



Low-Latency TCP/IP Stack for Data Center Applications

David Sidler, Zsolt István, Gustavo Alonso
Systems Group, Dept. of Computer Science, ETH Zürich

Original Architecture [1]

Scalable 10Gbps TCP/IP Stack Architecture for Reconfigurable Hardware

David Sidler, Gustavo Alonso, Michaela Bilitz, Kiron Kumar, Ken Visser
 Dpt. of Computer Science, ETH Zürich
 [dsidler, alonso]@inf.ethz.ch [mbilitz, kironk, kvisvisser]@ethz.ch

Raymond Carley
 Dept. of Electrical & Computer Engineering, Carnegie Mellon University
 rcarley@ece.cmu.edu

Abstract—TCP/IP is the predominant communication protocol in modern networks but also one of the most demanding. Consequently, TCP/IP offload is becoming increasingly popular with standard network interface cards. TCP/IP Offload Engines have also emerged for FPGAs, and are being offered by vendors such as Infineon, Framerite, HIL, PLDA and DSA Group. With the target application being high-frequency trading, these implementations focus on low latency and support a limited session count. However, many more applications beyond high-frequency trading can potentially be accelerated based on FPGAs using TCP with high session count in a scalable manner. This work is a re-architectured FPGA on-chip engine and aims to a CPU co-processor functions such as encryption, compression, checksum and more others in addition to handling the complete network stack.

This paper introduces a novel architecture for a 10Gbps line-rate TCP/IP stack for FPGAs that can scale with the number of sessions and thereby addresses these new applications. We prototyped the design as a VCSIP development board, demonstrating compatibility with existing network infrastructure, operating at full 10Gbps throughput full-duplex while supporting 10,000 sessions. Finally, the design has been described primarily using high-level synthesis, which accelerates development and improves maintainability.

1. INTRODUCTION

TCP/IP is the cornerstone of modern network communications with its support for reliable data transfer including flow control, congestion avoidance, duplicate data suppression,

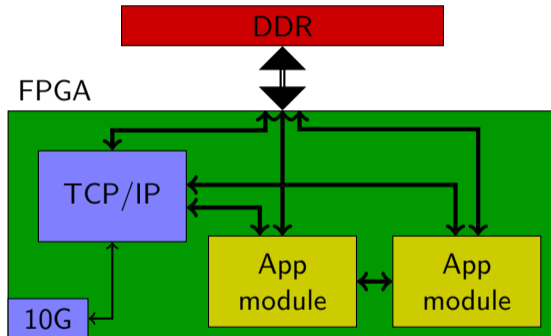
co-existence applications [2]. As a result, TCP/IP offload is increasingly integrated into network interface cards with many of the major vendors, such as Broadcom, Emulex, and Cavium offering full offload since 2011 [3]. TCP/IP Offload Engines (TOEs) have also emerged for FPGAs, offered by vendors such as Inphiq, Framerite, HIL, PLDA and DSA Group [4], [5], [6], [7]. However, these implementations target high-frequency trading, which is driven by latency requirements [8]. To overcome latency, these stacks constrain session support as we explain further in section II. Our design aims to expand the applicability of FPGAs by creating a flexible architecture that, in addition to delivering full line-rate throughput, can also scale to high session counts, an essential prerequisite to deployment in networked servers in data centers. Instead of offloading the TCP/IP hardware onto the host CPU, the solution aims to accelerate applications such as encryption, compression, checksum [9] or high-level protocol processing such as JMS [10] by pushing the TCP hardware complexity inside the FPGA. This enables session creation through reconfigurable logic for potentially thousands of sessions.

To maximize the stack's applicability, it was essential to create a flexible solution that allows to efficiently and easily adapt the design to different congestion avoidance schemes, potentially out-of-order processing, etc., while using a minimal resource footprint. To achieve this, we adopted a Co-based design flow using high-level synthesis (HLS) that simplifies the design, enables a more flexible and secure communication towards scalable

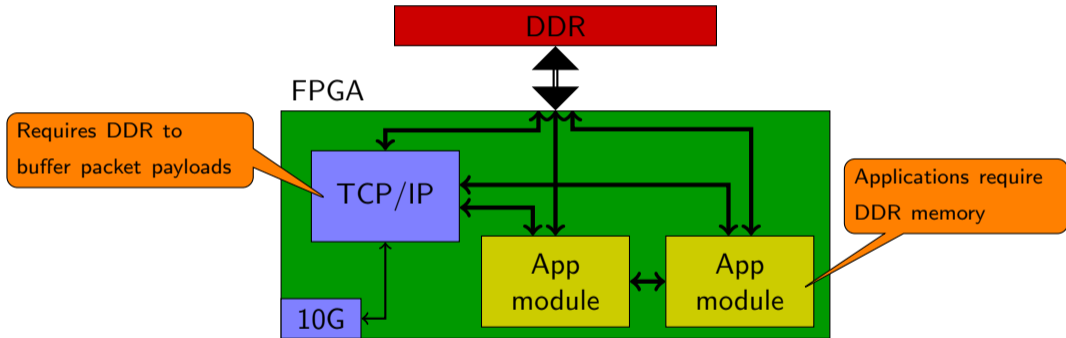
- 10 Gbps bandwidth TCP/IP stack
- Supporting thousands of concurrent connections
- Generic implementation as close to specification as possible
- Enables seamless integration of FPGA-based applications into existing networks

