

**T. HOEFLER**

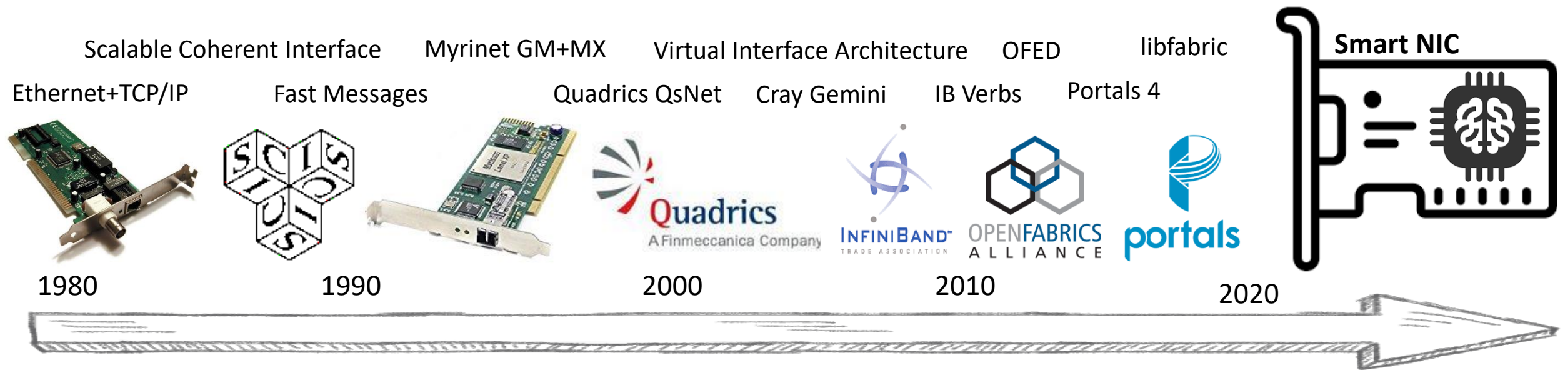
WITH S. DI GIROLAMO, K. TARANOV, D. DE SENSI, L. BENINI, R. E. GRANT, R. BRIGHTWELL, A. KURTH, M. SCHAFFNER, T. SCHNEIDER, J. BERÁNEK, M. BESTA, L. BENINI, D. ROWETH

# General in-network processing – time is ripe!

Keynote at the High-Performance Interconnects Forum with HPC China 2020



# The Development of High-Performance Networking Interfaces



sockets

(active) message based

protocol offload

remote direct memory access (RDMA)

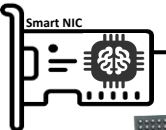


coherent memory access

OS bypass

zero copy

triggered operations



ARM cores  
(with full OS, outside packet pipe)

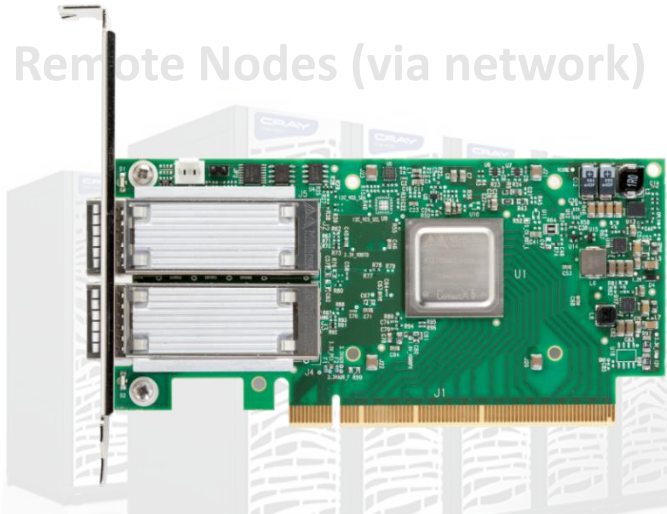
Flow Processors  
(limited flexibility, P4)



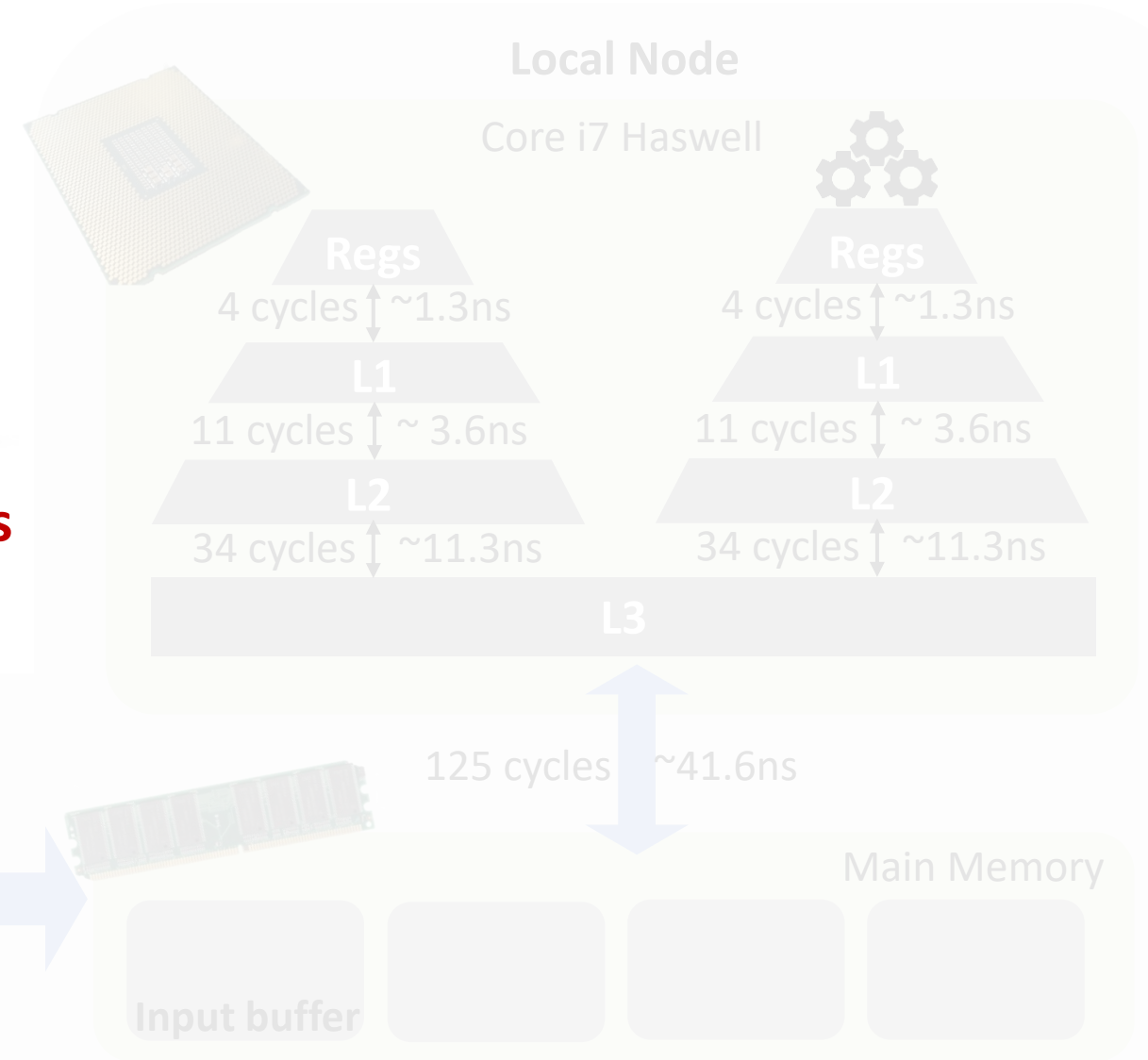
FPGAs  
(limited productivity,  
silicon efficiency)



# Data Processing in modern RDMA networks

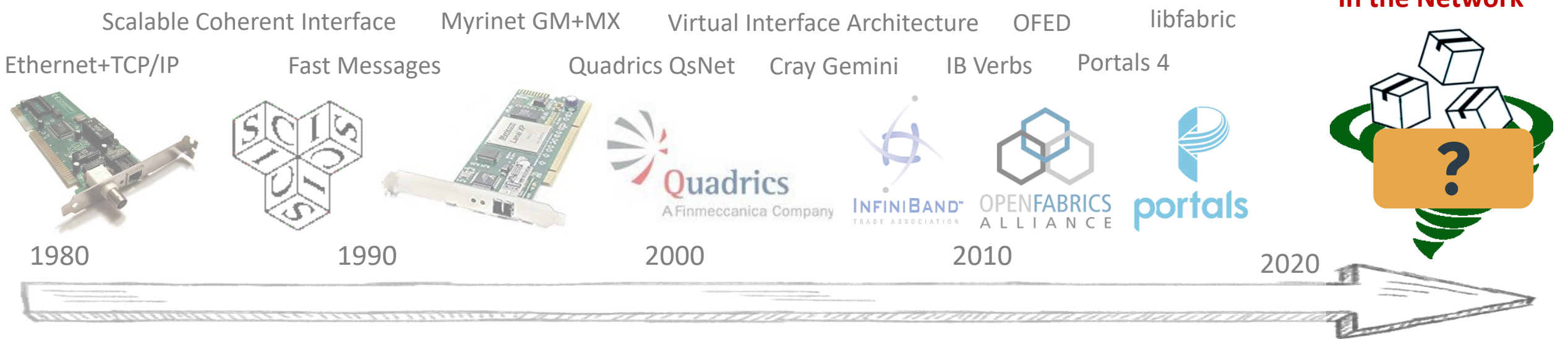


**Mellanox Connect-X5: 1 packet/5ns**  
**Tomorrow (400G): 1 packet/1.2ns**



# The future of High-Performance Networking Interfaces

**SPIN**  
Streaming Processing  
In the Network



sockets      (active) message based      protocol offload      remote direct memory access (RDMA)      ?  
 coherent memory access      OS bypass      zero copy      triggered operations      packet handlers

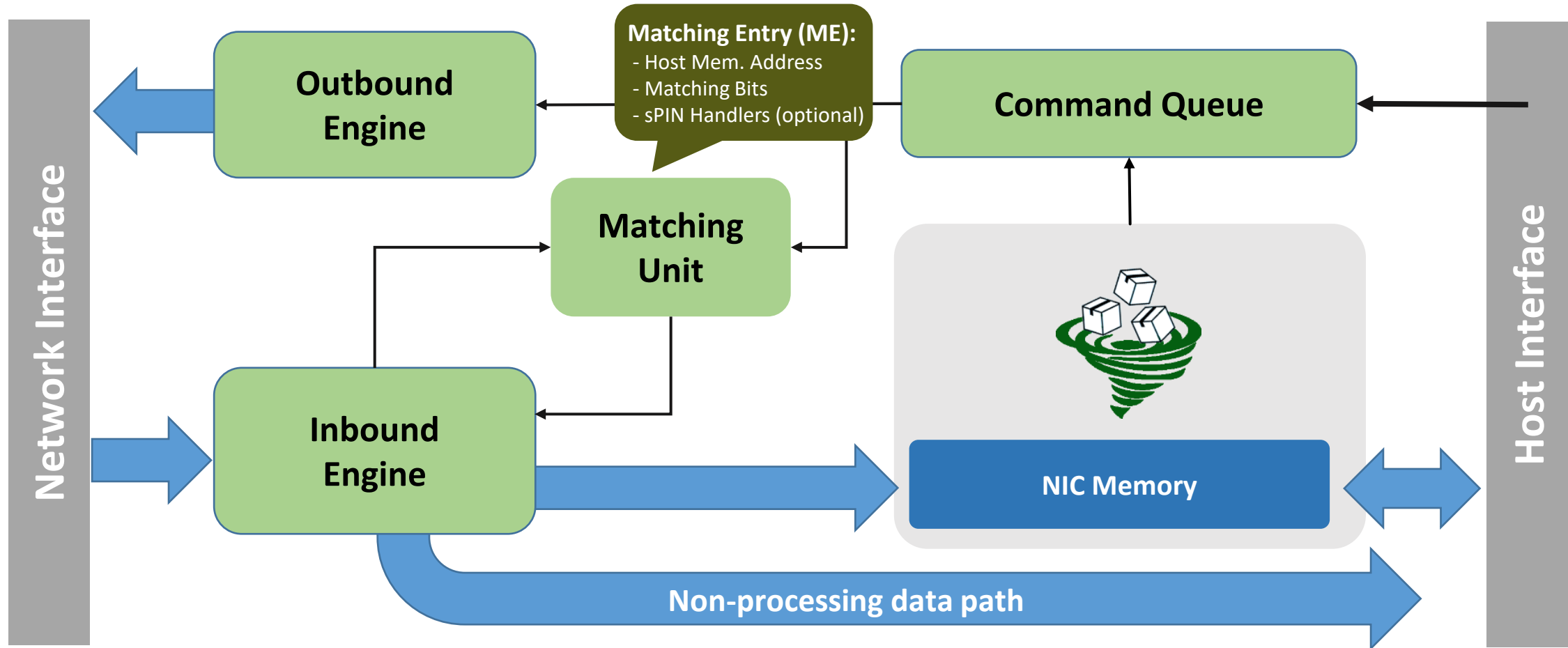
### Established Principles for Compute Acceleration

OpenGL, DirectX11 → Generalization → NVIDIA CUDA, OpenMP 4.0, OpenCL  
 Revolutionizes Acceleration

