# Consensus in a Box
## Inexpensive Coordination in Hardware
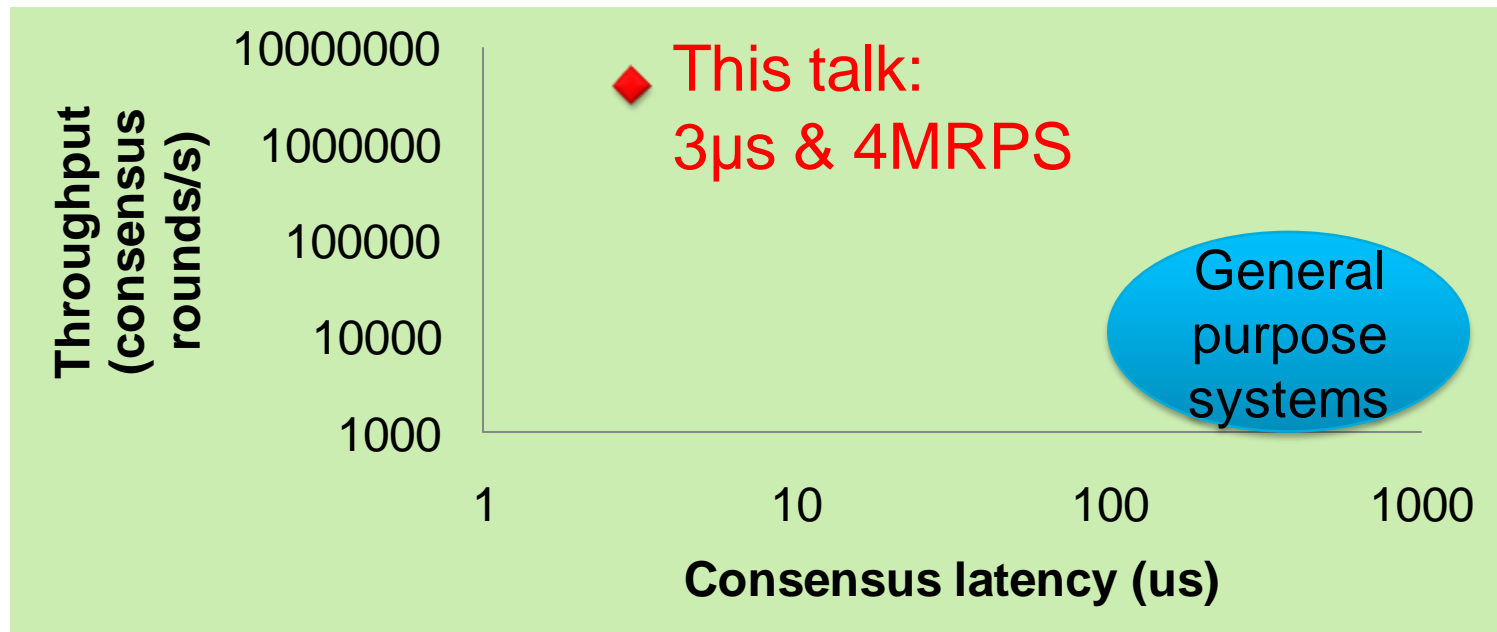
**Zsolt István, David Sidler, Gustavo Alonso, Marko Vukolic[*]**

Systems Group, Department of Computer Science, ETH Zurich

[*]IBM Research, Zurich

# Motivation: Cost of consensus

- Consensus is an essential function in datacenters
- How can consensus be made inexpensive?

# Related work

- Speeding up consensus is an important problem
  - Related work in networking, systems, HPC, etc.

