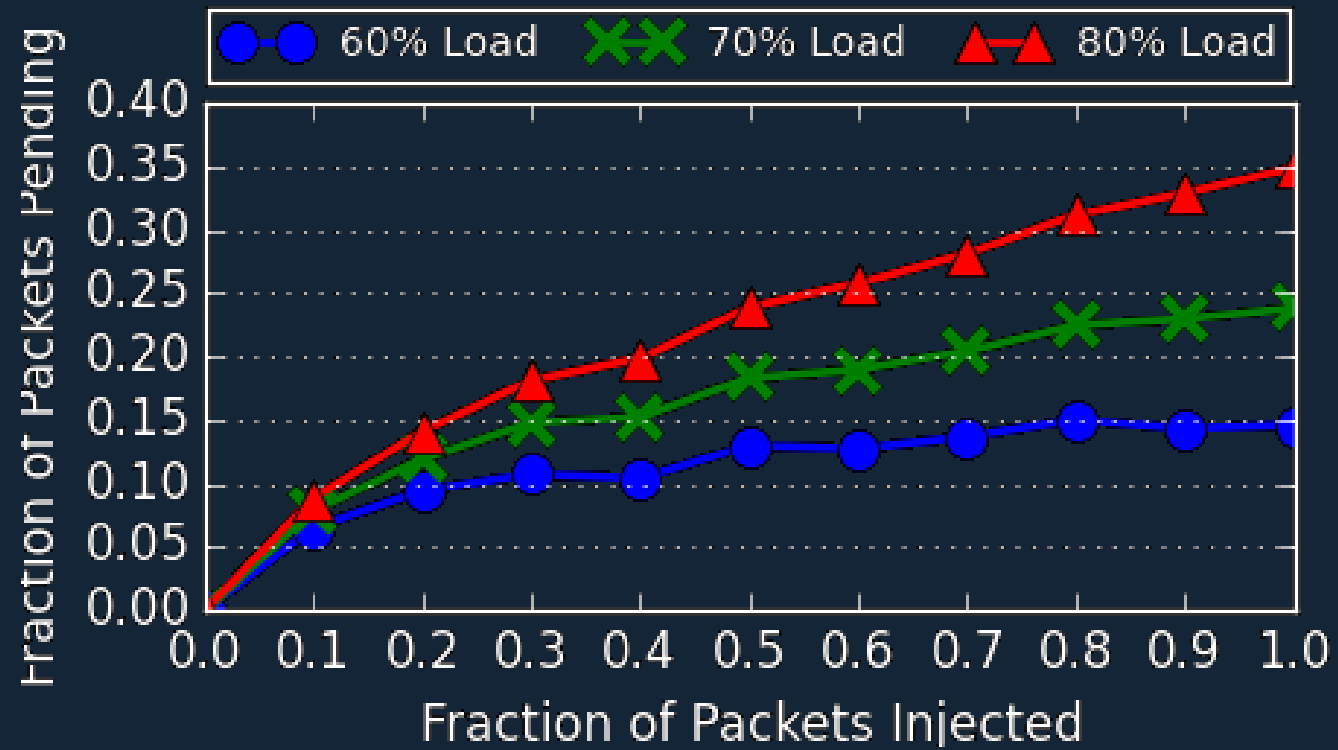


# Digging Deeper:

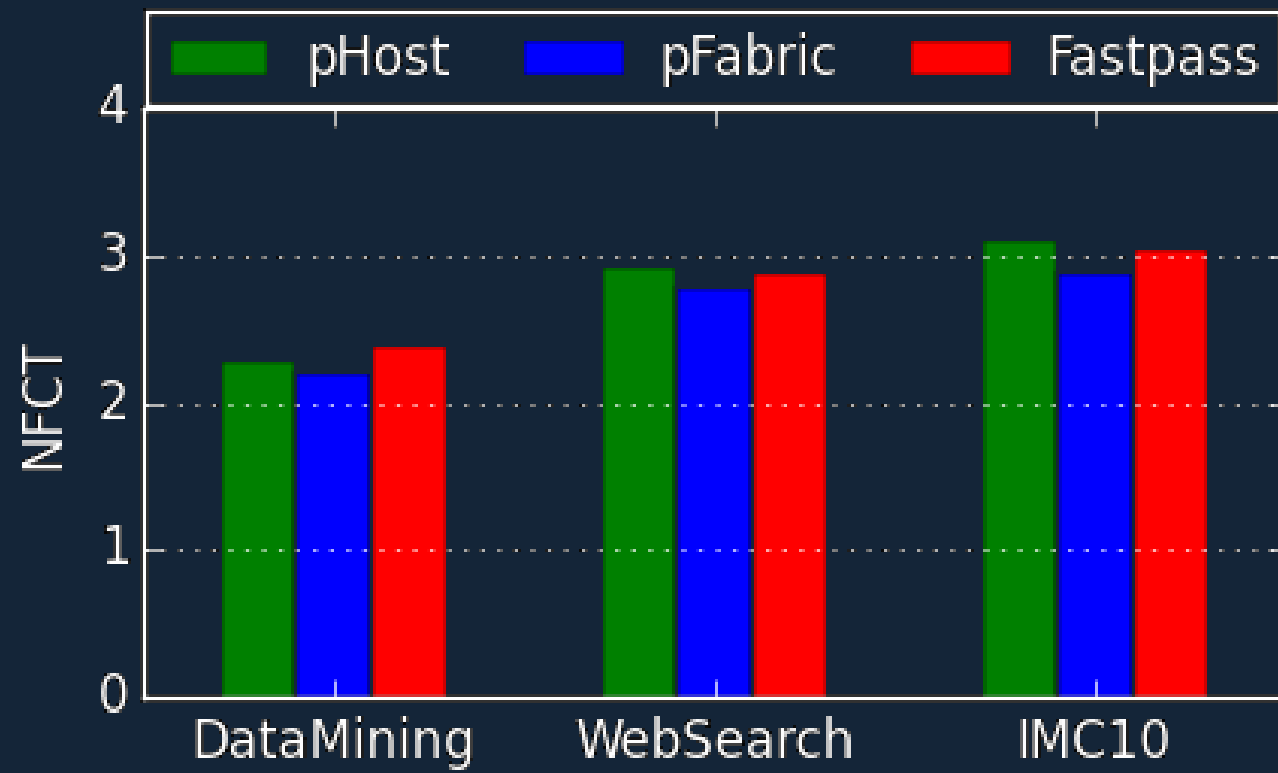
How pHost matches pFabric performance?



