# Complementing TCP with Homa
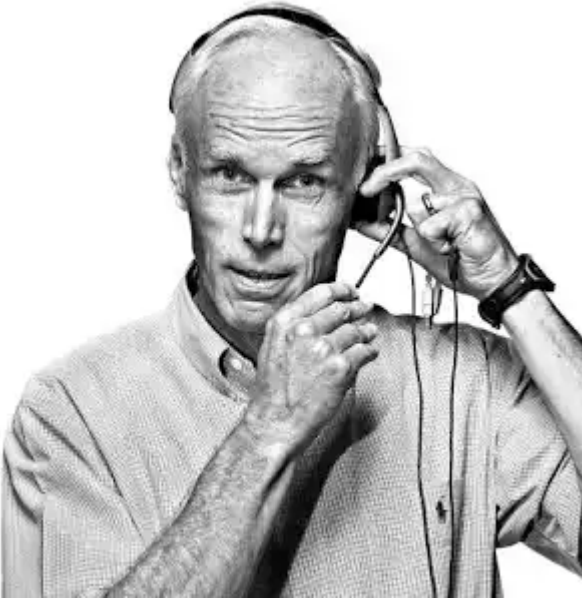
Stanford's Reliable, Rapid Request-Response Protocol

Endric Schubert, Ph.D. & Ulrich Langenbach

# Presentation Outline
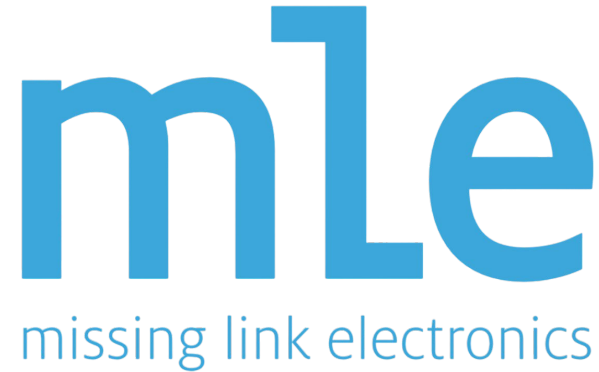
1. The Homa Protocol from John Ousterhout @Stanford University

2. Can Homa coexist peacefully with TCP/IP?

3. RRRRP - Homa, FPGA Accelerated

4. Call for Collaboration

# Acknowledgements



John Ousterhout, Stanford University

Team MLE in Berlin and Neu-Ulm, Germany

Björn Petersen, Institute for Micro Electronics, Ulm University, Germany

# Part 1

The Homa Protocol

# Homa - Started by John Ousterhout et al.

## It's Time to Replace TCP in the Datacenter

## A Linux Kernel Implementation of the Homa Transport Protocol

# Homa Reduces Tail Latencies in Loaded Networks

- Experimental results
- 25 GigE Network
- Compares Linux kernel space implementation of
  - TCP/IPv4
  - Homa/IPv4

- X-axis is distribution of message mengths in workload
- Y-axis is Slowdown RTT_loaded / RTT_unloaded



2 ms down to 100 µs? Nice!

(Gbps)

TCP P99
TCP P50
Homa P99
Homa P50

19x @ P99

7.5x @ P50

Slowdown

10
100
10
1

60    376    561    976    49.4K    1.0M

Message Length

# Courtesy of John Ousterhout, Stanford University

## 1. Homa is Message-Based

- **Dispatchable units are explicit in the protocol**

- **Enables efficient load balancing**
  - Multiple threads can safely read from a single socket
  - Future NICs can dispatch messages directly to threads

- **Enables run-to-completion (e.g. SRPT)**

## 2. Homa is Connectionless

- **Fundamental unit is a remote procedure call (RPC)**
  - Request message
  - Response message
  - RPCs are independent

- **No long-lived connection state**
  - (But there is long-lived per-peer state: ~200 bytes)

- **No connection setup overhead**
  - Use one socket to communicate with many peers

- **Homa ensures end-to-end RPC reliability**
  - No need for application-level timers

# Courtesy of John Ousterhout, Stanford University

## 3. Homa: Receiver-Driven Congestion Control



- **Receiver can delay grants to:**
  - Reduce congestion in TOR
  - Prioritize shorter messages

- **Message sizes allow receivers to predict the future:**
  - Faster, more accurate response to congestion

# Homa Uses Priority Queues

- **Modern switches: 8–16 priority queues per egress port**

- **Homa receivers select priorities for SRPT:**
  - Favor shorter messages

- **Achieve both high throughput and low latency**
  - Need buffering to maintain throughput (e.g. if sender doesn't respond to grant)
  - But buffers can result in delays
  - Solution: overcommitment:
    - Grant to multiple messages
    - Different priority for each message

## Overcommitment

Short messages use high priority queues (low latency)



Buffers accumulate in low-priority queues (ensure throughput)

## 4. Homa: SRPT

- **Combination of grants, priorities**

- **Run-to-completion improves performance for every message length!**

- **Starvation risk for longest messages?**
  - Use 5-10% of bandwidth for oldest message