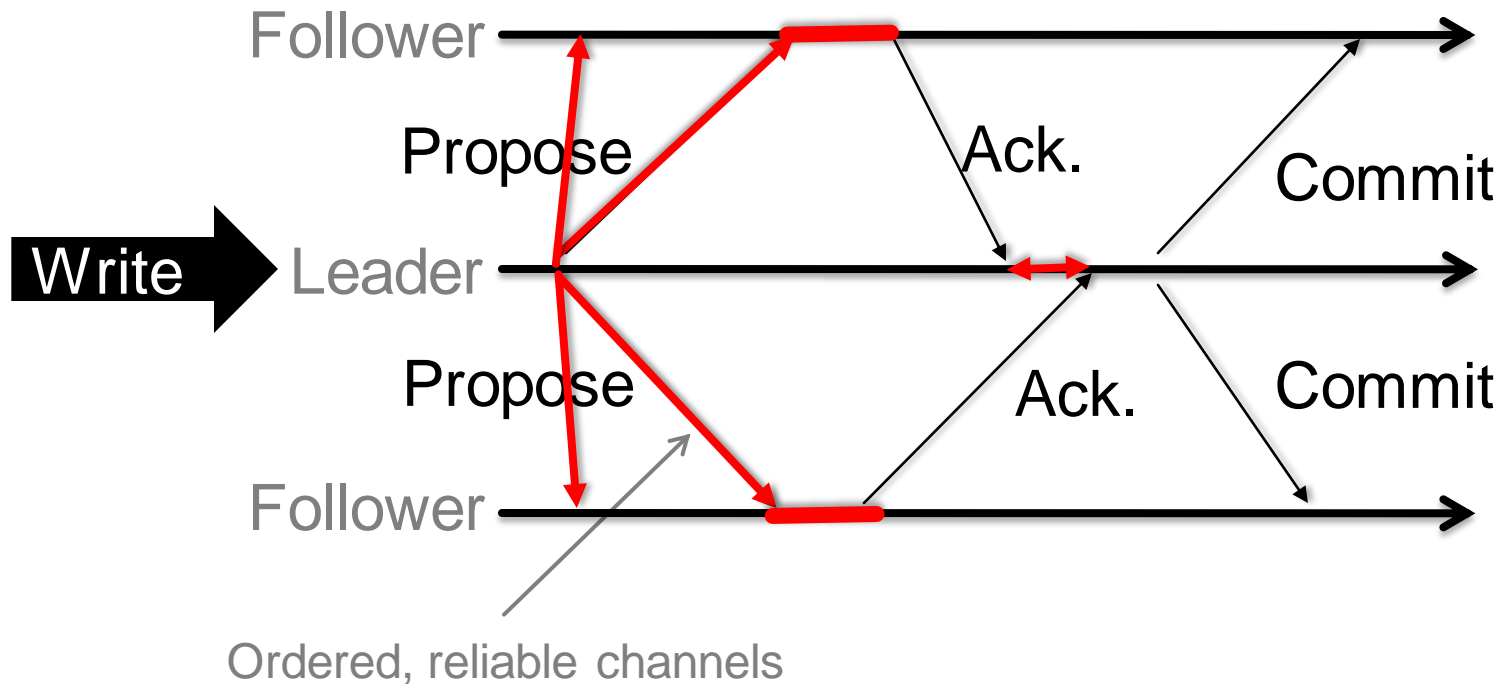- Specialized hardware can remove traditional limitations

[1] Zhang et al. Smartswitch: Blurring the line between network infrastructure & cloud applications. In HotCloud'14

[2] Mai et al. NetAgg: Using Middleboxes for Application-Specific On-path Aggregation in Data Centres. In CoNEXT'14

[3] Dang et al. NetPaxos: Consensus at Network Speed. In SOSR'15

[4] Poke et al. DARE: High-Performance State Machine Replication on RDMA Networks. In HPDC'15.
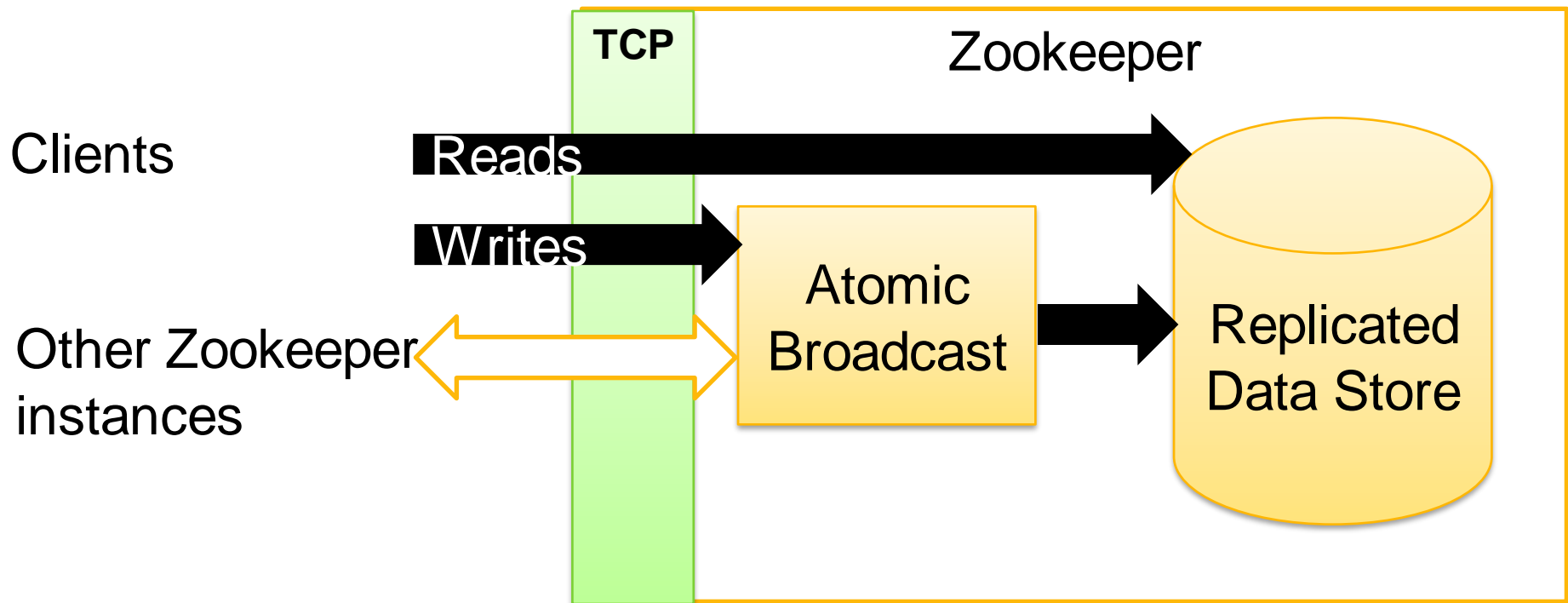
# Consensus in a Box