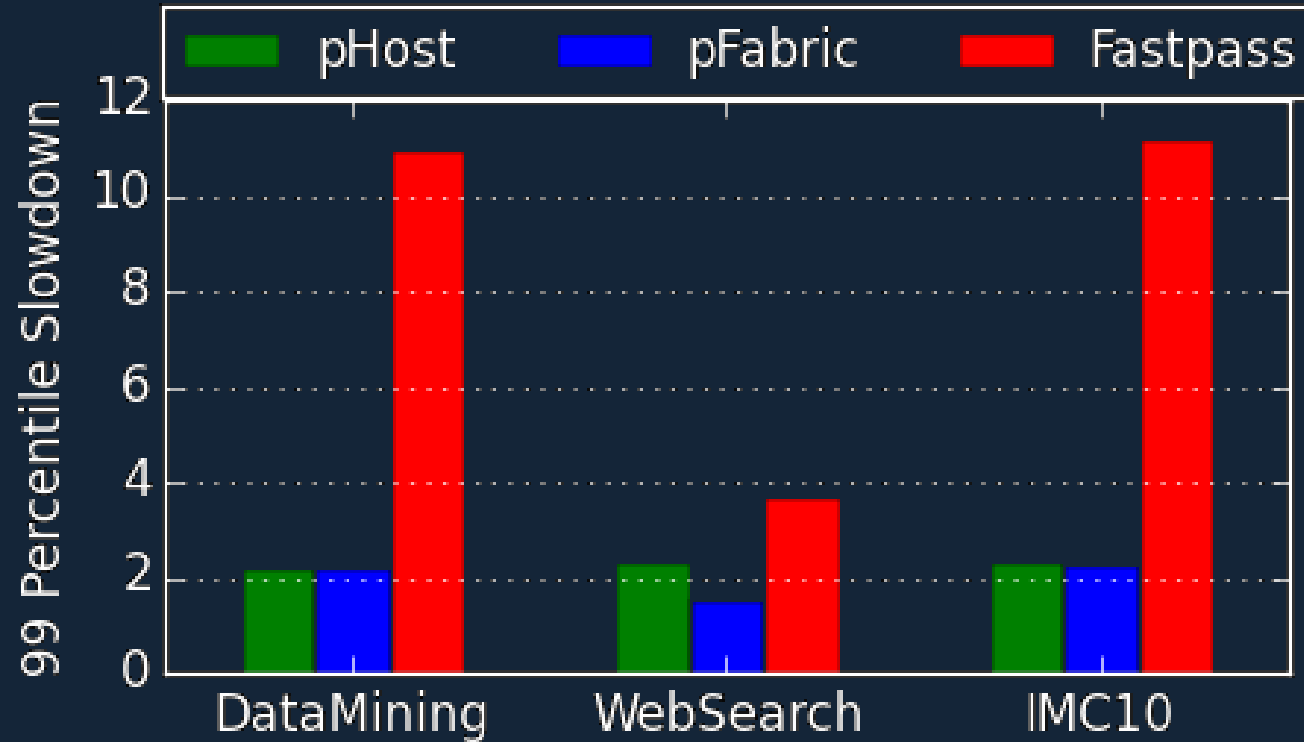
# Stability



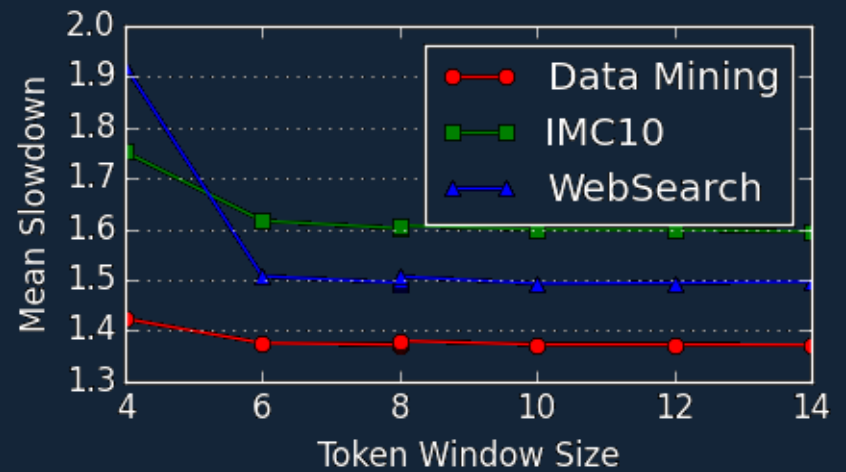