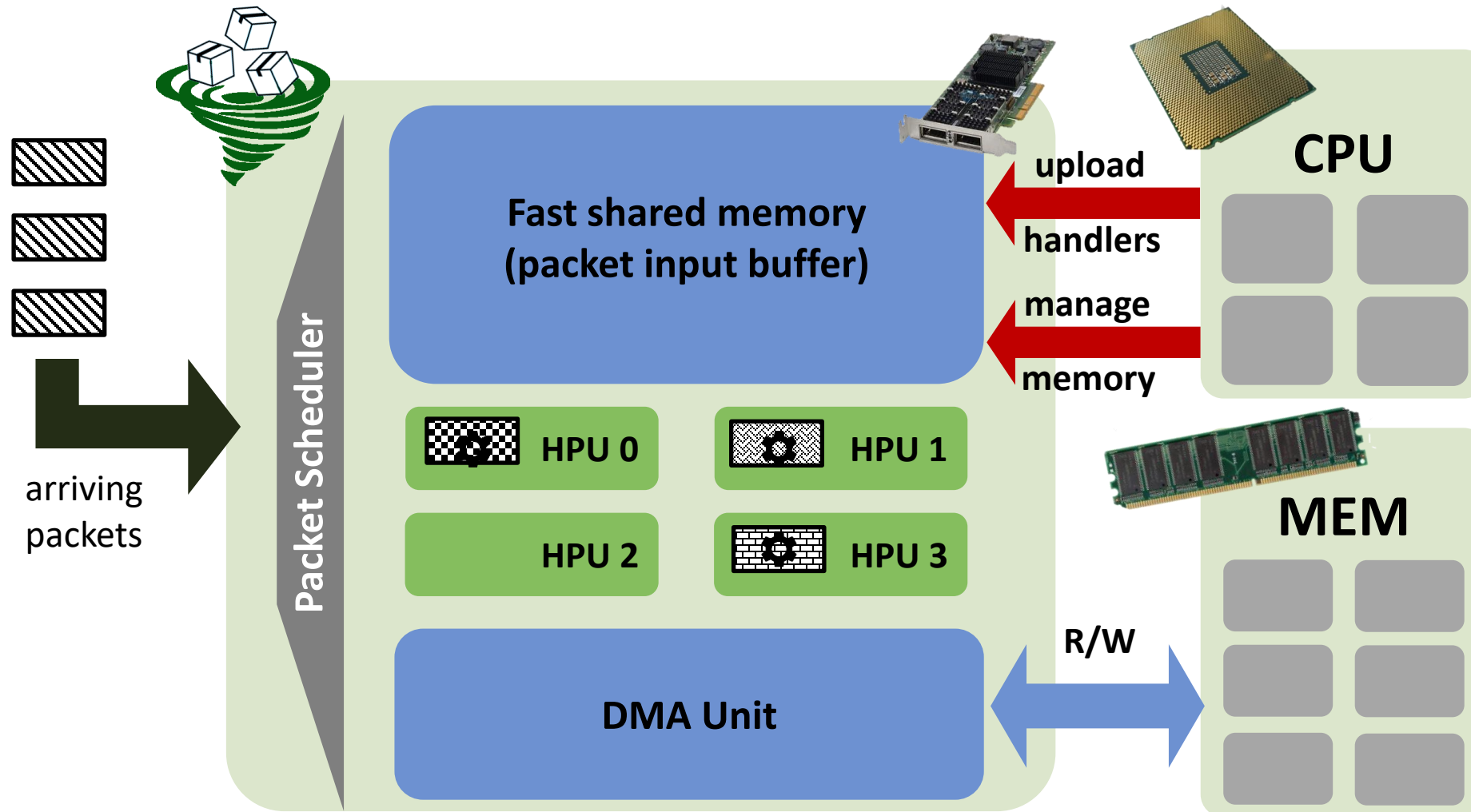
### Where do we stand in Network Acceleration?

P4, eBPF → Generalization → Data Acceleration → RISC-V

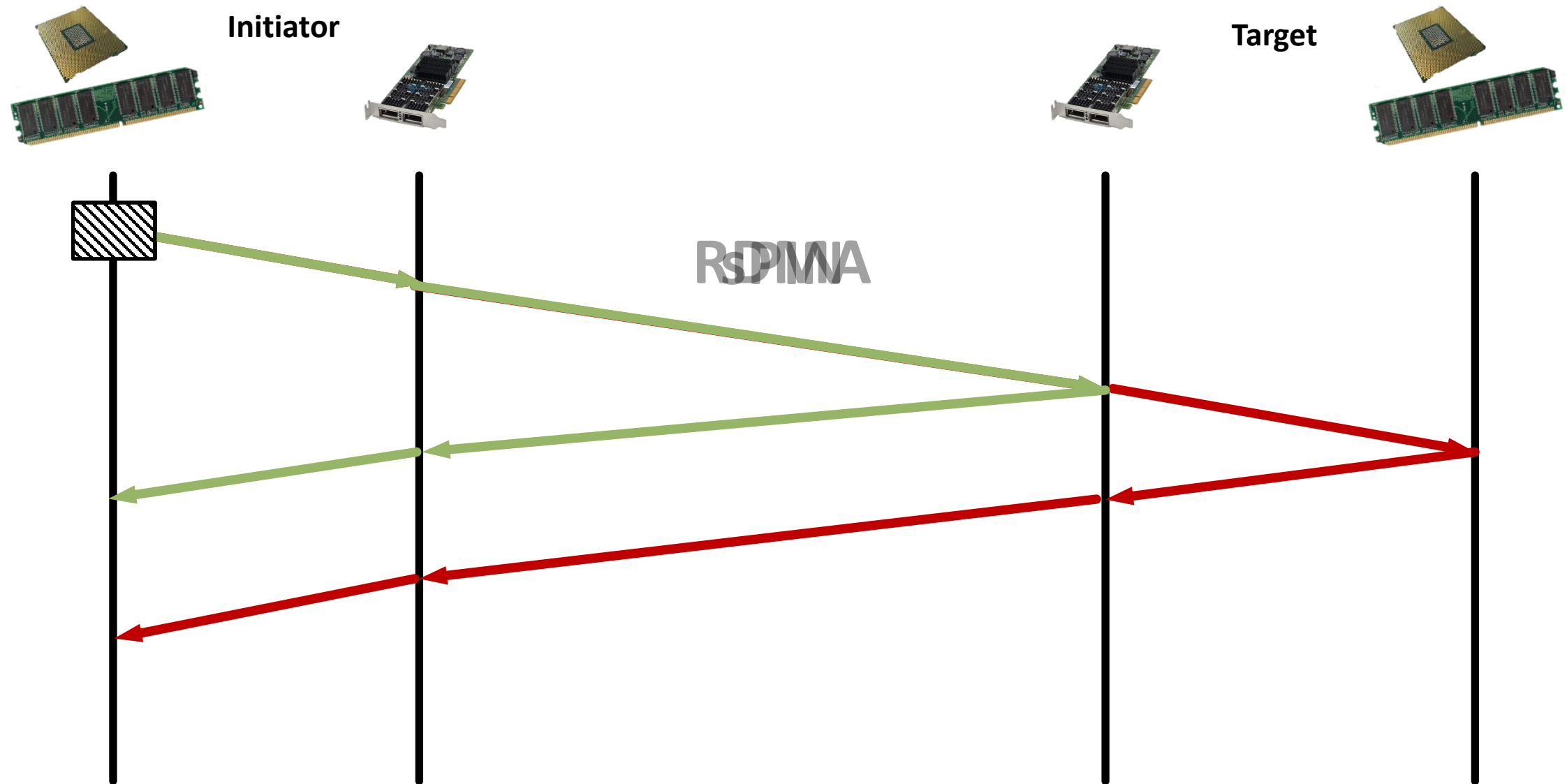
# sPIN NIC – Architecture for fast Network Processing



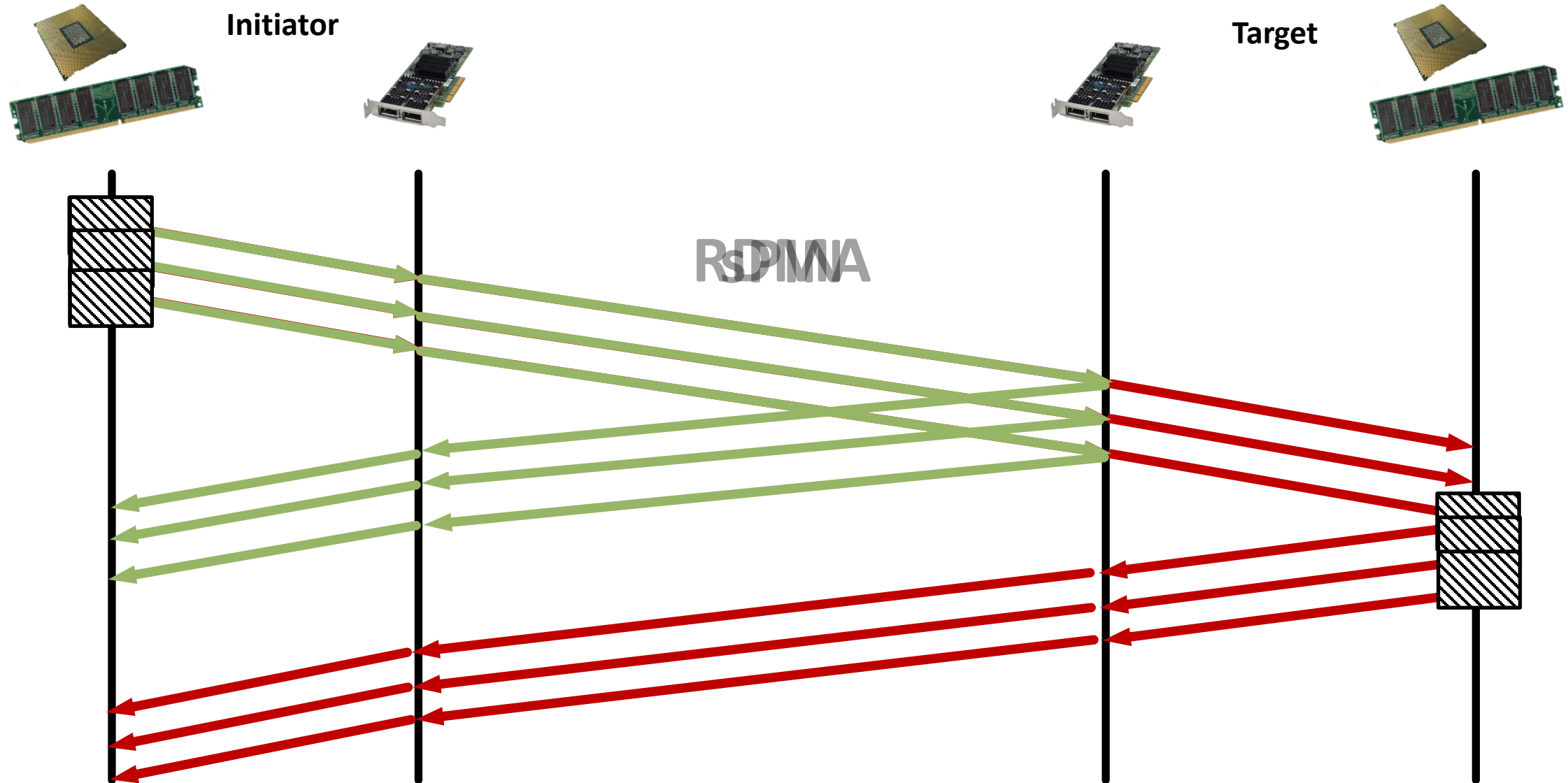
# sPIN NIC - Abstract Machine Model for Packet Processing



# RDMA vs. sPIN in action: Simple Ping Pong

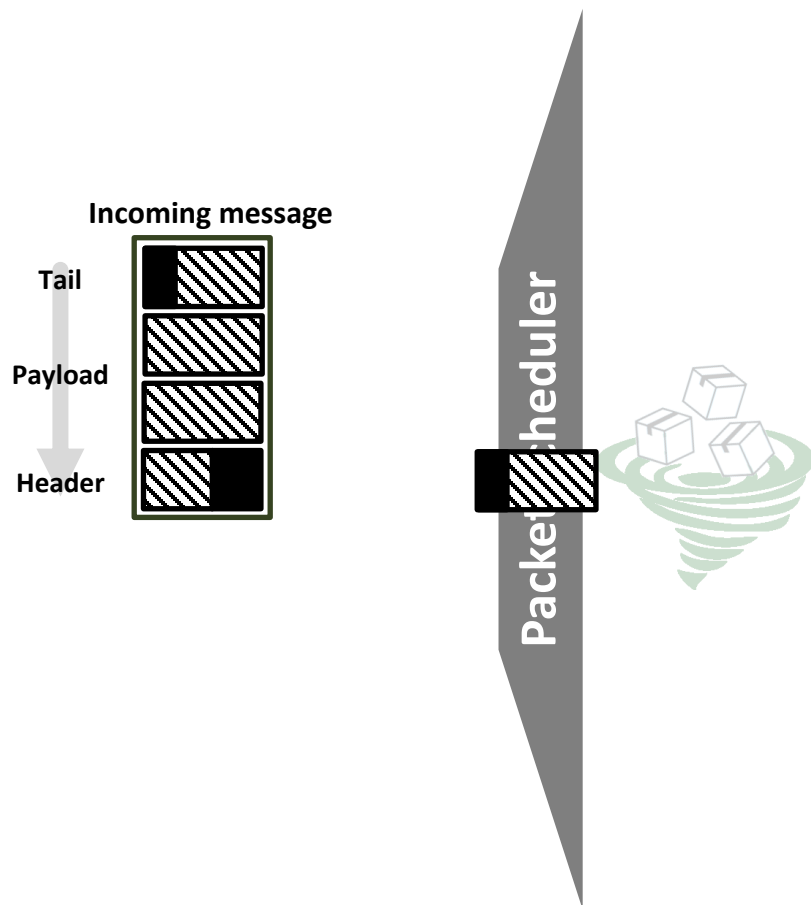


# RDMA vs. sPIN in action: Streaming Ping Pong





# sPIN – Programming Interface



## Header handler

```
__handler int pp_header_handler(const ptl_header_t h, void *state) {
    pingpong_info_t *i = state;
    i->source = h.source_id;
    return PROCESS_DATA; // execute payload handler to put from device
}
```