- **Why?** Consensus is expensive, but desired

- **What?** Atomic broadcast – Zookeeper's ZAB protocol

- **How?** Specialized processor on FPGA
  - Tight integration with 10Gbps networking + deep pipelining

- **Evaluation?** Drop-in replacement for Memcached with Zookeeper's replication
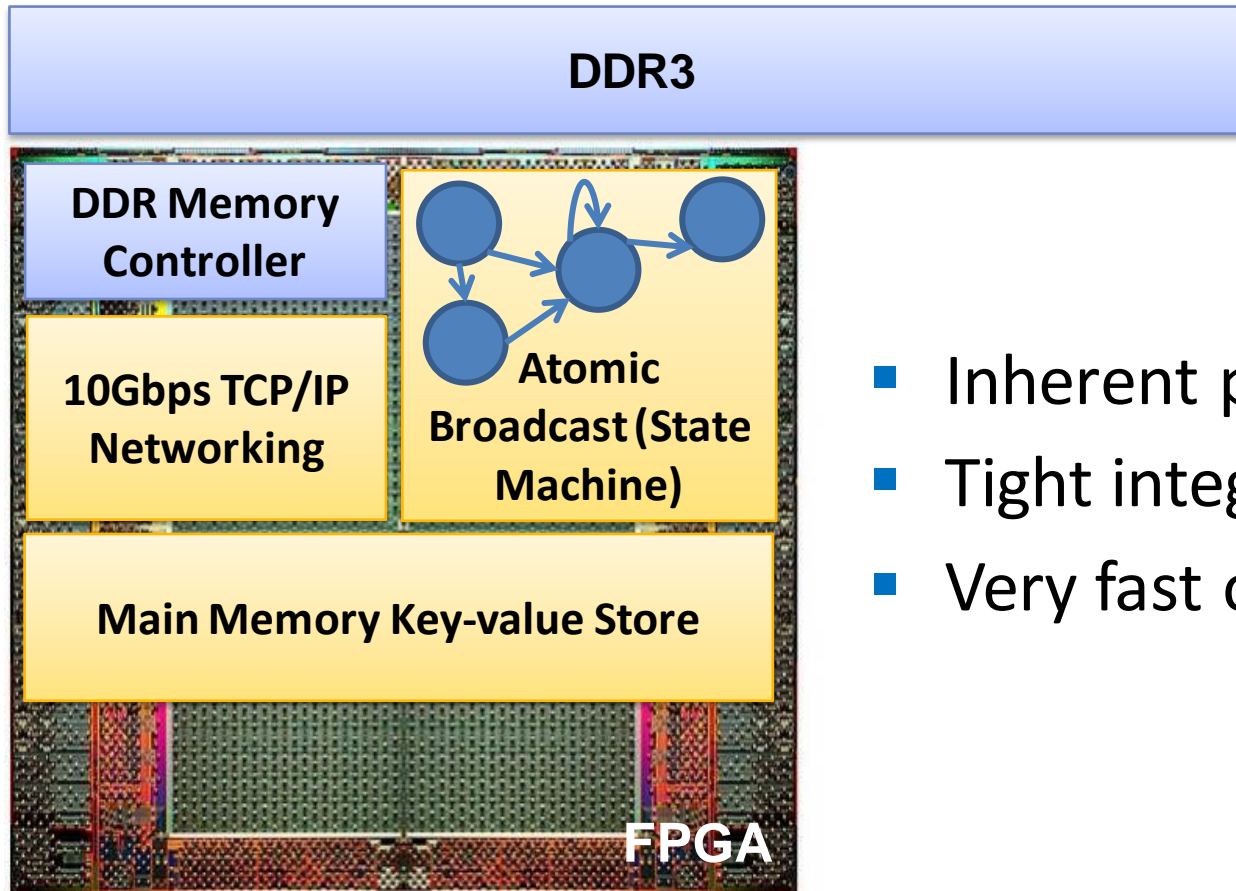
# Zookeeper's Atomic Broadcast



Follower

Propose

Write → Leader

Propose

Follower

Ack.

Commit

Ack.