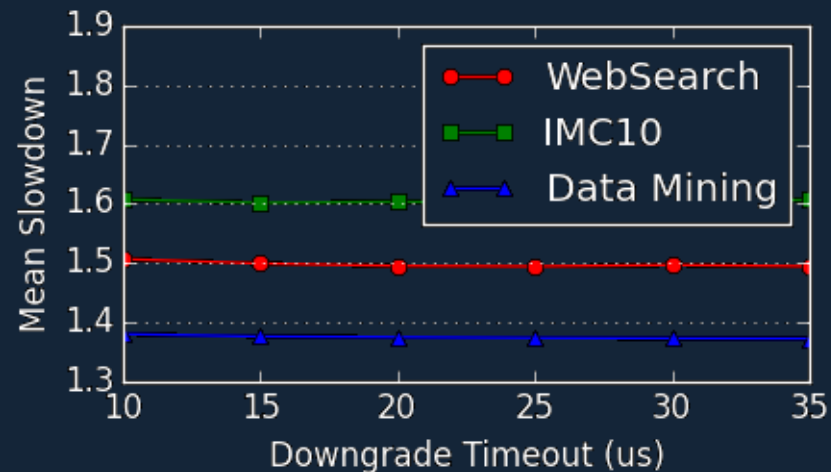
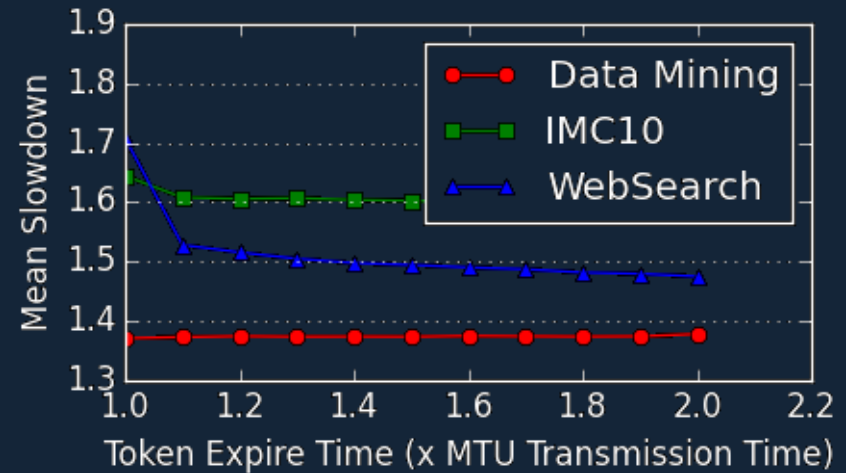
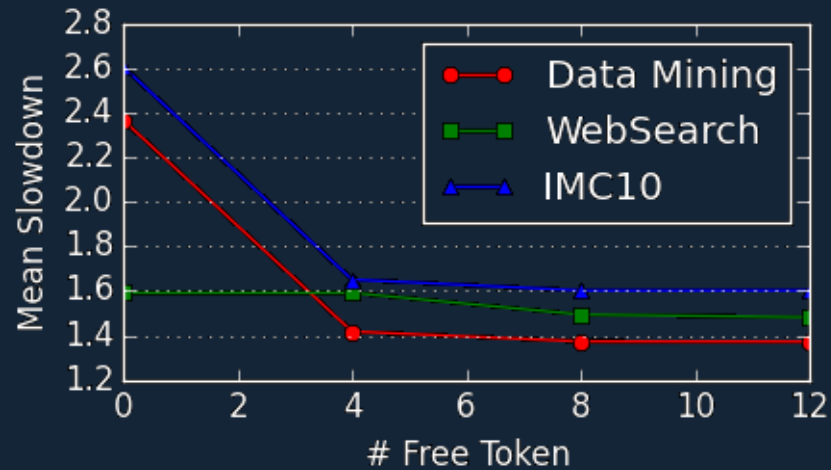
# NFCT