## Payload handler

```
__handler int pp_payload_handler(const ptl_payload_t p, void *state) {
    pingpong_info_t *i = state;
    PtlHandlerPutFromDevice(p.base, p.length, 1, 0, i->source, 10, 0, NULL, 0);
    return SUCCESS;
}
```

## Completion handler

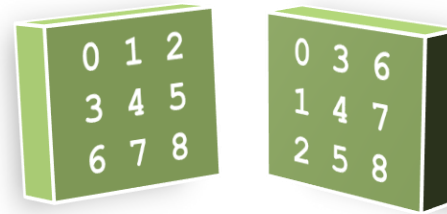
```
__handler int pp_completion_handler(int dropped_bytes,
                                     bool flow_control_triggered, void *state) {
    return SUCCESS;
}
```

```
connect(peer, /* ... */, &pp_header_handler, &pp_payload_handler, &pp_completion_handler);
```

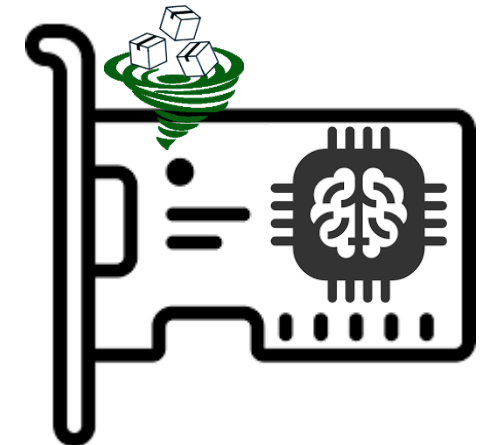
# Talk roadmap



Motivation and Overview

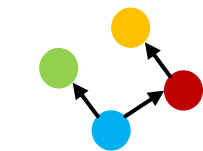


Data Layout Transformation



Hardware Implementation

## further use cases



Network Group Communication



Distributed Data Management

...

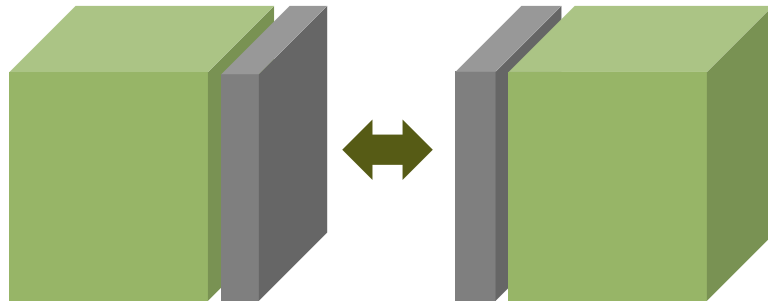
## Application domain



## Memory layout



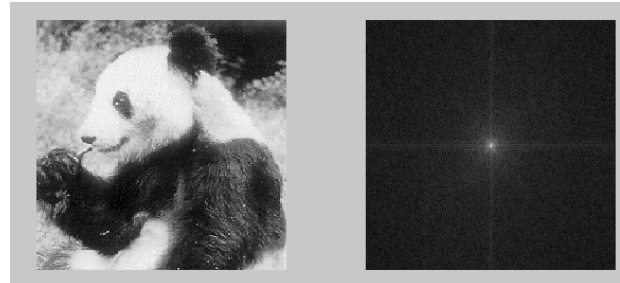
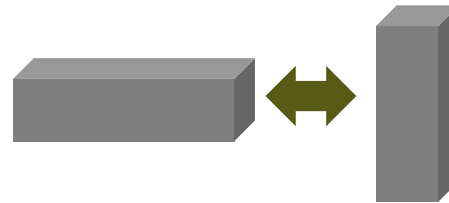
## Structured Exchange



<https://specfem3d.readthedocs.io/en/latest/>

L. Carrington et al. High-frequency simulations of global seismic wave propagation using SPEC-FEM3D\_GLOBE on 62K processors. SC 2008.

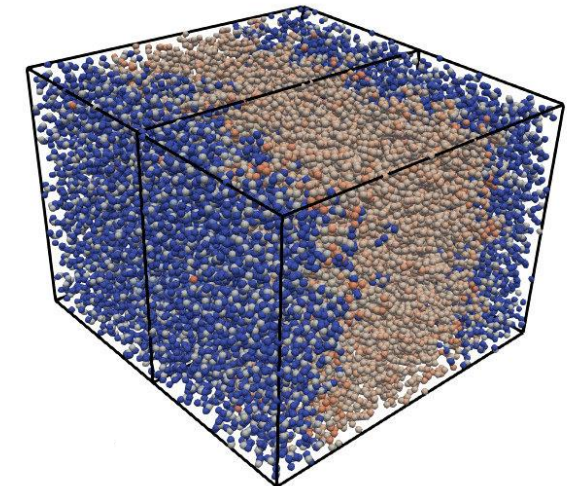
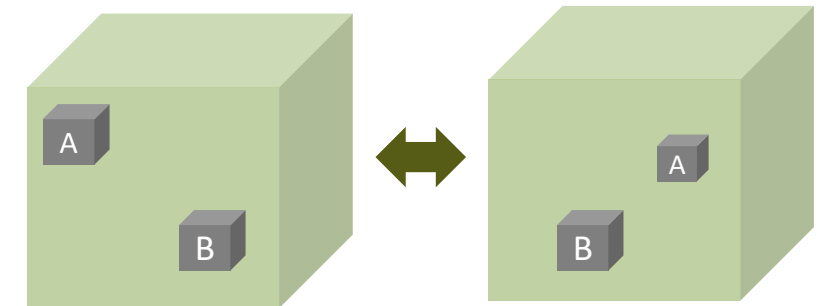
## Reshaping



<http://fourier.eng.hmc.edu/e161/lectures/fourier/node10.html>

T. Hoefler et al. Parallel zero-copy algorithms for fast Fourier transform and conjugate gradient using MPI datatypes. EuroMPI 2010.