# 5. Homa: No Order Requirement

- **Can use packet spraying in datacenter networks**
  - Hypothesis: will eliminate core congestion
    (unless core fabric systemically overloaded)

- **Better load balancing across CPU cores**

# HomaModule - Implemented as a Linux Kernel Module



## Uses A-Priori Knowledge

- Link Rate between NIC and ToR switch

- NIC Queue Length (SRPT), i.e. "estimated time until Tx buffer is empty"

- Coexistence w/ other protocols Interaction with Tx pacer in Linux netdev NAPI

- Distance between machines Handling non-uniform RTT

- Priority Queues in Switches

# Part 2

Can Homa coexist peacefully with TCP?

# Experimental Test Setup (at Cloudlab xl170)

## HW Setup

25 GigE network
4 or 12 nodes, each

  Intel E5-2640v4

  Mellanox ConnectX-4

## SW Setup

HomaModule v2023-12-20

util/cp_vs_tcp tests in parallel

with Linux iperf

## NW Traffic

util/cp_node for the
Homa vs TCP tests

      RTT and Slowdown

plus additive TCP "background noise" via iperf (any-to-any)

# Experimental Test Setup

| | % | (TCP) Stream target load *(using iperf)* | | | | | |
|---|---|---|---|---|---|---|---|
| | | 10 | 30 | 40 | 50 | 60 | 70 |
| Homa vs TCP RPC target load *(using cp_node)* | 10 | W2, W3, W4 | W2, W3, W4, | W2, W3, W4, | W2, W | W2, W3, W4, | W2, W3, W4, |
| | 30 | | | W4, W5 | | | |
| | 40 | W3, W4, W5 | W3, W4, W5 | W3, W4, W5 | W3, W4, W5 | W3, W4, W5 | N/A |
| | 50 | W3, W4, W5 | W3, W4, W5 | W3, W4, W5 | W3, W4, W5 | N/A | N/A |
| | 60 | W | W5 | W5 | N/A | N/A | N/A |
| | 70 | W | | | N/A | N/A | N/A |

*How much does TCP RPC disturb the TCP iperf traffic?*

*How does Homa RPC behave under increasing TCP "background noise" vs. TCP RPC?*

*How much does Homa RPC disturb the TCP iperf traffic?*

# Slowdown Results Workload W4



How much does TCP disturb the TCP iperf traffic?

How does Homa RPC behave under increasing TCP "background noise" vs. TCP RPC?

How much does Homa disturb the TCP iperf traffic?

# Slowdown Results Workload W4



W4 4 nodes, 4.0 Gbps (W4 Homa/cp_node: 2.0 Gbps, TCP/iperf: 2.0 Gbps)
RPC Slowdown, based on assumed minimum (15 us + 100% link throughput)

# Slowdown Results Workload W4



W4 4 nodes, 16.0 Gbps (W4 Homa/cp_node: 2.0 Gbps, TCP/iperf: 14.0 Gbps)
RPC Slowdown, based on assumed minimum (15 us + 100% link throughput)

# Slowdown Results Workload W4

# Slowdown Results Workload W4



W4 4 nodes, 18.0 Gbps (W4 Homa/cp_node: 14.0 Gbps, TCP/iperf: 4.0 Gbps)
RPC Slowdown, based on assumed minimum (15 us + 100% link throughput)

# RTT Results Workload 4



W4 4 nodes, 4.0 Gbps (W4 Homa/cp_node: 2.0 Gbps, TCP/iperf: 2.0 Gbps)
RPC RTT

W4 4 nodes, 16.0 Gbps (W4 Homa/cp_node: 2.0 Gbps, TCP/iperf: 14.0 Gbps)
RPC RTT

# RTT Results Workload 4



W4 4 nodes, 16.0 Gbps (W4 Homa/cp_node: 14.0 Gbps, TCP/iperf: 2.0 Gbps)
RPC RTT

# RTT Results Workload 4



W4 4 nodes, 18.0 Gbps (W4 Homa/cp_node: 14.0 Gbps, TCP/iperf: 4.0 Gbps)
RPC RTT

# Part 3

RRRRP - Homa, FPGA Accelerated

# HomaModule - A Linux Kernel Implementation of Homa

A layer just above IP, parallel to TCP and UDP

Uses GRO (Generic Receive Offloading)



**Figure 2:** Structure of Homa/Linux. Homa components are shown in blue; existing Linux kernel modules are in yellow. Gray areas represent different cores. Only the primary sending and receiving paths are shown; other Homa elements such as the pacer thread and timer thread also transmit packets.

# RRRRP - Two Acceleration Approaches

## Offload Engines

- Runs as SW on CPU
- Offload engines in FPGA
- Uses "golden" reference implementation (i.e. HomaModule)
- Low to medium eng. work
- Instant benefits

## Full Acceleration

- Entire stack runs in FPGA
- No SW on CPU
- Major eng. work (mostly in testing correctness)
- Integrated with MLE NPAP, a TCP/UDP/IP FPGA Stack
- Maybe later

# RRRRP - Offload Approach with Corundum.io

Open source FPGA NIC ported to many FPGA cards (http://www.corundum.io)
Good PCIe subsystem which supports many Rx/Tx FPGA queues.