# 99%-ile slowdown for short flows

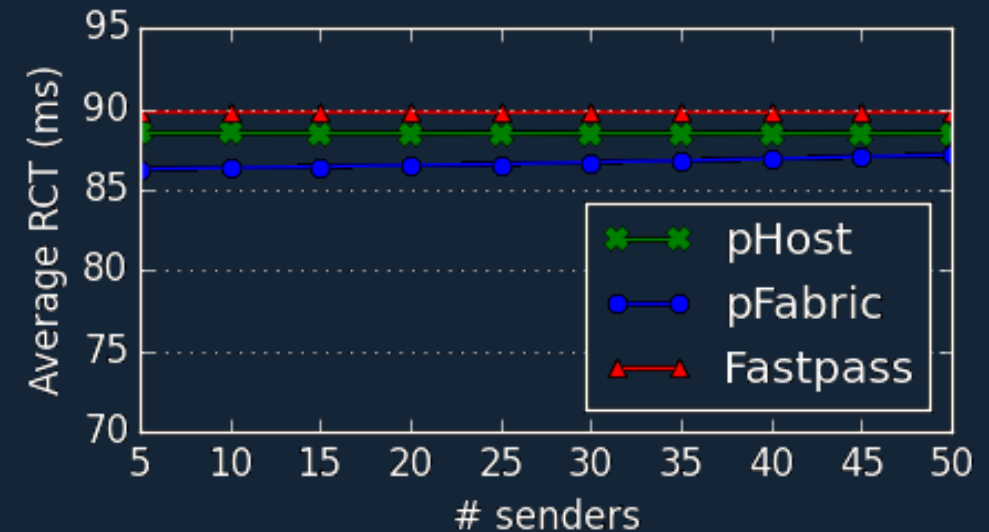
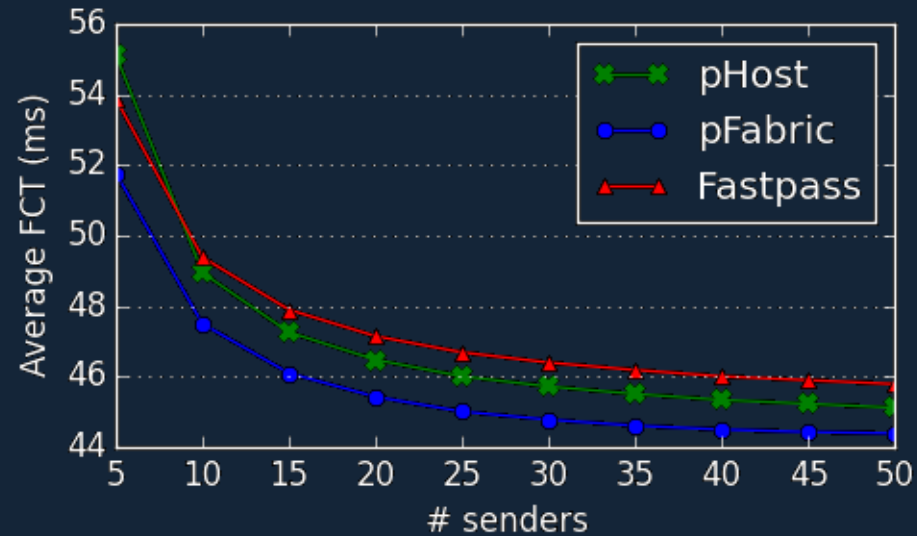