## Unstructured Exchange



W. Usher et al. libIS: a lightweight library for flexible in transit visualization. ISAV 2018.

# Programming Support for Non-Contiguous Transfers

**ARMCI**

**CAF**

**Chapel**

**Portals 4**

**SHMEM**

**UPC**

**X10**

**MPI**

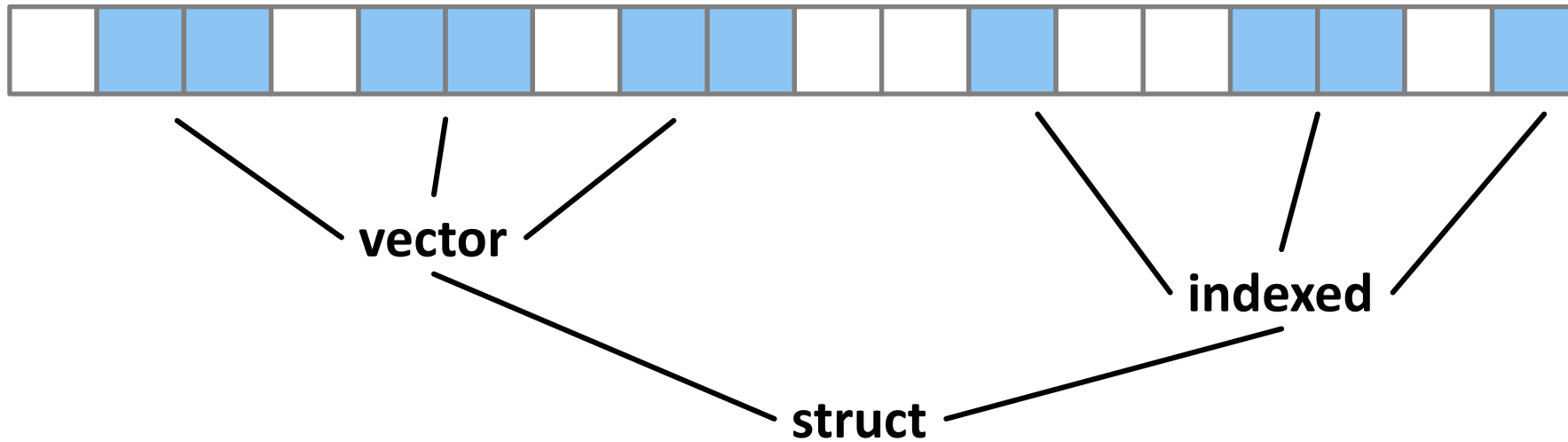
I/O Vectors

Strided transfers

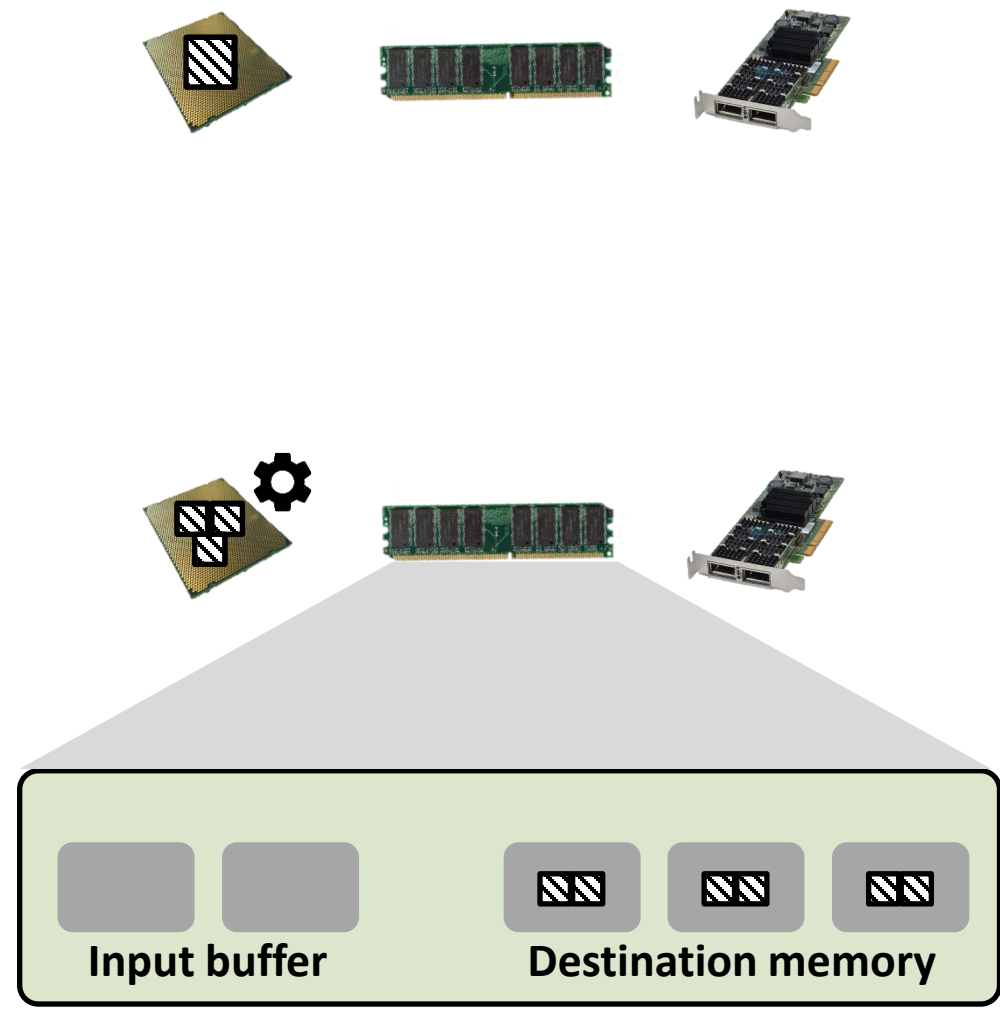
Compiler-Assisted Aggregation

Derived Datatypes

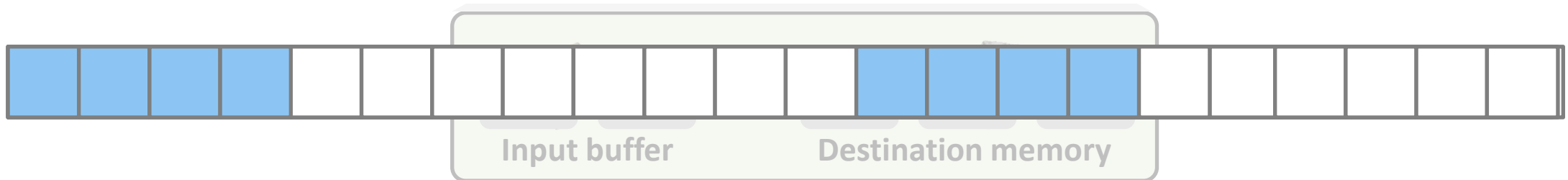
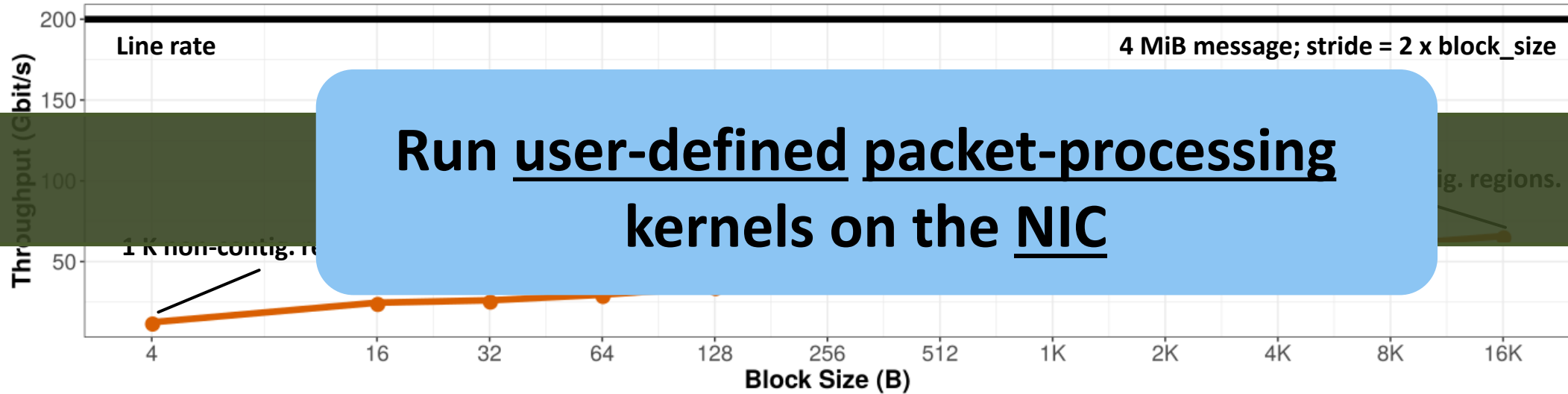
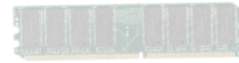
Support for multiple strides (e.g., 3D faces)



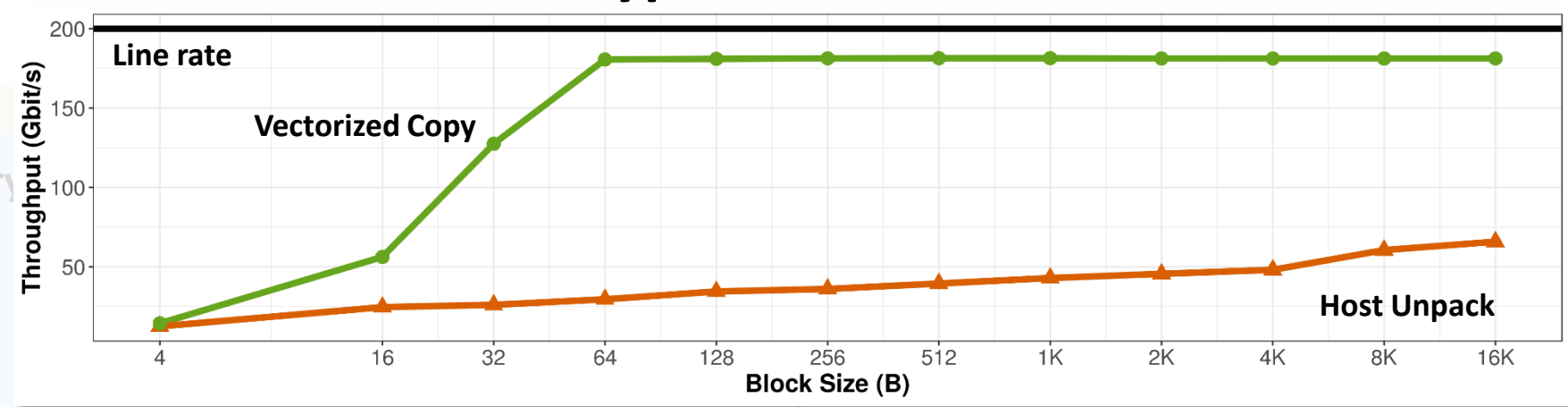
# State of the Art in MPI Datatypes Processing



# State of the Art in MPI Datatypes Processing



# A simple vectorize scatter datatype



vector

indexed

struct

```

_handler vector payload_handler(handler_args_t *args)
{
    spin_vec_t v;
    uint32_t num;
    uint32_t sti;

    uint8_t *pkt;

    uint8_t *hos;
    uint32_t hos;
    uint8_t *hos;

    for (uint32_t i = 0; i < num; i++)
    {
        PtlHandl
        pkt_payl
        host_adc
    }

    return SPIN_
}
    
```

DMA all contig. regions contained in the packet

Can we define a general handler to process arbitrary datatypes?



# Porting the MPI Types Library [1] to sPIN



## NIC Memory

Index{ #blocks: 2, blocklen: 1,  
offsets: {0, x}, basetype: \* }



Vector{ #blocks: 3, blocklen: 2,  
stride: 3, basetype: □ }

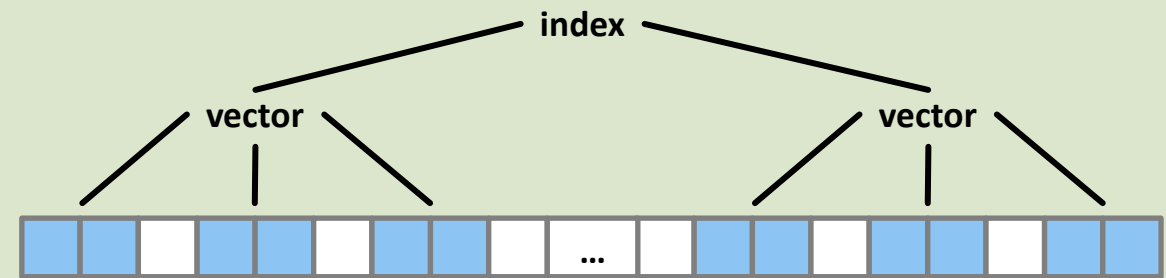
Dataloops

Index: ██████████  
Vector: ██████████

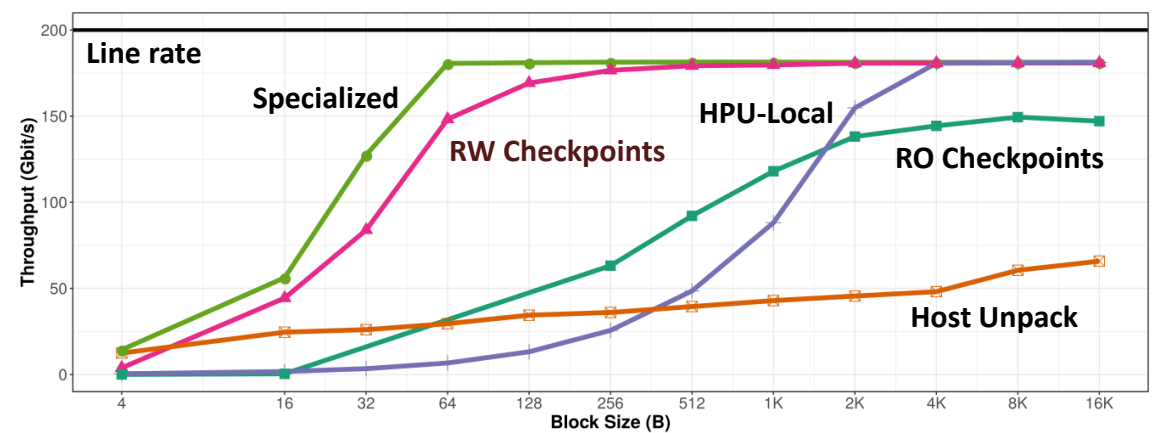
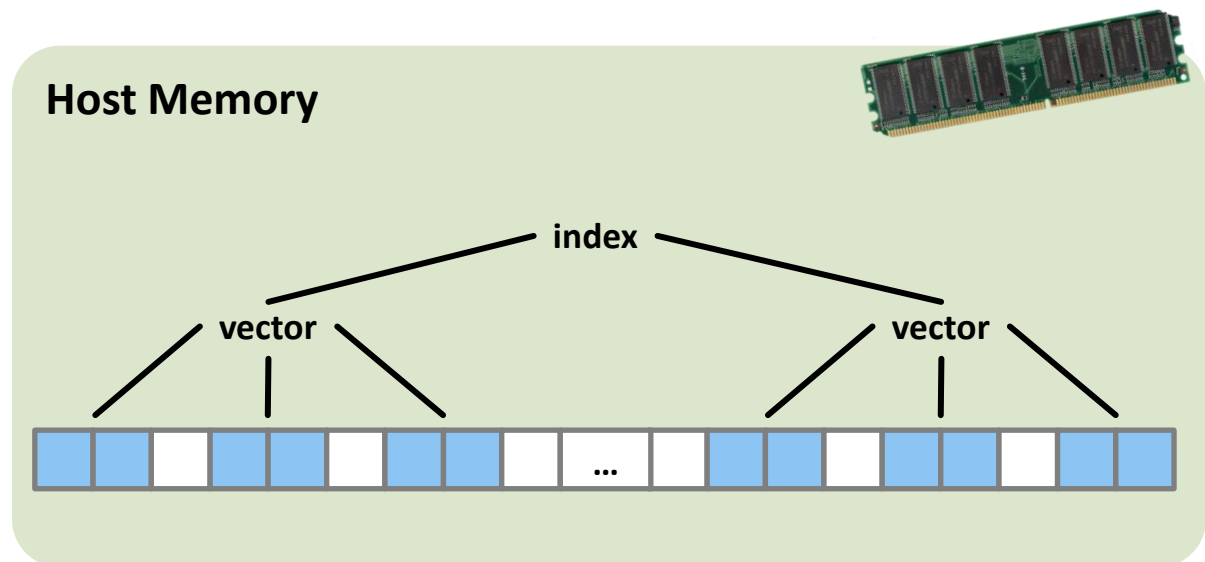
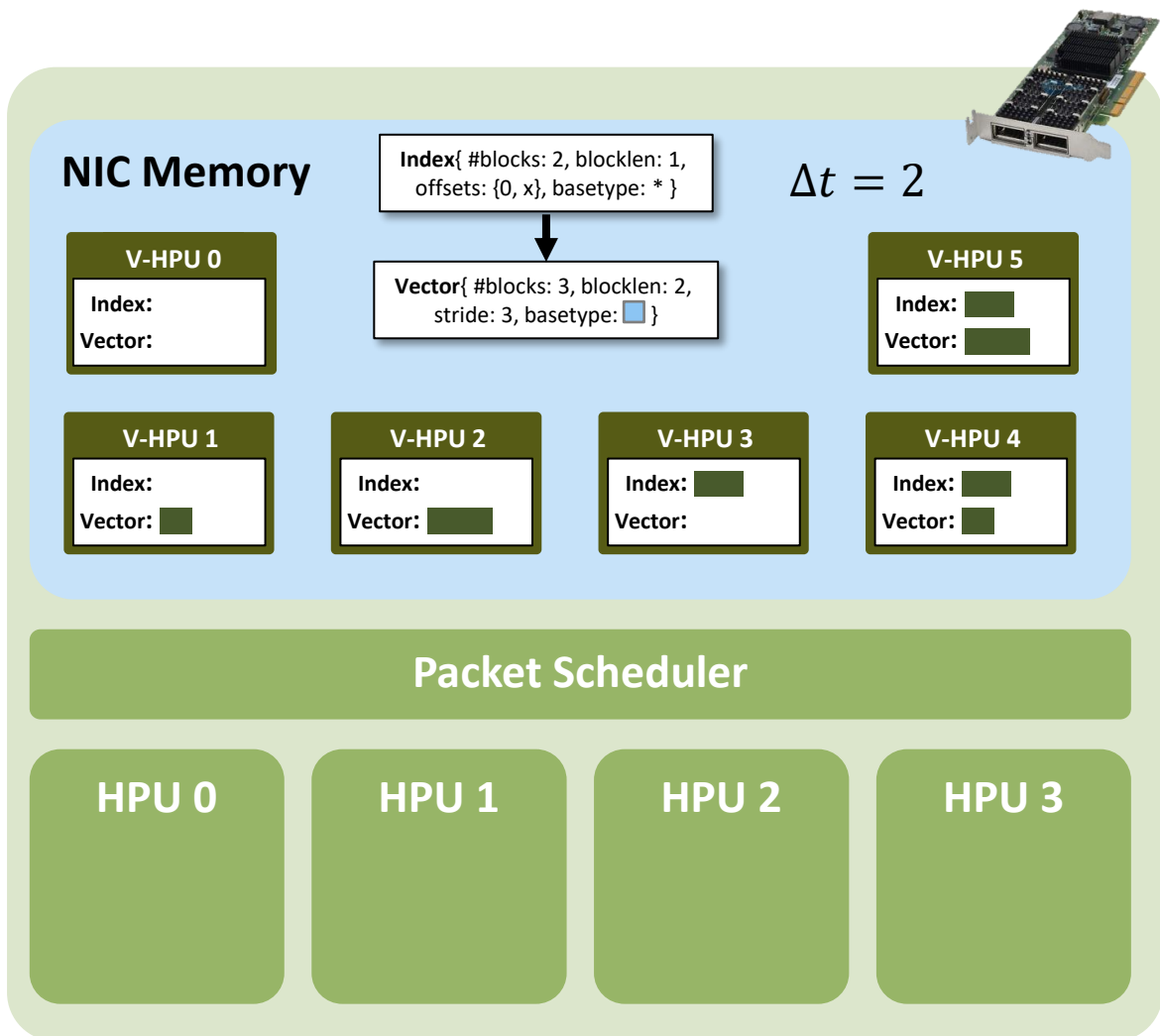
Segment