Commit

Ordered, reliable channels

# Zookeeper from 10000ft



Clients

Other Zookeeper instances

TCP

Zookeeper

Reads

Writes

Atomic Broadcast

Replicated Data Store

# Specialized processor architecture



- Inherent parallelism
- Tight integration
- Very fast on-chip memory

# What makes it go fast…

## Latency

- **Networking optimizations**
  - Low-latency on-chip buffers for RX path
  - Datacenter and application-specific knowledge
- **Predictable behavior**
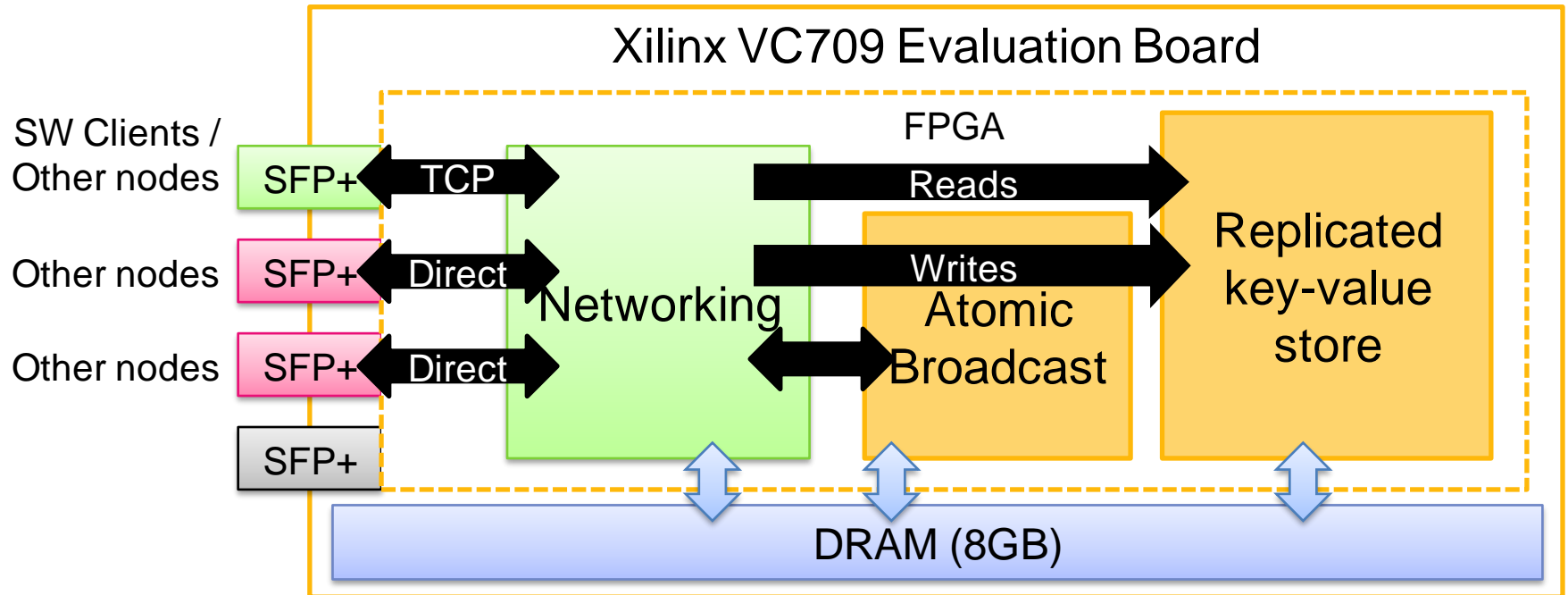  - Fast local caches for common case behavior
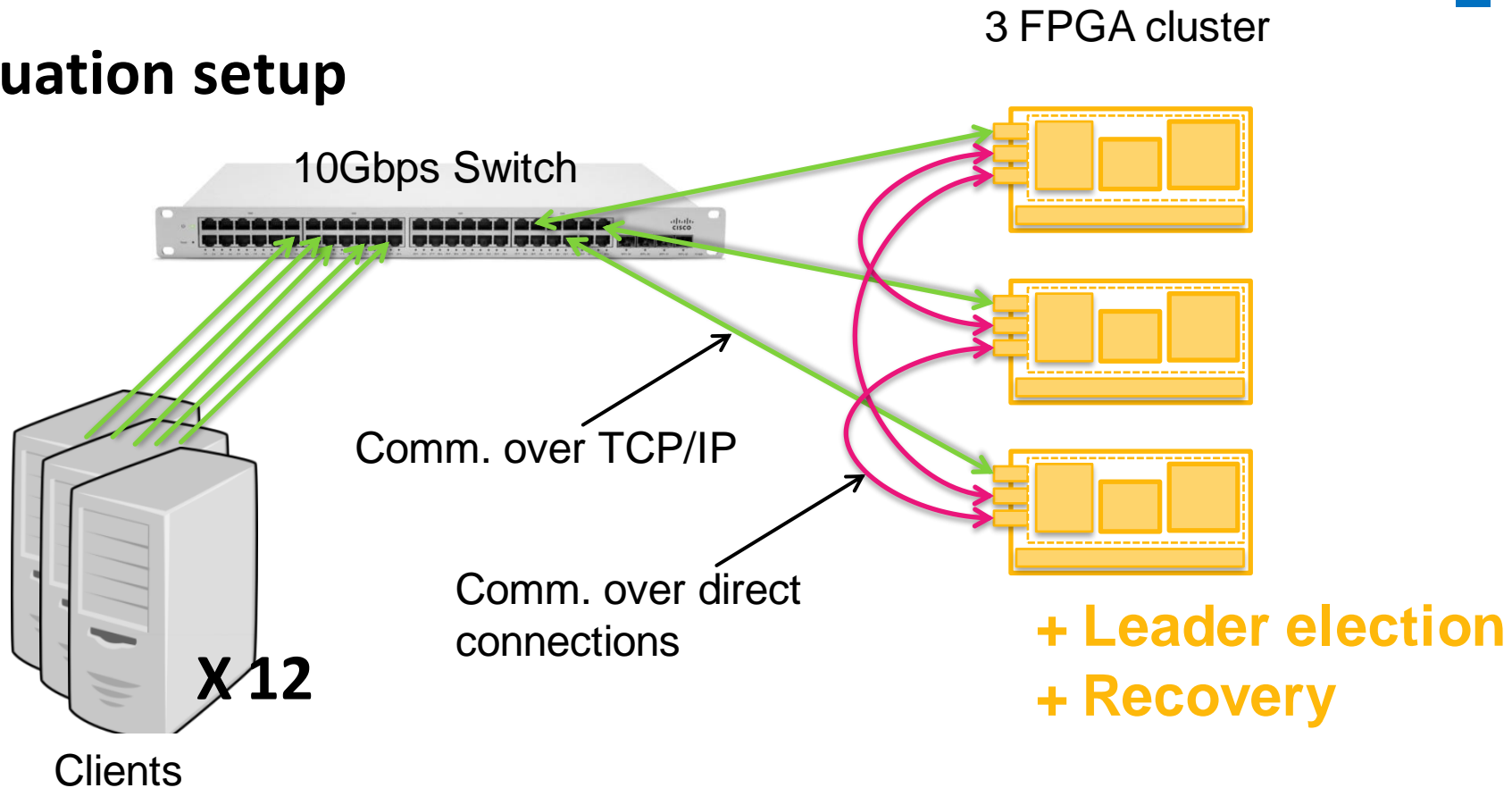
## Throughput