# Parameter Sensitivity Analysis



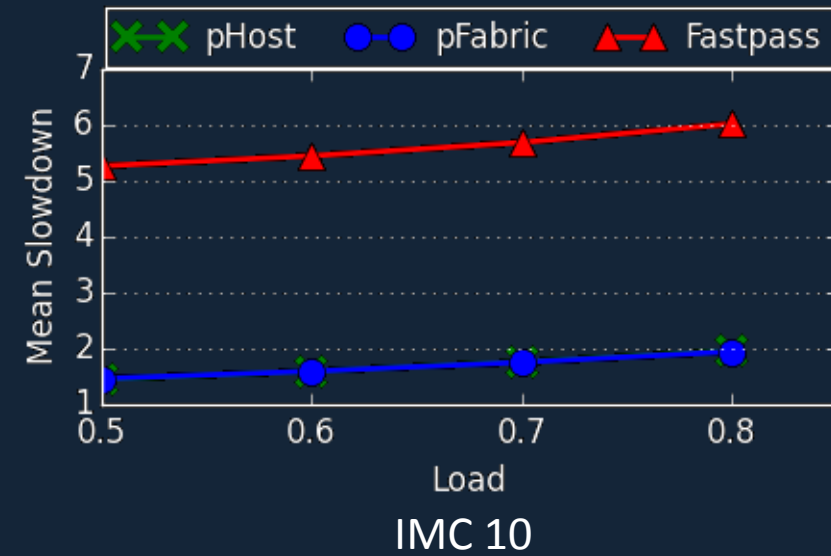
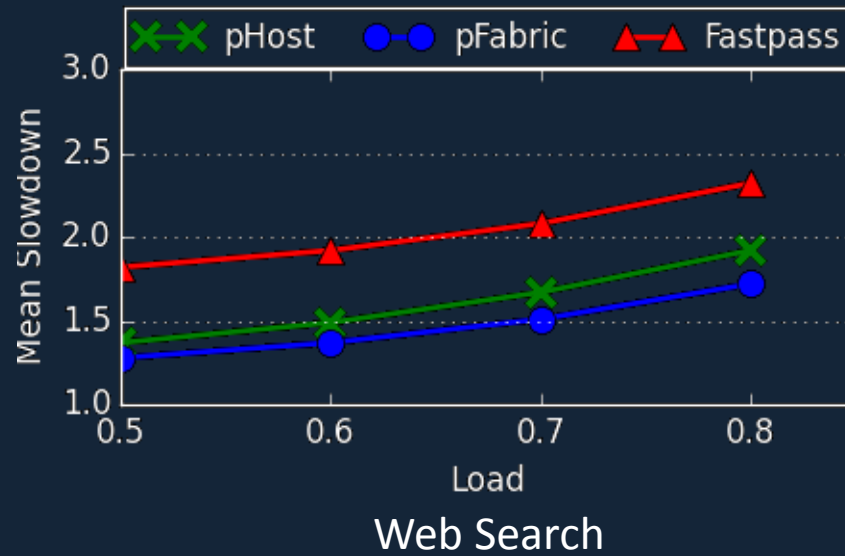
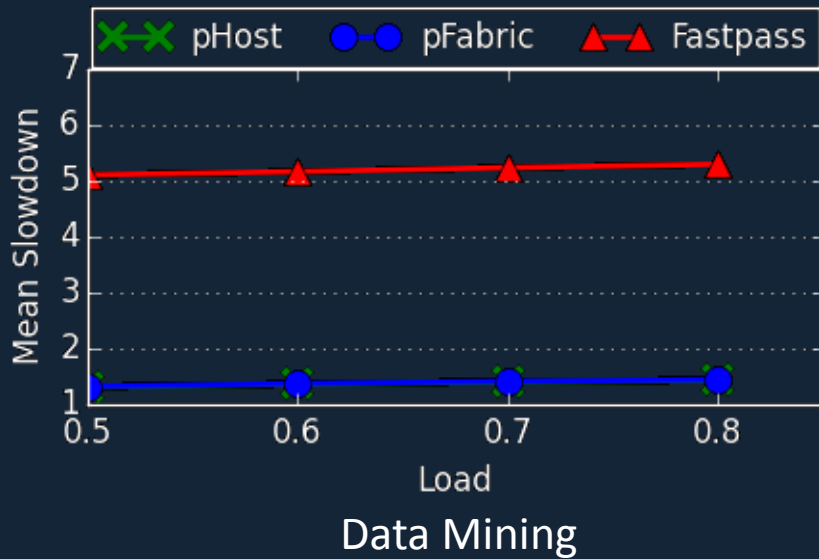
# Incast Traffic

