TAS: TCP Acceleration as an OS Service

Antoine Kaufmann¹, **Tim Stamler²**, Simon Peter², Naveen Kr. Sharma³, Arvind Krishnamurthy³, Thomas Anderson³

MPI-SWS¹ The University of Texas at Austin² University of Washington³

RPCs are Essential in the Datacenter

Remote procedure calls (RPCs) are a common building block for datacenter applications

Scenario: An efficient key-value store in a datacenter

- 1. Low tail latency is crucial
- 2. Thousands of connections per machine
- 3. Both the application writer and datacenter operator want the full feature set of TCP
 - a) Developers want the convenience of *sockets* and *in-order delivery*
 - b) Operators want *flexibility* and strong *policy enforcement*

You might want to simply go with Linux...

Linux provides the features we want *sockets in-order delivery* But at what cost?

flexibility

policy enforcement

You might want to simply go with Linux...

Linux provides the features we want

sockets But at what cost? in-order delivery

flexibility

policy enforcement

A simple KVS model:

256B RPC request/response over Linux TCP 250 application cycles per RPC

You might want to simply go with Linux...

Linux provides the features we want

sockets But at what cost?

A simple KVS model:

256B RPC request/response over Linux TCP 250 application cycles per RPC

flexibility

8,300 Total CPU Cycles per RPC

App Processing: 3%

Kernel Processing: 97%

policy enforcement

We're only doing a small amount of useful computation!

in-order delivery



Why not kernel-bypass?

NIC interface is optimized, bottlenecks are in OS
Arrakis (OSDI '14), mTCP (NSDI '14), Stackmap (ATC '16)
Do network processing in userspace
Expose the NIC interface to the application

Hardware I/O virtualization

Why not kernel-bypass?

NIC interface is optimized, bottlenecks are in OS *Arrakis* (OSDI '14), *mTCP* (NSDI '14), *Stackmap* (ATC '16) Do network processing in userspace Expose the NIC interface to the application Hardware I/O virtualization

Avoid OS overheads, can specialize stack

Why not kernel-bypass?

NIC interface is optimized, bottlenecks are in OS Arrakis (OSDI '14), mTCP (NSDI '14), Stackmap (ATC '16) Do network processing in userspace

Expose the NIC interface to the application

Hardware I/O virtualization

- Avoid OS overheads, can specialize stack
- Operators have to trust application code
- E Little flexibility for operators to change or update network stack

Why not RDMA?

Remote Direct Memory Access:

Interface: **one-sided** and **two-sided** operations in NIC hardware

RPCs and sockets implemented on top of basic RDMA primitives

Why not RDMA?

Remote Direct Memory Access:

Interface: one-sided and two-sided operations in NIC hardware

RPCs and sockets implemented on top of basic RDMA primitives

Minimize or bypass CPU overhead

Why not RDMA?

Remote Direct Memory Access:

Interface: **one-sided** and **two-sided** operations in NIC hardware

RPCs and sockets implemented on top of basic RDMA primitives

- Minimize or bypass CPU overhead
- **Ex** Lose software procotol flexibility
- Bad fit for many-to-many RPCs
- **RDMA** congestion control (DCQCN) doesn't work well at scale

TAS: TCP Acceleration as an OS Service

An open source, drop-in, highly efficient RPC acceleration service

No additional NIC hardware required Compatible with all applications that already use sockets Operates as a userspace OS service using dedicated cores for packet processing Leverages the benefits and flexibility of kernel bypass with better protection

TAS accelerates TCP processing for RPCs while providing all the desired featuresSocketsIn-order deliveryFlexibilityPolicy enforcement



How does TAS fix it?

System call and cache pollution overheads

Dedicate cores for network stack

Complicated data path



Poor cache efficiency, unscalable



Minimize and localize connection state



Dividing Functionality

Linux Kernel TCP Stack

• Open/close connections

Per packet:

- Socket API, locking
- IP routing, ARP
- Firewalling, traffic shaping
- Generate data segments
- Congestion control
- Flow control
- Process & send ACKs
- Re-transmission timeouts

• Socket API, locking

Fast Path

Per packet:

- Generate data segments
- Process & send ACKs
- Flow control
- Apply rate-limit



Slow Path

Per connection:

- Open/close connections
- IP routing, ARP
- Firewalling, traffic shaping
- Compute rate
- Re-transmission timeouts



• Socket API, locking

Data packet payloads

Fast Path

Per packet:

- Generate data segments
- Process & send ACKs
- Flow control
- Apply rate-limit



Slow Path

Per connection:

- Open/close connections
- IP routing, ARP
- Firewalling, traffic shaping
- Compute rate
- Re-transmission timeouts



Congestion statistics Retransmissions Control packets

• Socket API, locking

Data packet payloads

Fast Path

Per packet:

- Generate data segments
- Process & send ACKs
- Flow control
- Apply rate-limit

Minimal Connection State



Slow Path

Per connection:

- Open/close connections
- IP routing, ARP
- Firewalling, traffic shaping
- Compute rate
- Re-transmission timeouts



Congestion statistics Retransmissions Control packets



• Socket API, locking

Data packets Payload buffers

Fast Path

Per packet:

- Generate data segments
- Process & send ACKs
- Flow control
- Apply rate-limit

Minimal Connection State



Slow Path

Per connection:

- Open/close connections
- IP routing, ARP
- Firewalling, traffic shaping
- Compute rate
- Re-transmission timeouts



Congestion statistics Retransmissions Control packets



Congestion Control

Inspired by CCP (SIGCOMM '18)



Periodically check/update connection state























Evaluation

Evaluation Questions

What is our throughput, latency, and scalability for RPCs?

Do real applications scale with # of cores and have low tail latency?

Do we distribute throughput fairly under network congestion?

(See paper for more in-depth analysis)

Systems for Comparison

We evaluate TAS against 3 other systems:

- 1. Linux
 - a) Full kernel, trusted congestion control
 - b) Sockets interface
- 2. mTCP (not in this talk, see paper)
 - a) Pure kernel bypass approach, untrusted congestion control

3. IX

- a) Replace Linux with optimized data path, run in privileged mode
- b) Uses *batching* to reduce overhead
- c) No sockets interface
- d) Requires kernel modifications

Experimental Setup

Intel Xeon Platinum 8160 CPU 24 cores @ 2.10GHz 196GB of RAM Intel XL710 40Gb Ethernet Adapter

Benchmarks:

- Single direction RPC benchmark
- RPC echo server
- A scalable key-value store
- Connection throughput fairness under congestion

Linux vs TAS on RPCs (1 App Core)

- Single direction RPC benchmark
- 32 RPCs per connection in flight

RX Pipelined RPC Throughput

• 250 cycle application workload

TX Pipelined RPC Throughput

• 64 bytes realistic small RPC



Linux vs TAS on RPCs (1 App Core)

- Single direction RPC benchmark
- 32 RPCs per connection in flight

RX Pipelined RPC Throughput

• 250 cycle application workload

TX Pipelined RPC Throughput

• 64 bytes realistic small RPC



Linux vs TAS on RPCs (1 App Core)

- Single direction RPC benchmark
- 32 RPCs per connection in flight

RX Pipelined RPC Throughput

• 250 cycle application workload

TX Pipelined RPC Throughput

• 64 bytes realistic small RPC



Connection Scalability

- 20 core RPC echo server
- 64B requests/responses
- Single RPC per connection

Key factor: minimized connection state



Key-value Store

Increasing server cores with matching load (~2000 connections per core)

IX and TAS provide ~6x speedup over Linux across all cores

TAS: 9 app cores, 7 TAS cores

IX, Linux: 16 app/stack cores



Key-value Store

Increasing server cores with matching load (~2000 connections per core)

IX and TAS provide ~6x speedup over Linux across all cores

TAS: 9 app cores, 7 TAS cores

IX, Linux: 16 app/stack cores

TAS has a 15-20% performance improvement over IX without sockets TAS: 8 app cores, 8 TAS cores



KVS Throughput

Key-value Store Latency

KVS latency measure with single application core, 15% server load



Key-value Store Latency

KVS latency measure with single application core, 15% server load



Key-value Store Latency

KVS latency measure with single application core, 15% server load



IX has 50% higher latency in the 90p case



IX has 50% higher latency in the 90p case



IX has 50% higher latency in the 90p case Latency is 20us (27%) higher in the 99.99p case



IX has 50% higher latency in the 90p case Latency is 20us (27%) higher in the 99.99p case In addition, IX has a 2.3x higher maximum latency

CDF





IX has 50% higher latency in the 90p case Latency is 20us (27%) higher in the 99.99p case In addition, IX has a 2.3x higher maximum latency



Why long IX tail? *Batching*



CDF

Fairness Under Incast

We want to see how TAS distributes throughput under congestion Incast scenario, with four 10G machines all sending to one 40G server TAS on average maintains fair throughput, while Linux is unstable



Fairness Under Incast

We want to see how TAS distributes throughput under congestion Incast scenario, with four 10G machines all sending to one 40G server TAS on average maintains fair throughput, while Linux is unstable



Fairness Under Incast

We want to see how TAS distributes throughput under congestion Incast scenario, with four 10G machines all sending to one 40G server TAS on average maintains fair throughput, while Linux is unstable



Conclusion

TAS has the convenience and features of Linux, with better performance & stability

Achieved by

- 1. Separating TCP packet processing into a fast and slow path
- 2. Minimizing connection state
- 3. Dedicating cores to the network stack

TAS is a purely software solution that is easy to deploy and operate

Try it yourself!

https://github.com/tcp-acceleration-service