- **Pipelined execution**

Networking > Consensus > Key-value store

# Deployment and Evaluation

# Hardware platform

# Evaluation setup

3 FPGA cluster



10Gbps Switch

Comm. over TCP/IP

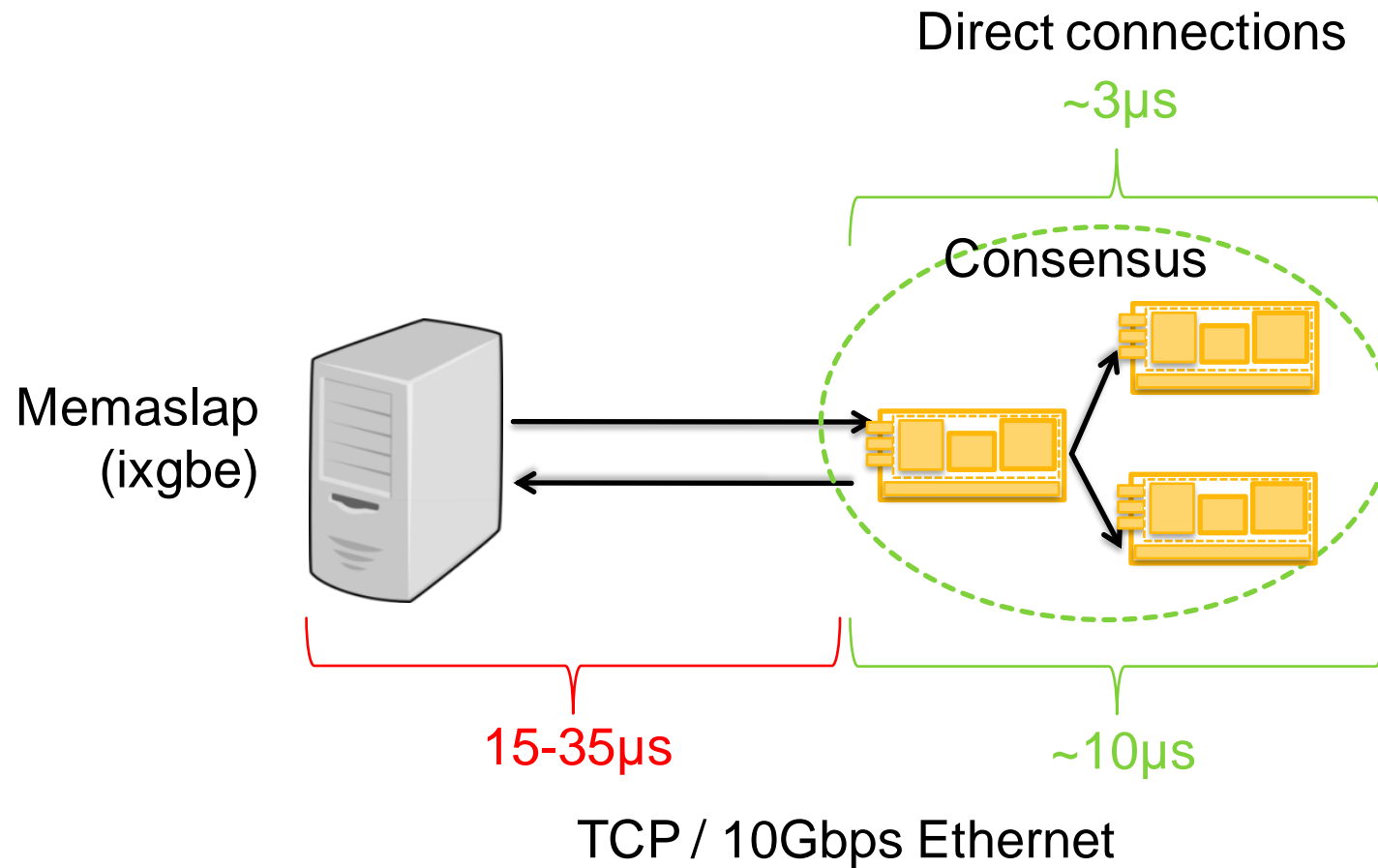Comm. over direct connections

X 12

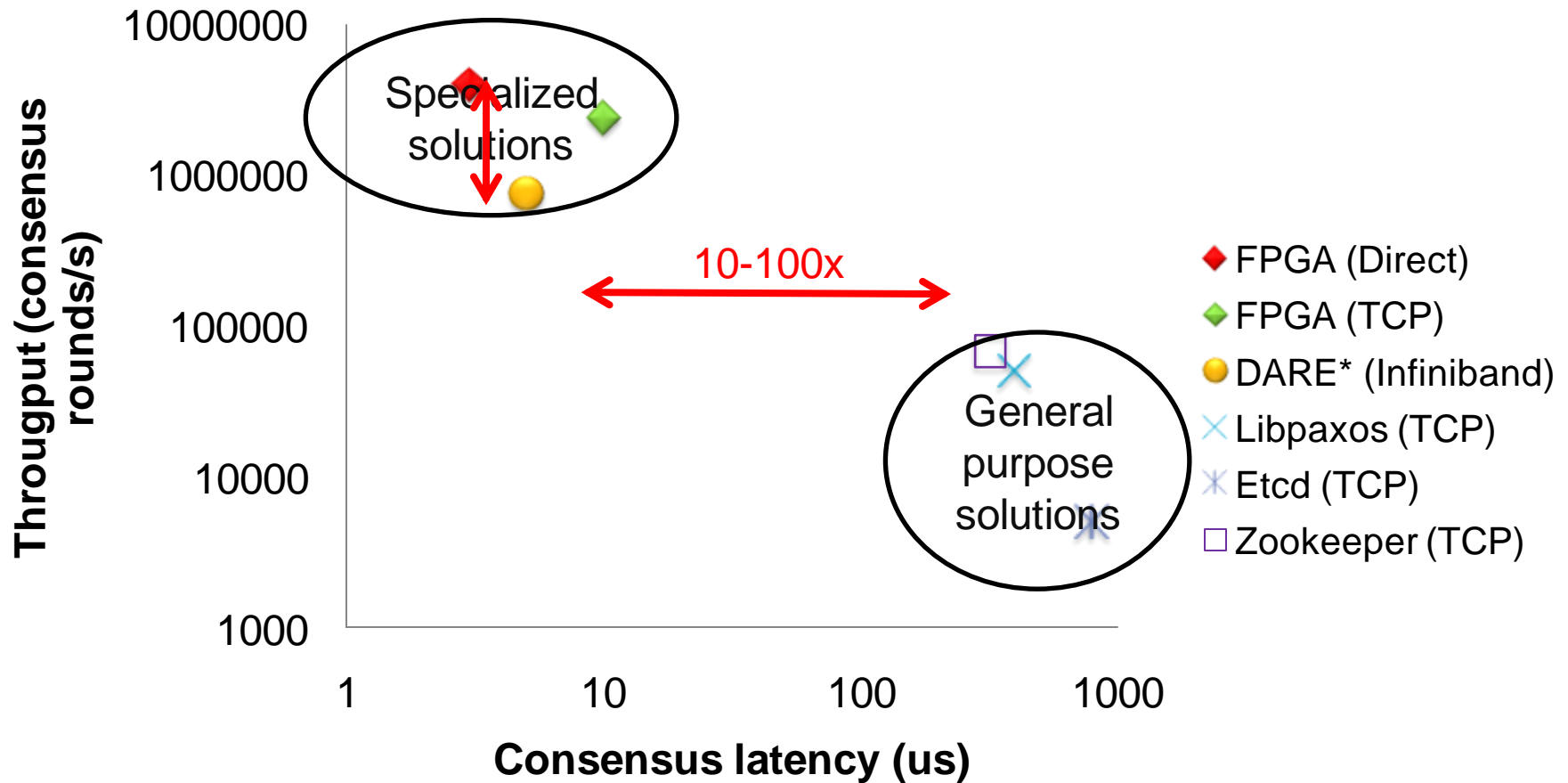Clients