- Each Request: N sources send 100MB data to one receiver
- 10000 requests



The performance difference between all three protocols are less than 4% and 7% in terms of FCT and RCT

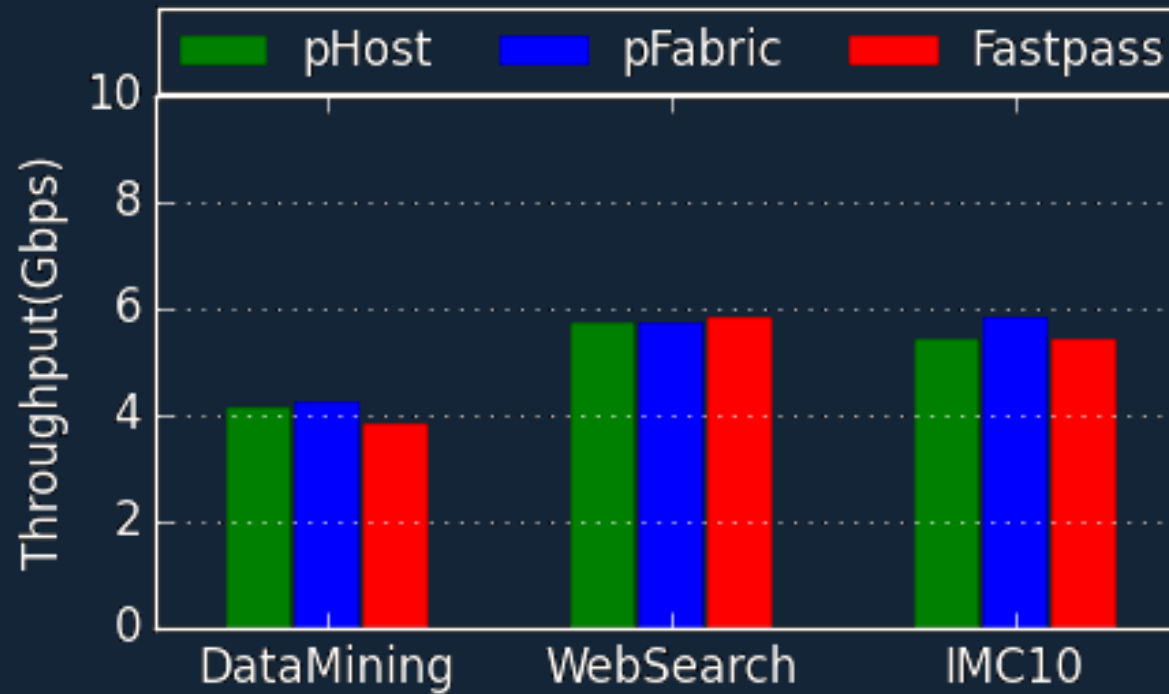
# Varying load



There is no relative performance change with varying load

# Throughput

- Load: Average byte rate at sender