[1] Sidler et al., *Scalable 10 Gbps TCP/IP Stack Architecture for Reconfigurable Hardware*, FCCM'15, <http://github.com/dsidler/fpga-network-stack>

Application Integration

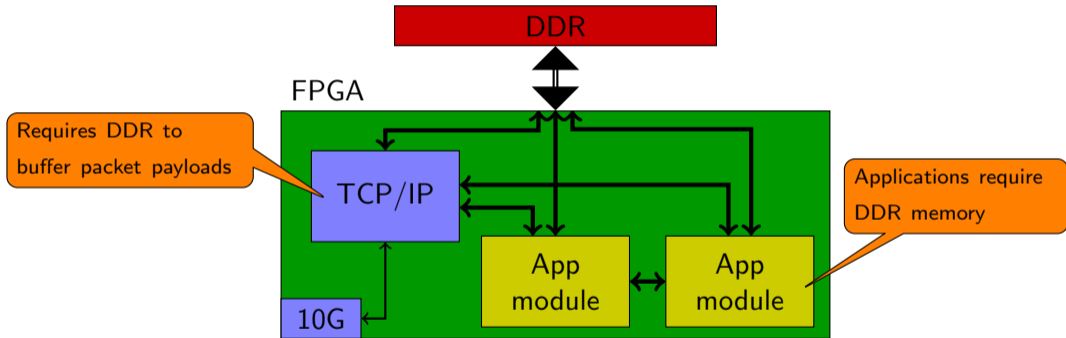


Application Integration



- Memory bandwidth is shared among multiple modules → potential bottleneck

Application Integration



- Memory bandwidth is shared among multiple modules → potential bottleneck
- Distributed systems rely on very low latency → to guarantee latency bounds to clients

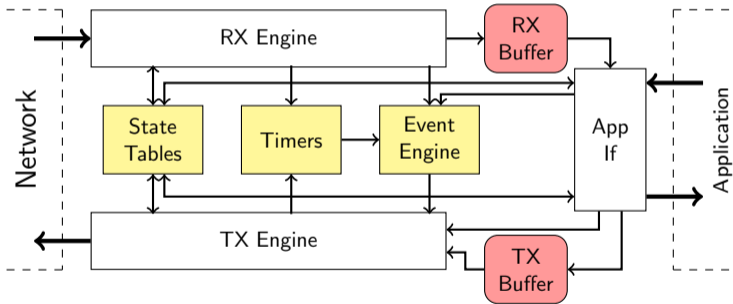
Assumptions

- Application
 - Client requests fit into an MTU (maximum transfer unit)
 - Synchronous clients
 - Application logic consumes data at line-rate

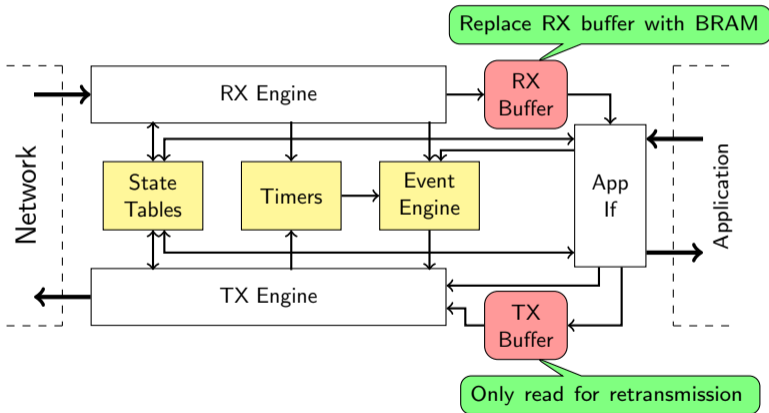
Assumptions

- Application
 - Client requests fit into an MTU (maximum transfer unit)
 - Synchronous clients
 - Application logic consumes data at line-rate
- Data center network
 - High reliability and structured topology
 - Data loss less common → fewer retransmission
 - Packets are rarely reordered