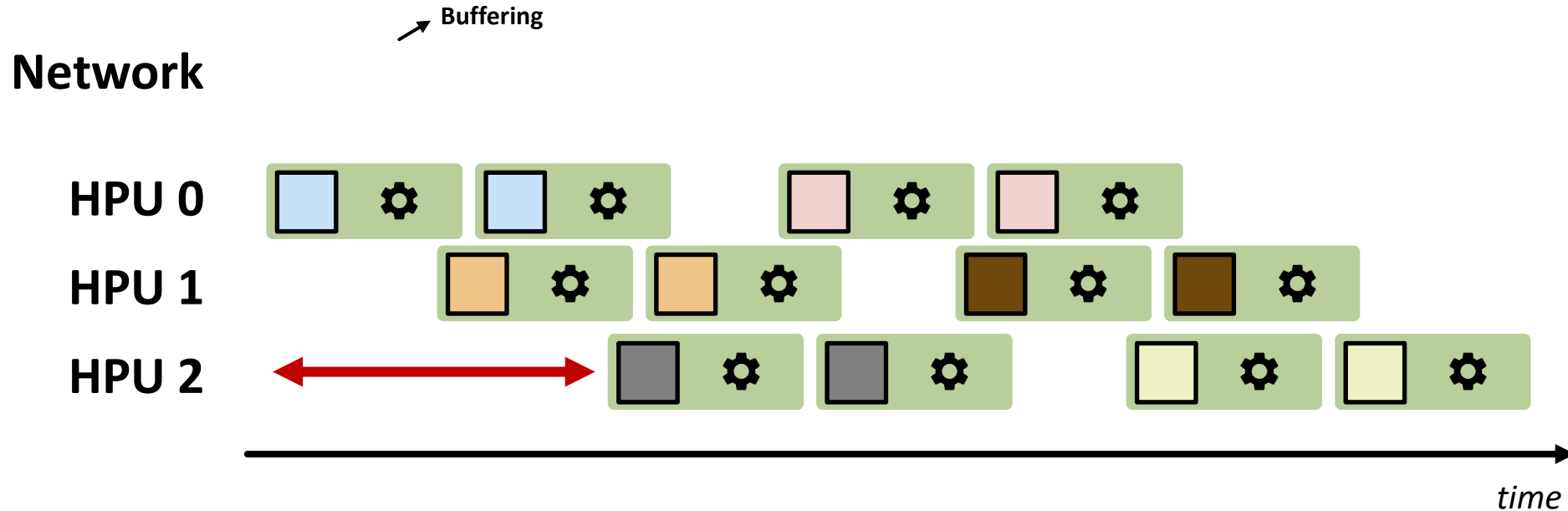
## Host Memory



# MPI Types Library on sPIN: Read-Write Checkpoints



# Checkpoint Interval Selection



$$T_C = T_{pkt} + \left\lceil \frac{\Delta r}{k} \right\rceil \cdot (P - 1) \cdot T_{pkt} + \left\lceil \frac{n_{pkt}}{P} \right\rceil \cdot T_{PH}(\gamma)$$

1

Limit the impact of the scheduling overhead

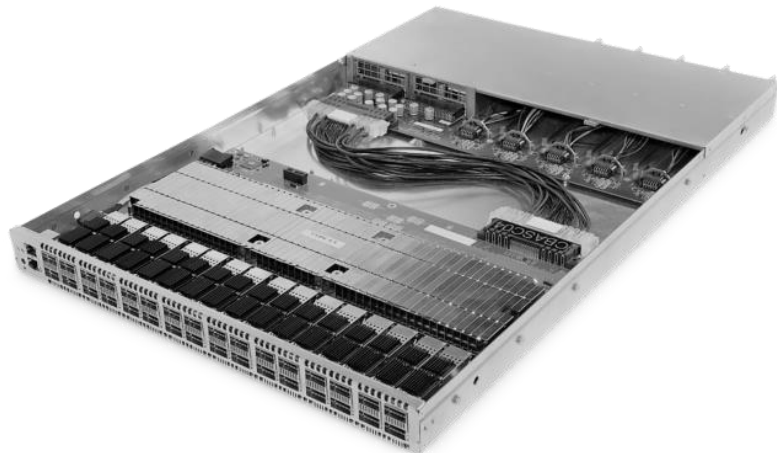
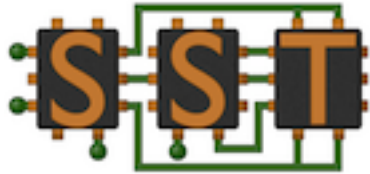
2

Do not saturate NIC memory with checkpoints

3

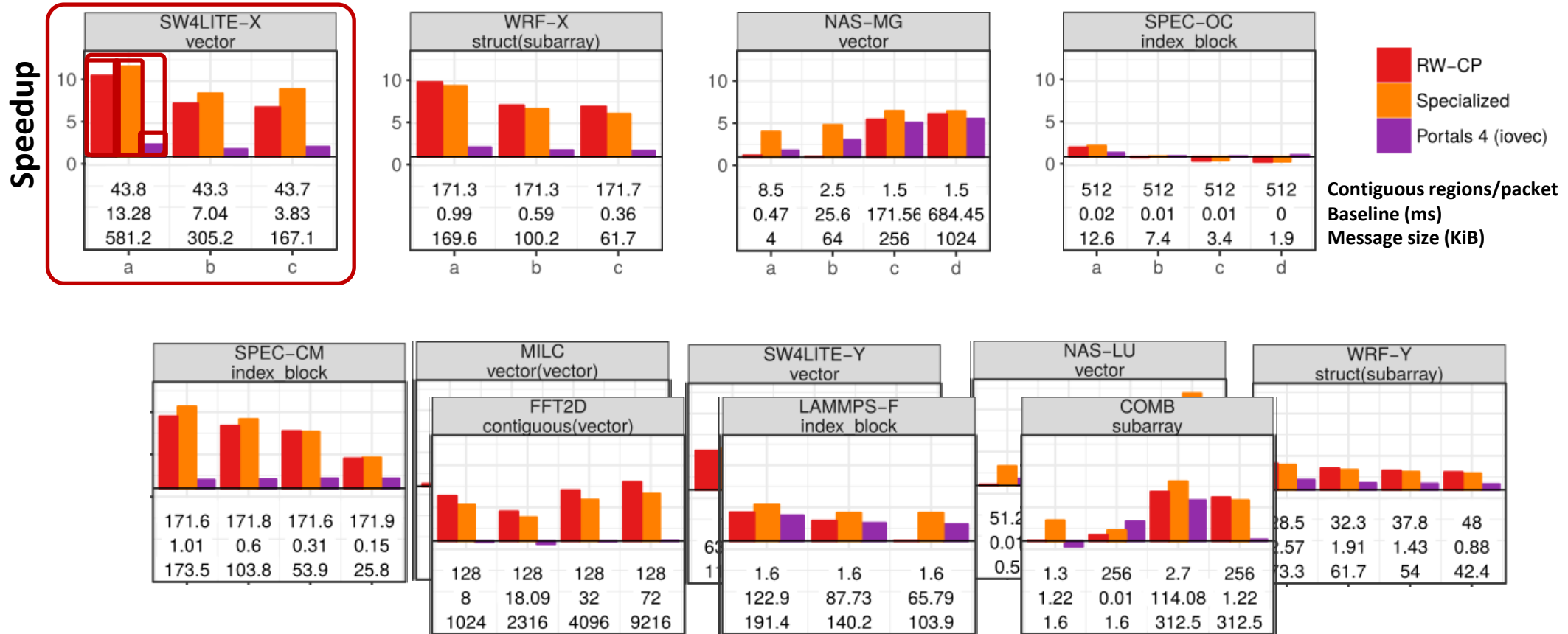
Do not saturate the packet buffer

## Cray Slingshot Simulator



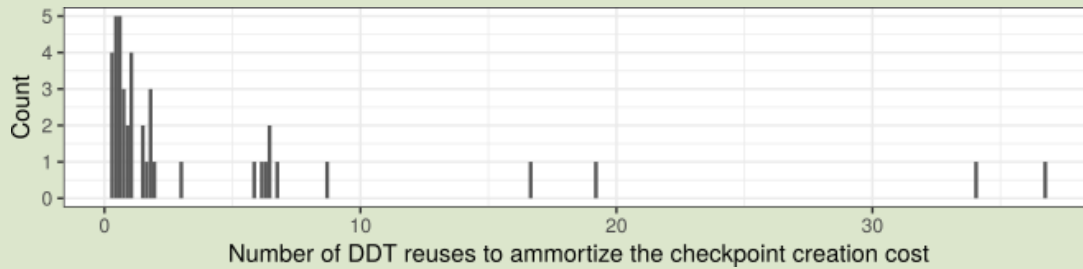
**32 Cortex A15 @800 MHz,  
single-cycle access memory**

# Real Application DDTs

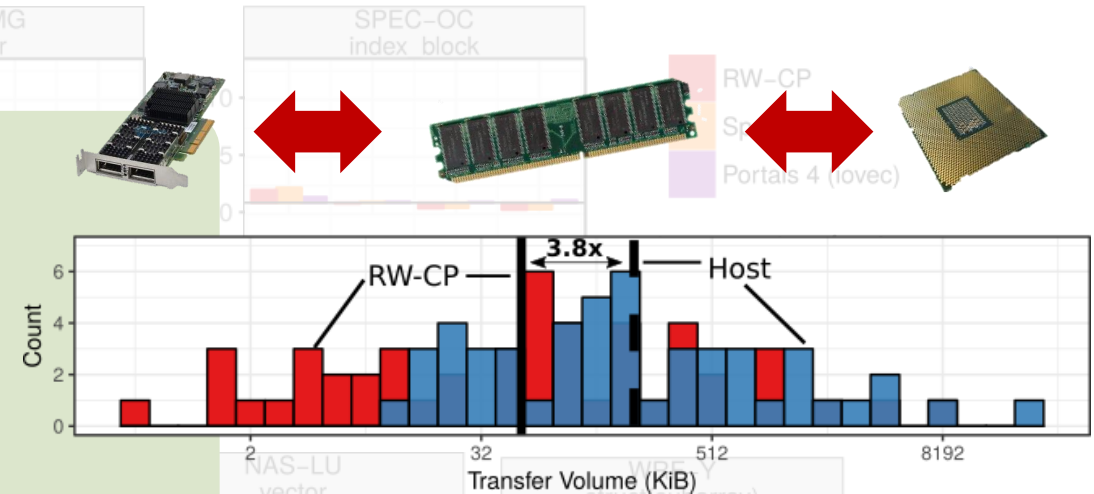


# Real Applications DDTs

## Checkpointing Overhead Data Movement

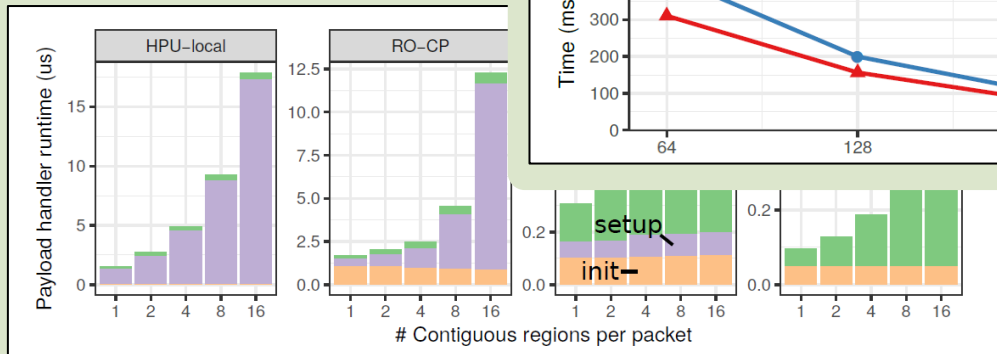


75% of the analyzed DDTs amortized after 4 reuses

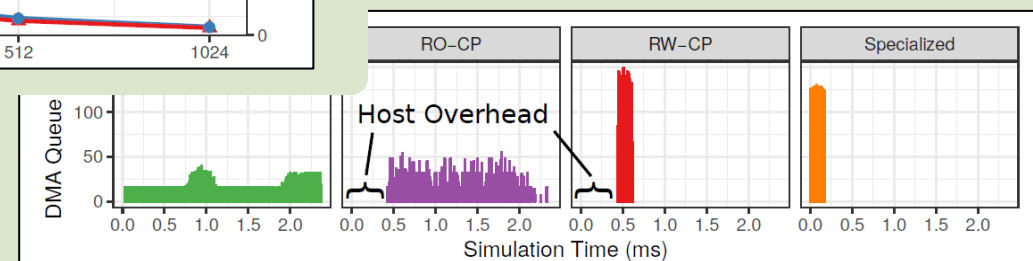
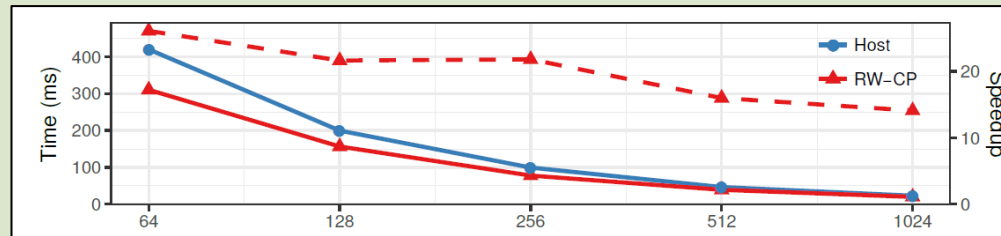