- Throughput: Average byte rate at receiver

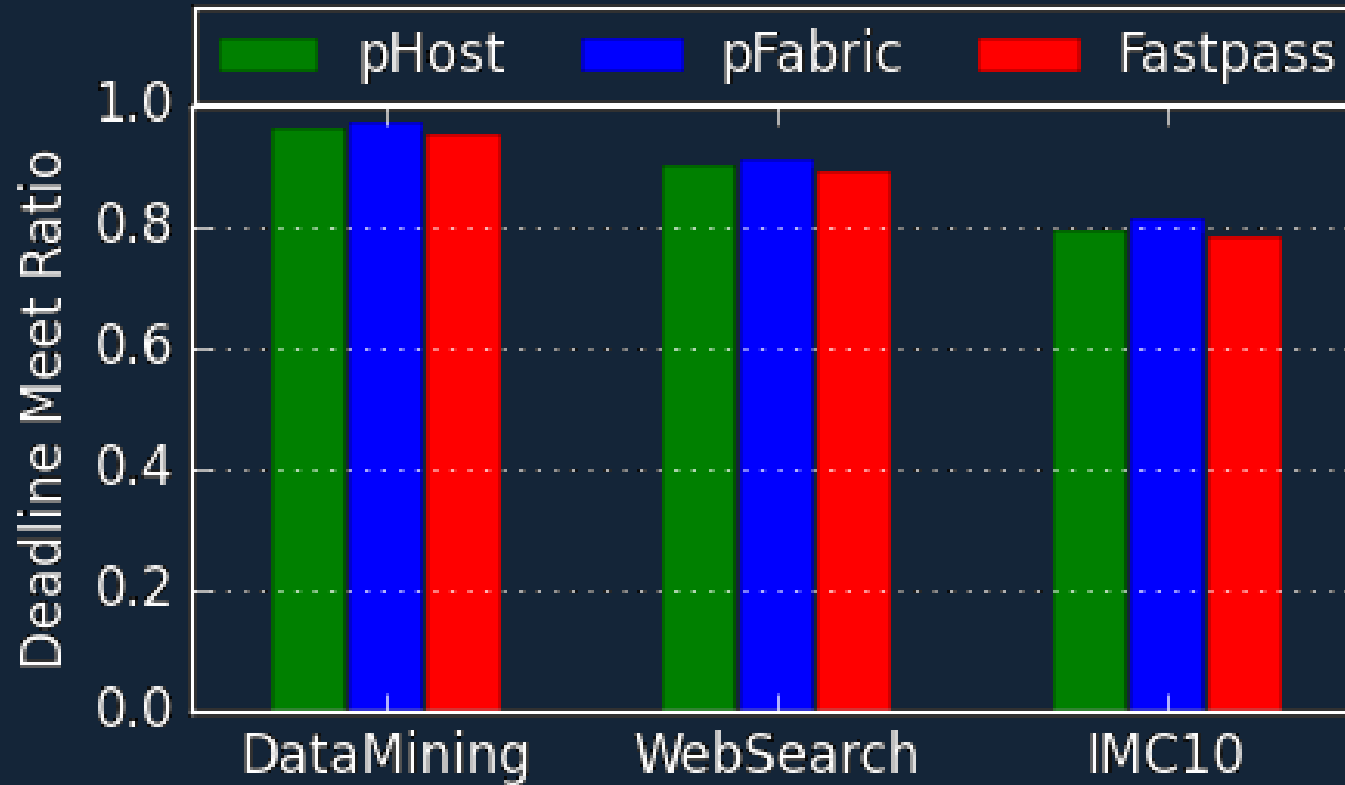


Similar performance for all flows (even for Fastpass)  
Because throughput is more related to long flow performance



# Deadline

- Exponential deadline distribution with mean 1000us
- At least 1.25x optimal FCT



All protocols have similar performance