# Offload Engines

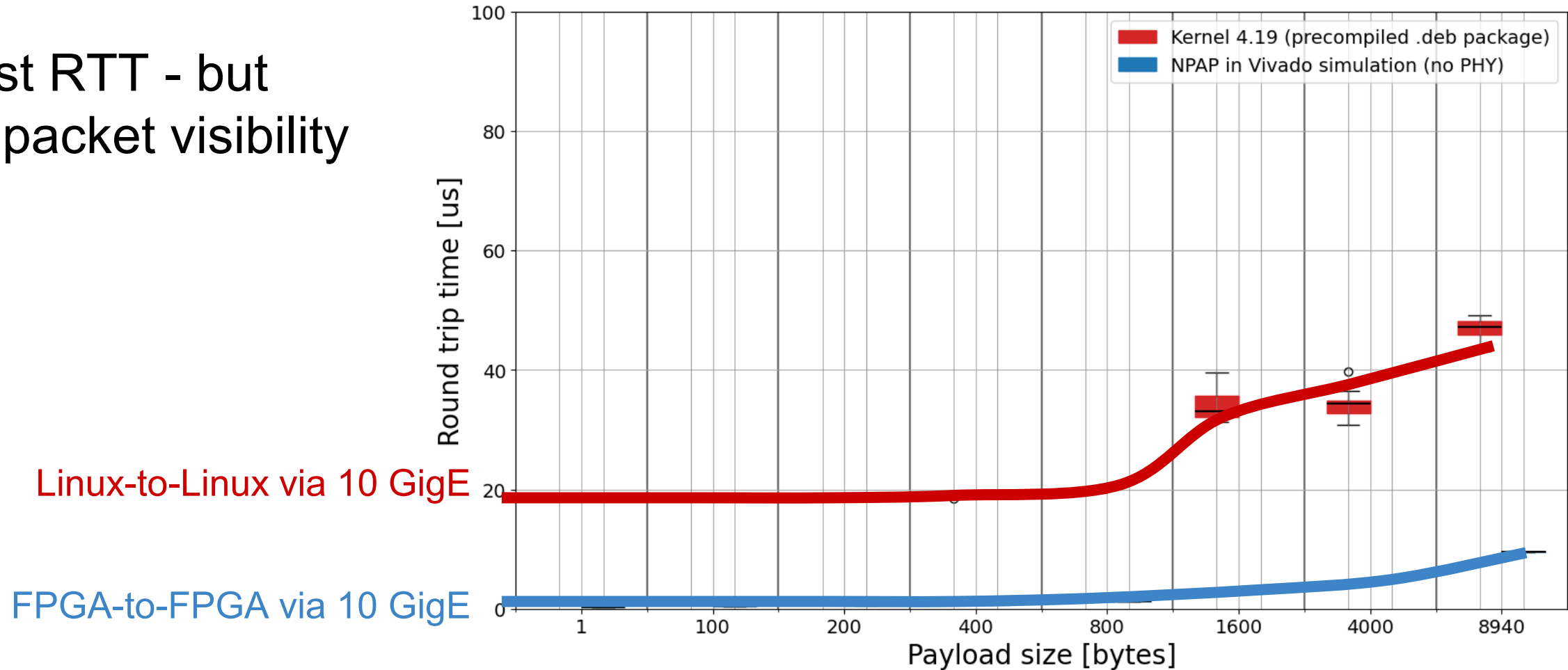| TCP (MLX 5) | HomaModule | RRRRP |
|---|---|---|
| Checksum Offload | N/A | N/A |
| Receive Side Scaling | Receive Flow Steering | Receive Side Scaling |
| Large Segmentation Offload | Generic Segmentation Offload | Large Segmentation Offload |
| Large Receive Offload | Generic Receive Offload | Large Receive Offload |

# HomaModule vs RRRRP - Slowdown for Workload W4

- Similar results for other workloads

- Current implementation runs on 10 GigE NIC
- Next: 25/50/100 GigE



W4 2 nodes, 9.0 Gbps
Cumulative % of Messages

# Motivation for Homa Full Acceleration

Lowest RTT - but
lacks packet visibility



Linux-to-Linux via 10 GigE

FPGA-to-FPGA via 10 GigE

Legend:
- Kernel 4.19 (precompiled .deb package)
- NPAP in Vivado simulation (no PHY)

Y-axis: Round trip time [us]
X-axis: Payload size [bytes]

# Part 4

## Call for Collaboration

⇒ Run on larger FPGA cluster
to check scalability

⇒ Try other workloads

⇒ Look at applications
Networked storage systems?
AI clusters?

# Please take a moment to rate this session.

Your feedback is important to us.

SDC 24