Up to 3.8x less moved data volume

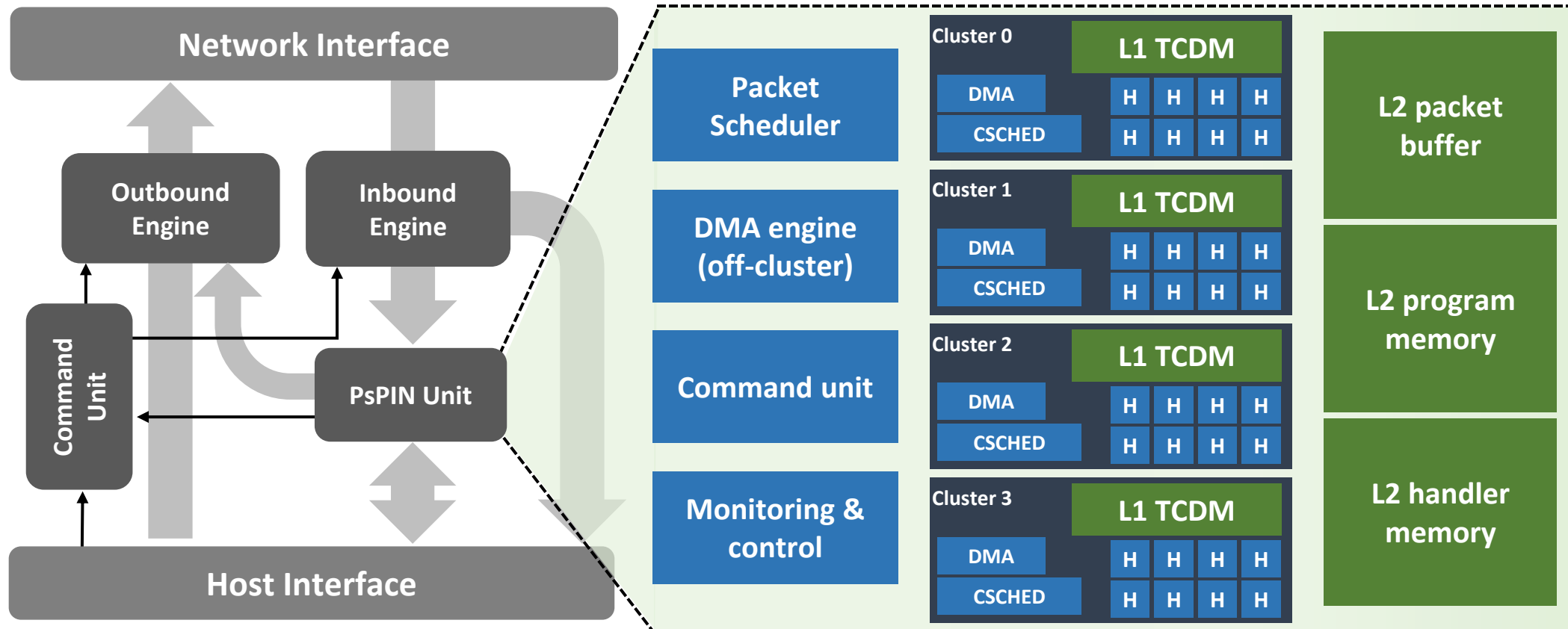
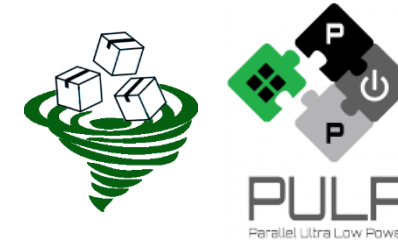
## Handler Analysis



## Full app speedup (FFT2D)



# PsPIN hardware implementation: sPIN on PULP



# Circuit Complexity and Power Estimations

- Processor synthesized in GlobalFoundries 22nm fully depleted silicon on insulator (FDSOI) technology

- Timing: 1 GHz

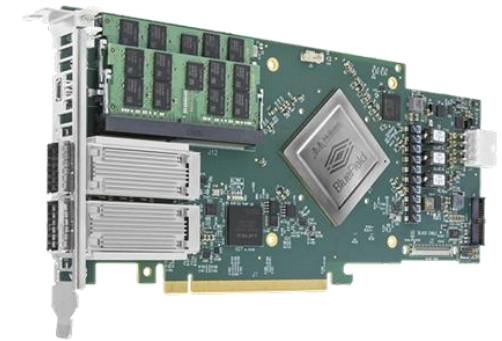
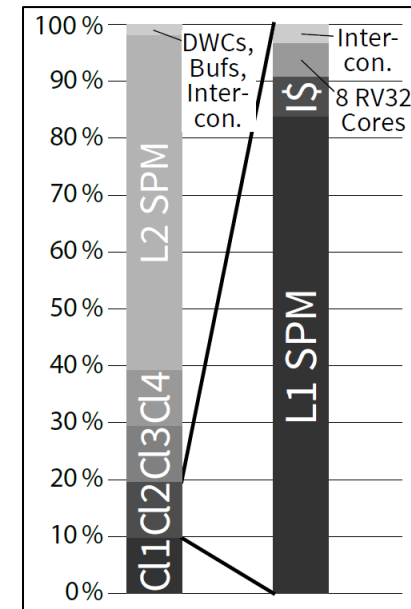
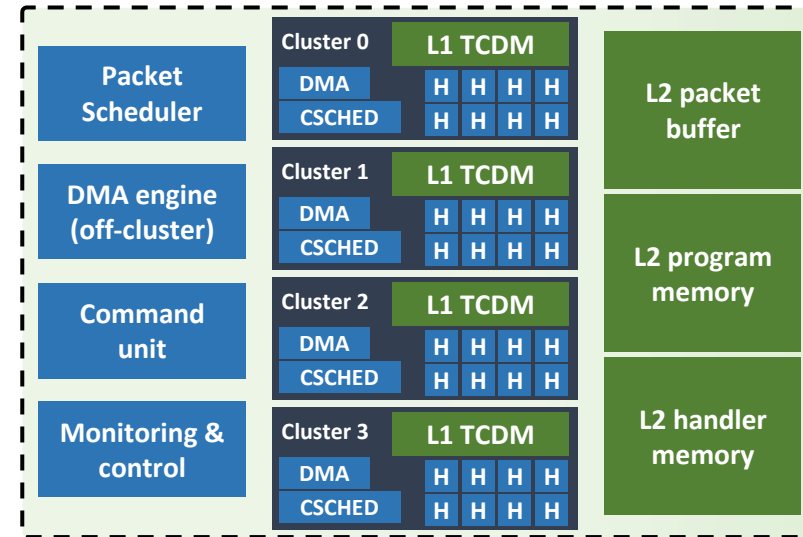
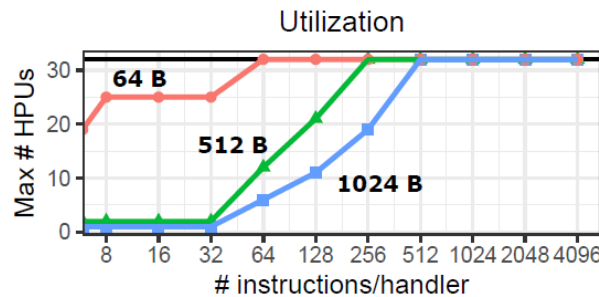
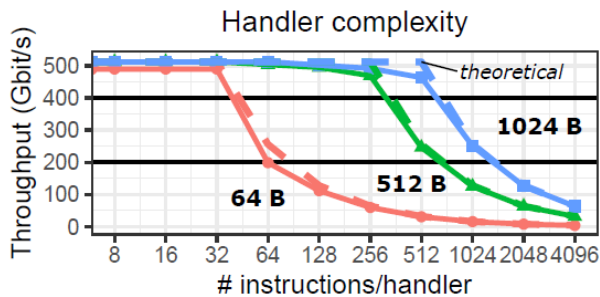
- Accelerator complexity: ~95 MGE

- 18.5 mm<sup>2</sup> area (assuming layout density 85%)
- Mellanox BlueField: 16 A72 64bit cores

*Estimated area: 51 mm<sup>2</sup>*

We could have up to 64 cores and 18 MiB of memory for that area.

- Power consumption (100% toggle rate): 6 W (not including I/O and PHY power).





# Why choosing PULP for sPIN?

## Architectures:

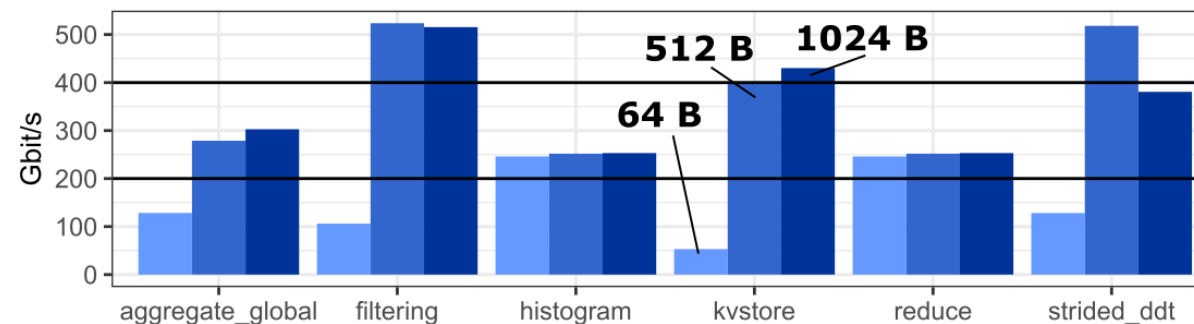
- zynq:** ARM Cortex-A53, 64-bit, 2-way superscalar, 1.2 GHz
- aolt:** Intel Skylake Gold 6154, 64-bit, out-of-order execution, 3 GHz
- PsPIN/RISCV:** RISC-V based, 32-bit, in-order, 1 GHz

Arch.	Tech.	Die area	PEs	Memory	Area/PE	Area/PE (scaled)
aolt	14 nm	485 mm <sup>2</sup> [4]	18	43.3 MiB	17.978 mm <sup>2</sup>	35.956 mm <sup>2</sup>
zynq	16 nm	3.27 mm <sup>2</sup> [3]	4	1.125 MiB	0.876 mm <sup>2</sup>	1.752 mm <sup>2</sup>
<b>PsPIN</b>	22 nm	18.5 mm <sup>2</sup>	32	12 MiB	0.578 mm <sup>2</sup>	0.578 mm <sup>2</sup>