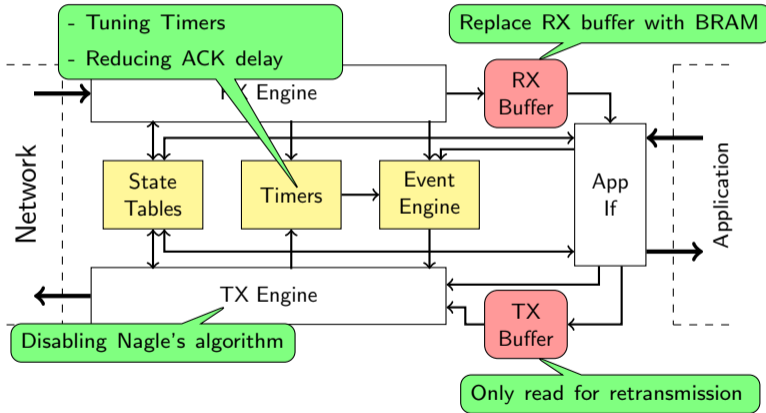
Optimizations for Data Center Applications



Optimizations for Data Center Applications

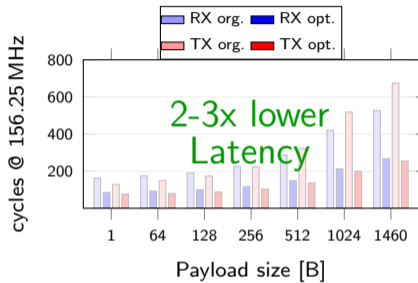


Optimizations for Data Center Applications

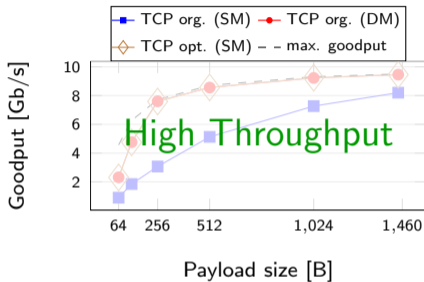
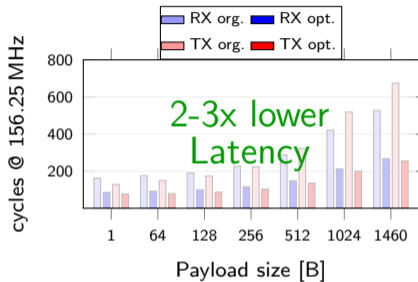


Results

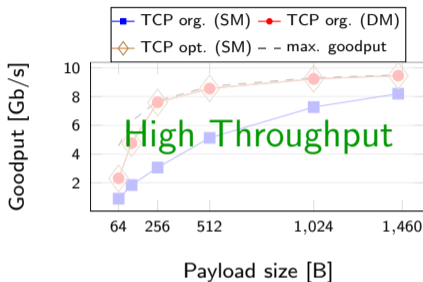
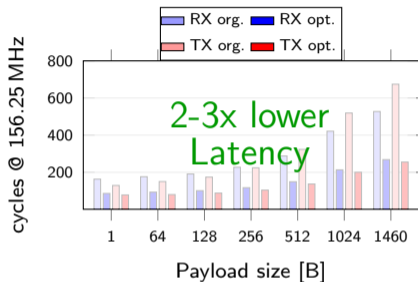
Results



Results

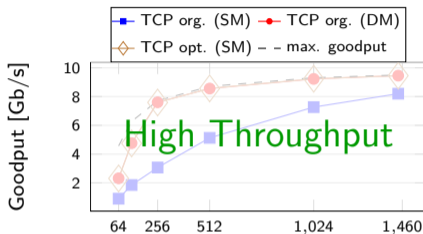
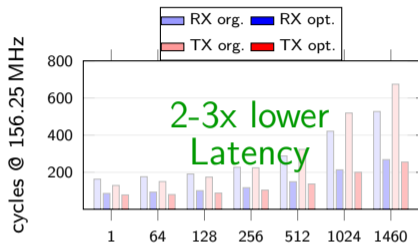


Results



	Mem. allocated	Mem. bandwidth
TCP org.	1,300 MB	40 Gbps
TCP opt.	650 MB	10 Gbps
Diff	-50%	-75%

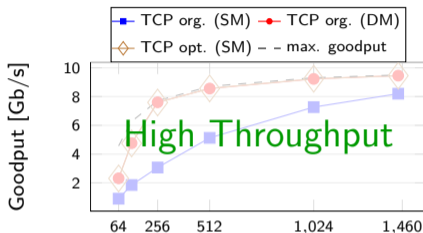
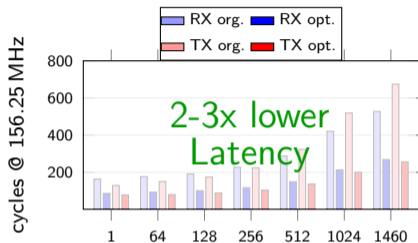
Results



These results enabled a consistent distributed key-value store [2]

	Mem. allocated	Mem. bandwidth
TCP org.	1,300 MB	40 Gbps
TCP opt.	650 MB	10 Gbps
Diff	-50%	-75%

Results



Visit our poster for more results and details!

Find the source at: <http://github.com/dsidler/fpga-network-stack>

	Mem. allocated	Mem. bandwidth
TCP org.	1,300 MB	40 Gbps
TCP opt.	650 MB	10 Gbps
Diff	-50%	-75%