**+ Leader election**
**+ Recovery**

- Drop-in replacement for Memcached with Zookeeper's replication
- Standard tools for benchmarking (libmemcached)
  - Simulating 100s of clients

# Latency of KVS writes (consensus)



Direct connections
~3µs

Consensus

Memaslap
(ixgbe)

15-35µs

~10µs

TCP / 10Gbps Ethernet

# The benefit of specialization...



Specialized solutions

10-100x

General purpose solutions

◆ FPGA (Direct)
◆ FPGA (TCP)
● DARE* (Infiniband)
✕ Libpaxos (TCP)
✳ Etcd (TCP)
□ Zookeeper (TCP)

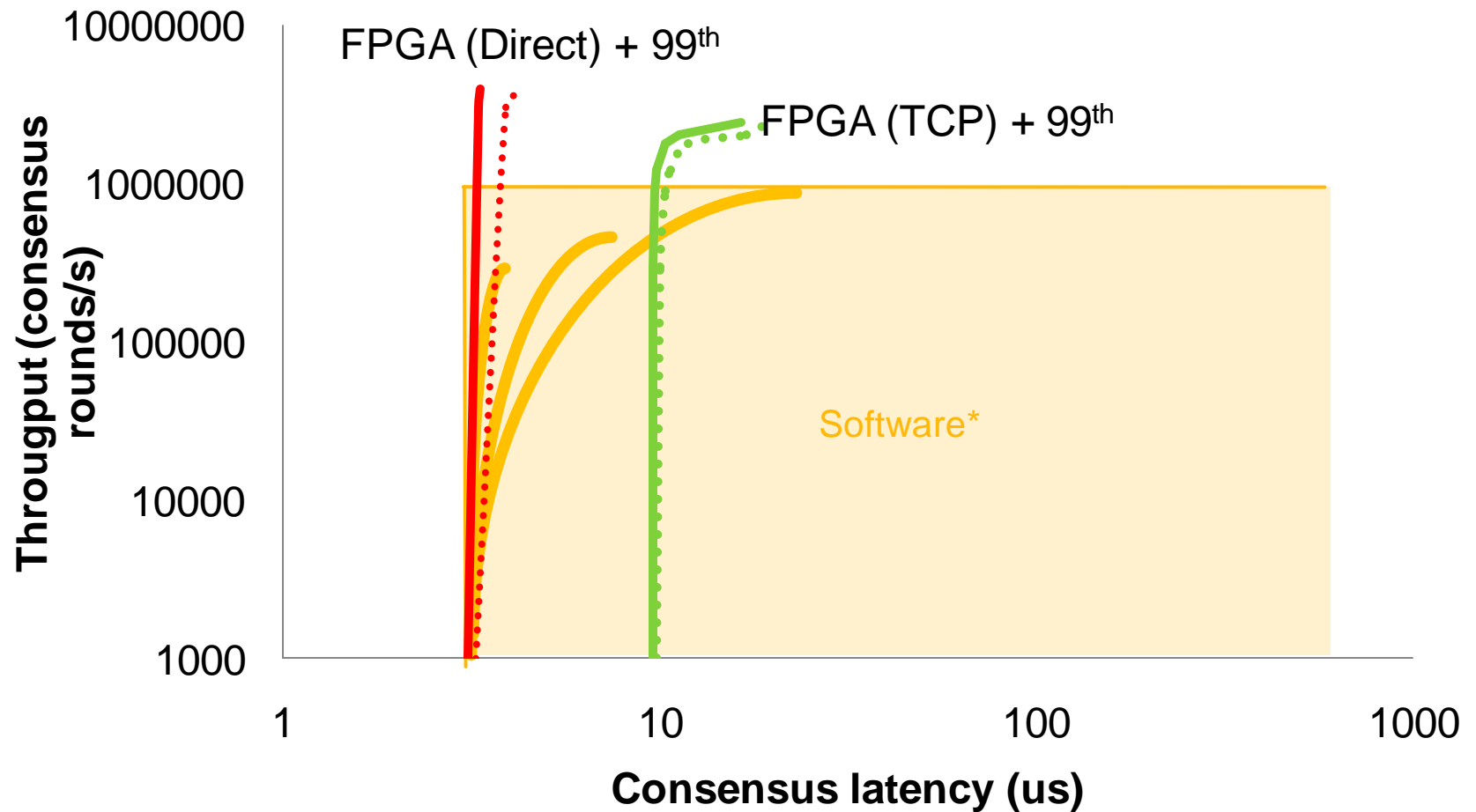Throughput (consensus rounds/s)

Consensus latency (us)

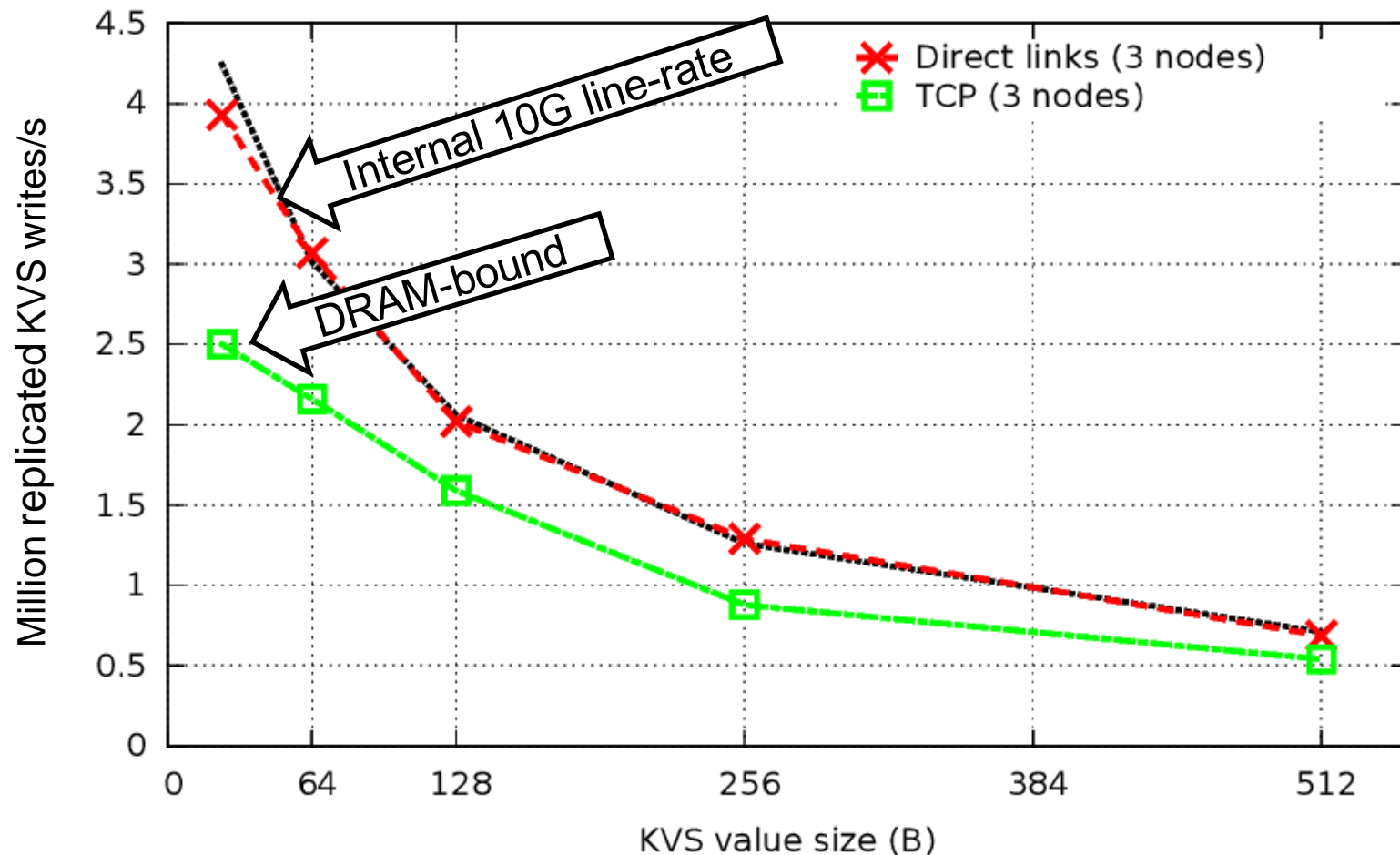[1] Dragojevic et al. FaRM: Fast Remote Memory. In NSDI'14.
[2] Poke et al. DARE: High-Performance State Machine Replication on RDMA Networks. In HPDC'15.
*=We extrapolated from the 5 node setup for a 3 node setup and removed estimated client overhead.

# Predictable hardware performance

# Throughput overview

- Consensus is expensive but often necessary
  - Solution: specialization and tight integration with networking

- We built high-throughput low-latency consensus in hardware



- Specialized hardware opens up new opportunities for smarter networks

{zsolt.istvan},{david.sidler}@inf.ethz.ch