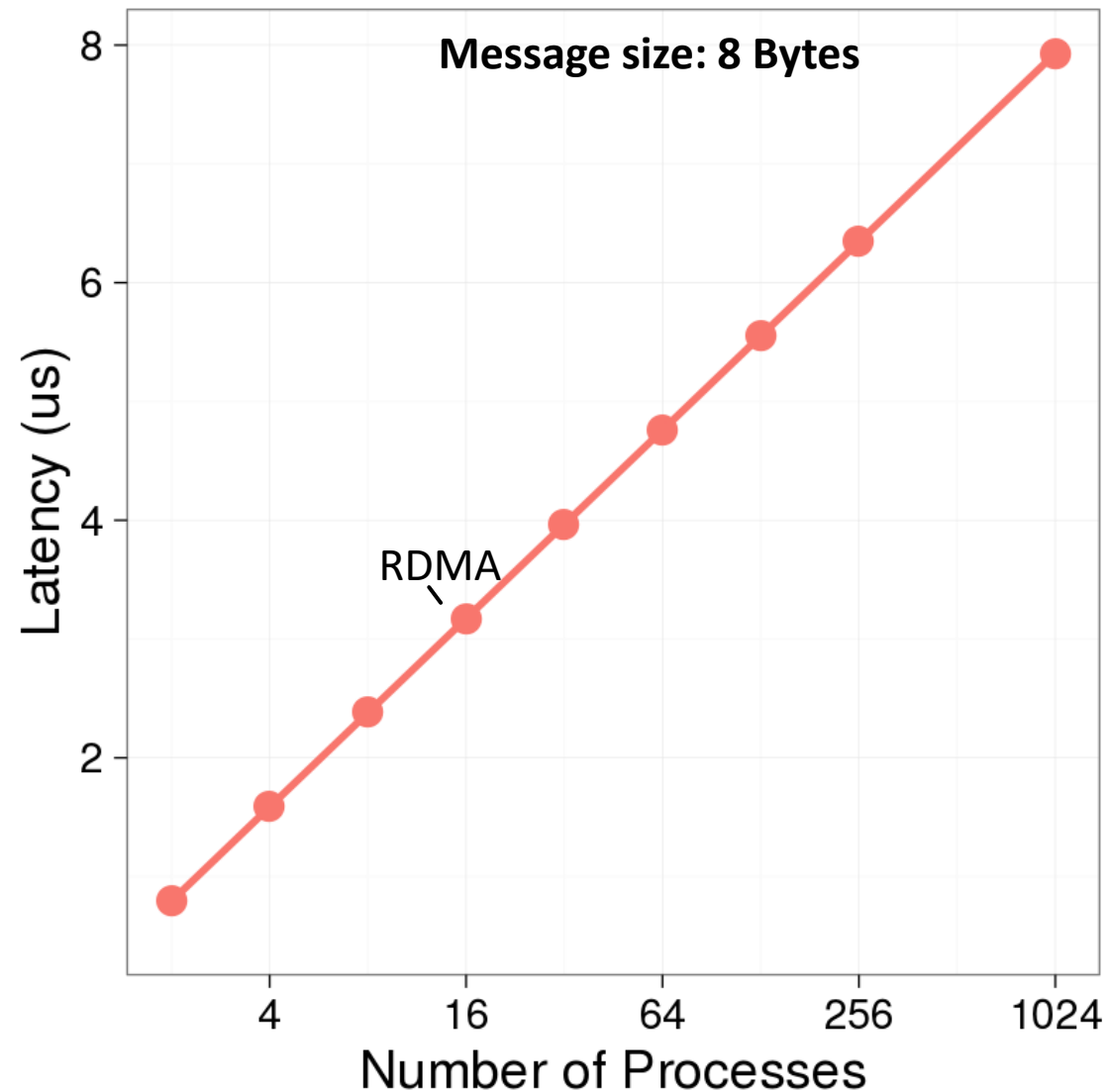
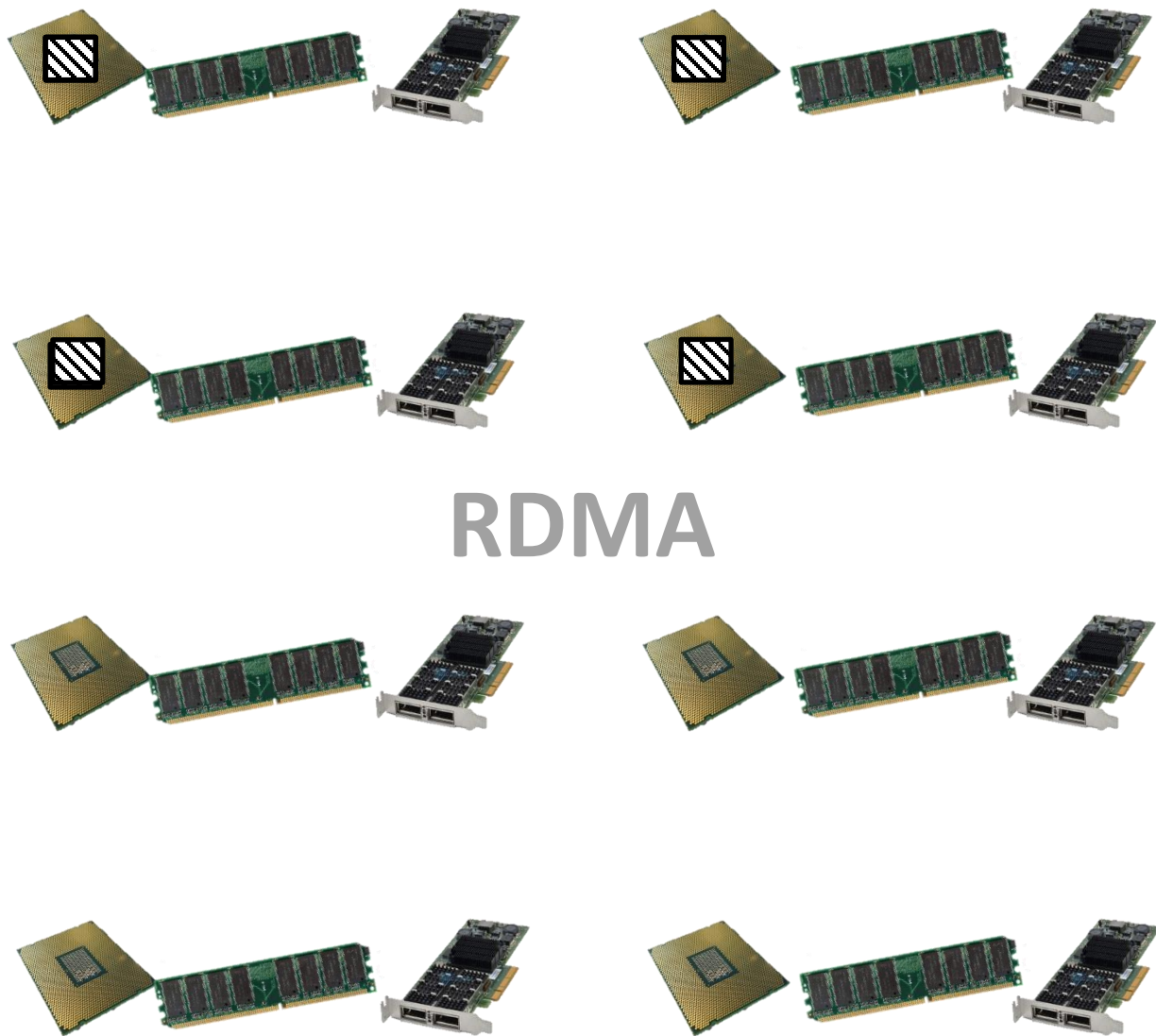
## Use cases:

- Data reduction*
- Single message aggregation*
- Packet filtering/rewriting*
- KV store cache*
- Strided datatypes*
- Histogram*

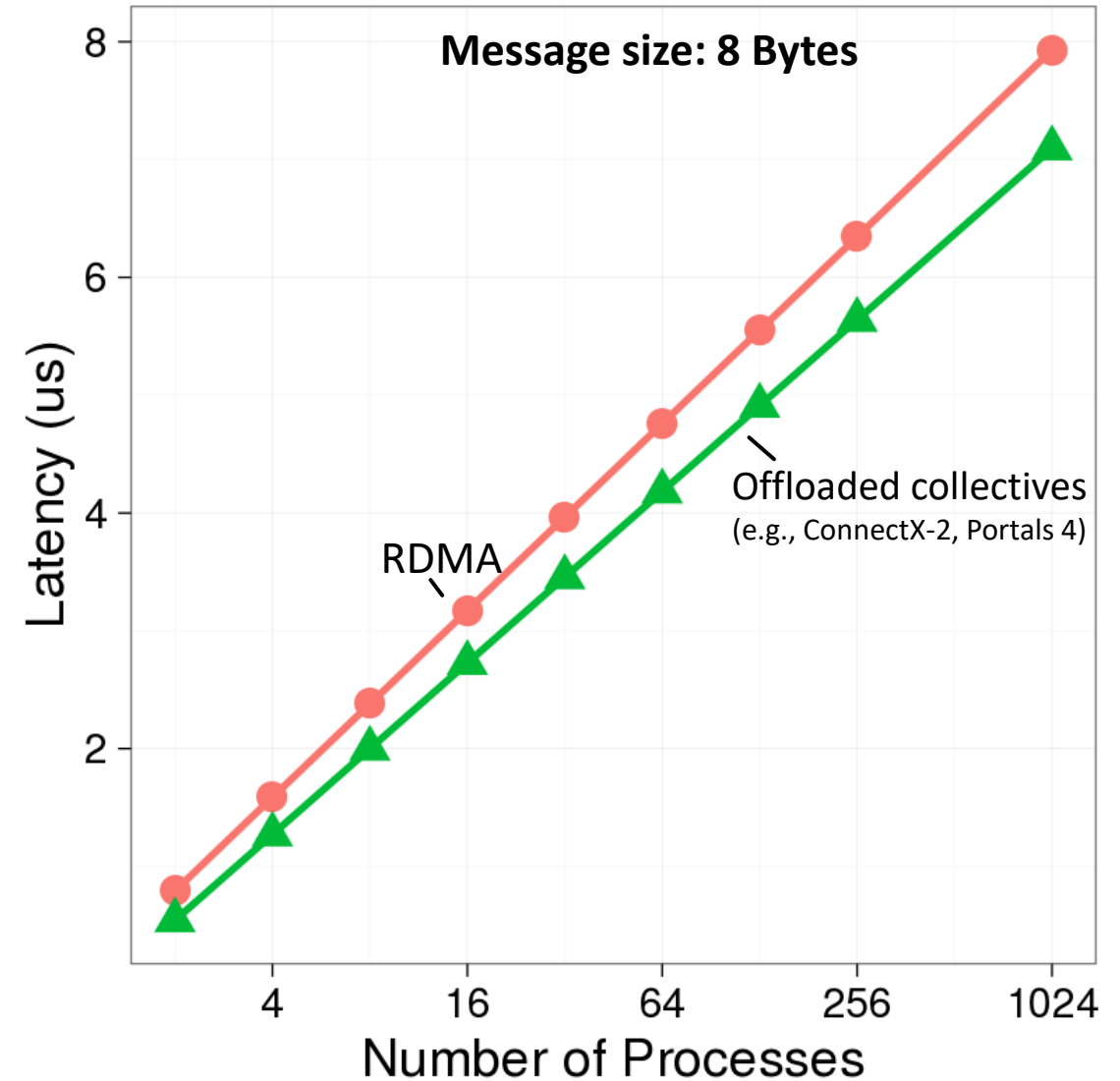
Actual throughput on PsPIN:



# Illustrating broadcast acceleration with sPIN



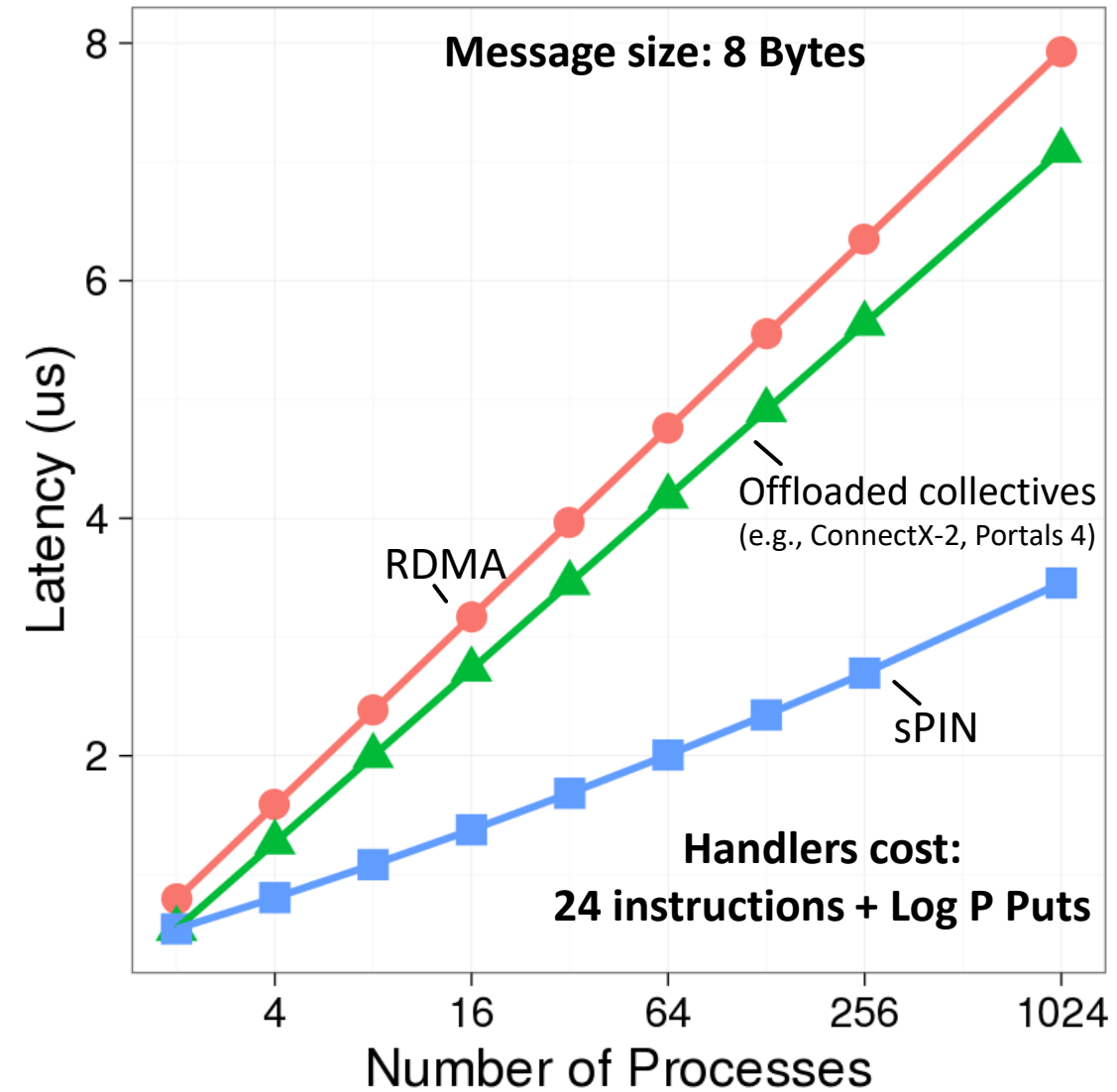
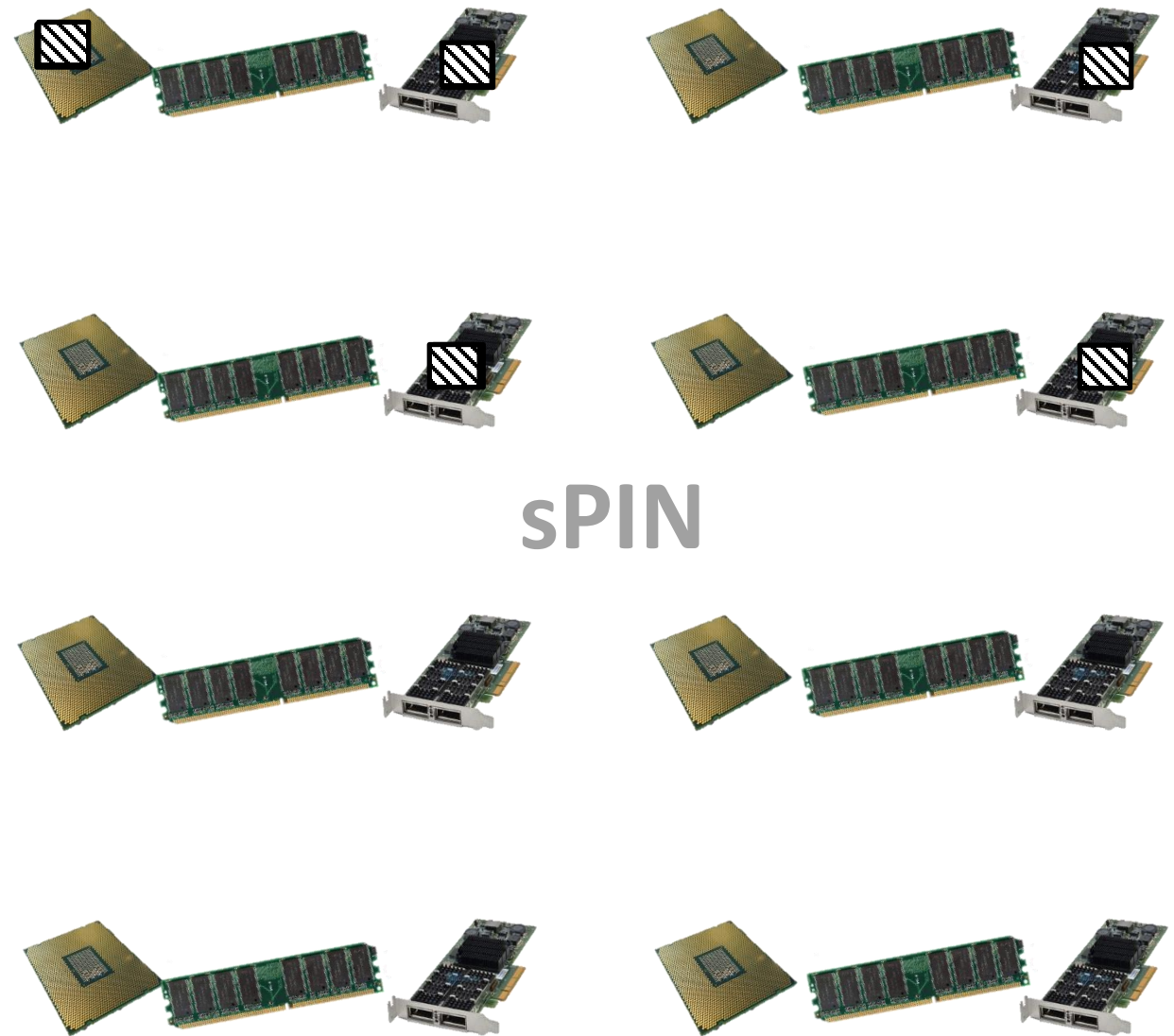
# Illustrating broadcast acceleration with sPIN



Underwood, K.D., et al., Enabling flexible collective communication offload with triggered operations. *HOTI'11*

Liu, J., et al., High performance RDMA-based MPI implementation over InfiniBand. *International Journal of Parallel Programming* 2004

# Illustrating broadcast acceleration with sPIN



Underwood, K.D., et al., Enabling flexible collective communication offload with triggered operations. *HOTI'11*


Liu, J., et al., High performance RDMA-based MPI implementation over InfiniBand. *International Journal of Parallel Programming* 2004

## Further results and use-cases

# Further results and use-cases

SPCL ETH zürich


### Use Case 4: MPI Rendezvous Protocol



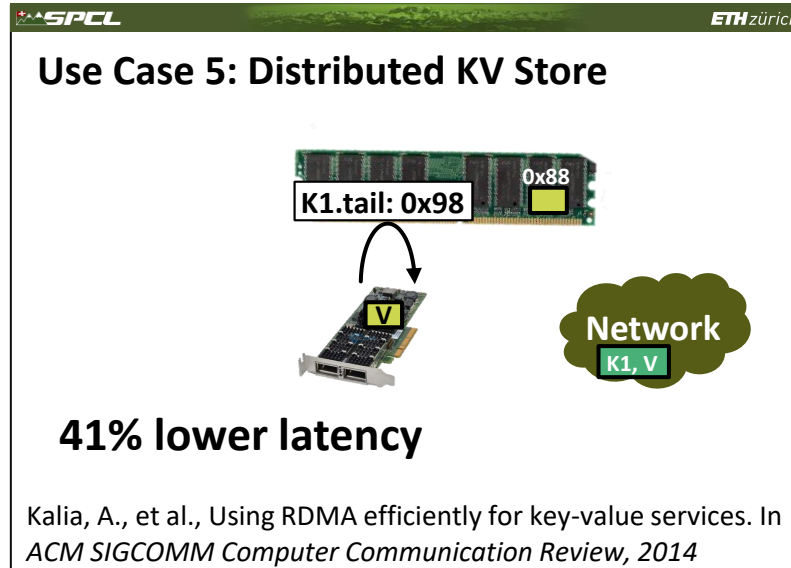
program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

# Further results and use-cases

**Use Case 4: MPI Rendezvous Protocol**



program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%



# Further results and use-cases

**Use Case 4: MPI Rendezvous Protocol**

program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

**Use Case 5: Distributed KV Store**

**41% lower latency**

Kalia, A., et al., Using RDMA efficiently for key-value services. In *ACM SIGCOMM Computer Communication Review*, 2014

**Use Case 6: Conditional Read**

Discarded data: 80%

Speedup

Data Size

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017



# Further results and use-cases

**Use Case 4: MPI Rendezvous Protocol**

program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

**Use Case 5: Distributed KV Store**

**41% lower latency**

Kalia, A., et al., Using RDMA efficiently for key-value services. In *ACM SIGCOMM Computer Communication Review*, 2014

**Use Case 6: Conditional Read**

Discarded data: 80%

Speedup

Data Size

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017

**Use Case 7: Distributed Transactions**

data pkts

log pkts

Dragojević, A, et al., No compromises: distributed transactions with consistency, availability, and performance. *SOSP'15*

# Further results and use-cases

**Use Case 4: MPI Rendezvous Protocol**

program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

**Use Case 5: Distributed KV Store**

**41% lower latency**

Kalia, A., et al., Using RDMA efficiently for key-value services. In *ACM SIGCOMM Computer Communication Review*, 2014

**Use Case 6: Conditional Read**

Discarded data: 80%

Speedup vs Data Size (32B, 4KiB, 512KiB, 64MiB)

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017

**Use Case 7: Distributed Transactions**

Dragojević, A, et al., No compromises: distributed transactions with consistency, availability, and performance. *SOSP'15*

**Use Case 8: FT Broadcast**

Bosilca, G., et al., Failure Detection and Propagation in HPC systems. *SC'16*

# Further results and use-cases

### Use Case 4: MPI Rendezvous Protocol

program	p	msgs	ovhd	ovhd	red
MILC	64	5.7M	5.5%	1.9%	65%
POP	64	772M	3.1%	2.4%	22%
coMD	72	5.3M	6.1%	2.4%	60%
coMD	360	28.1M	6.5%	2.8%	58%
Cloverleaf	72	2.7M	5.2%	2.4%	53%
Cloverleaf	360	15.3M	5.6%	3.2%	42%

### Use Case 5: Distributed KV Store

## The Next 700 sPIN use-cases

... just think about sPIN graph kernels ...

**41% lower latency**

Kalia, A., et al., Using sPIN for distributed KV Store services. In ACM SIGCOMM Conference on Computer Communications Security, 2017

### Use Case 6: Conditional Read

Barthels, C., et al., Designing Databases for Future High-Performance Networks. *IEEE Data Eng. Bulletin*, 2017

### Use Case 7: Distributed Transactions

Dragojević, A, et al., No compromises: distributed transactions with consistency, availability, and performance. SOSP'15

### Use Case 8: Distributed Consensus

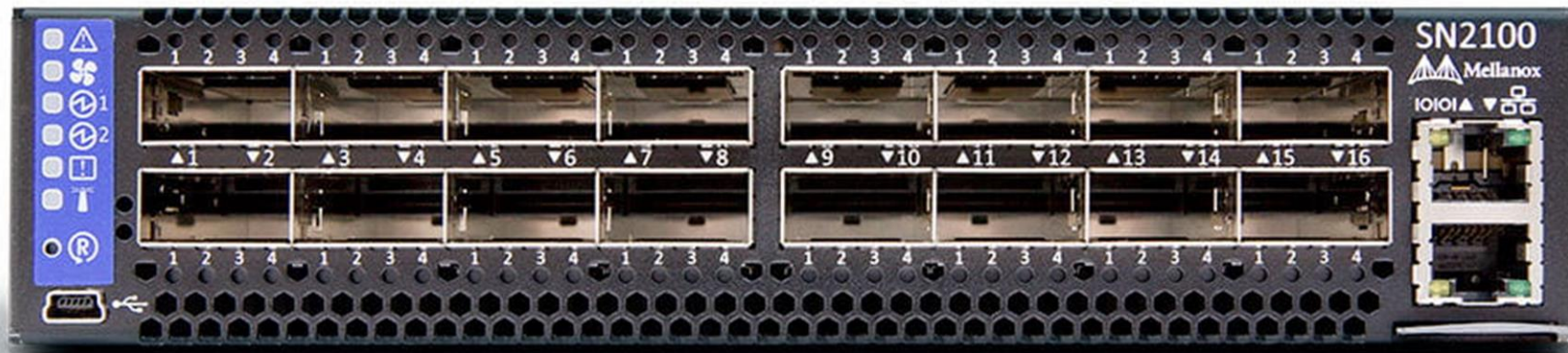
Bosilca, G., et al., Failure Detection and Propagation in HPC systems. SC'16

### Use Case 9: Distributed Consensus

István, Z., et al., Consensus in a Box: Inexpensive Coordination in Hardware. NSDI'16

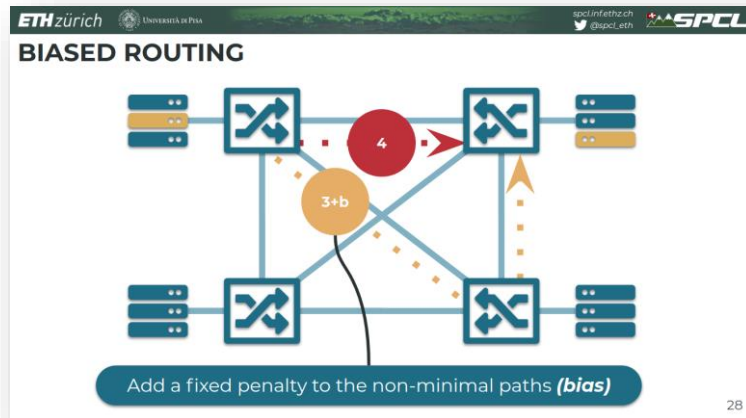
## Next step - pushing sPIN into network switches?

- **Needs to be carefully vetted!**
  - Can we achieve our goals with P4?
- **What else needs to be fixed before we go into the network?**
- **We chose to investigate network noise first**
  - Not enough time here but let me give you a brief overview.



# Network noise analysis and mitigation

## Analysis of the impact of adaptive routing on network noise



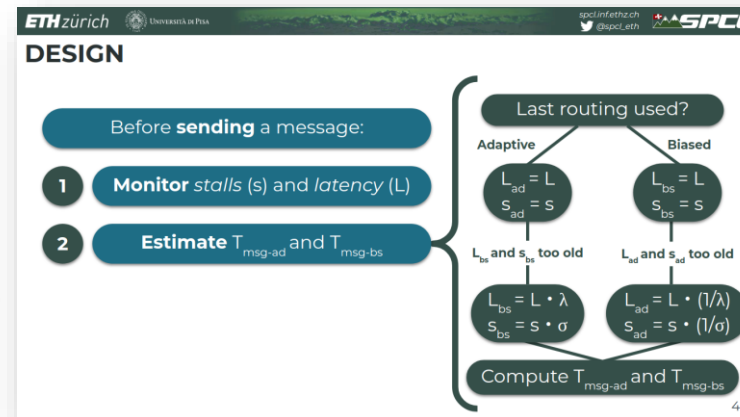
## Mitigating Network Noise on Dragonfly Networks through Application-Aware Routing

Daniele De Sensi  
desensi@di.unipi.it  
University of Pisa  
Pisa, Italy  
ETH Zurich  
Zurich, Switzerland

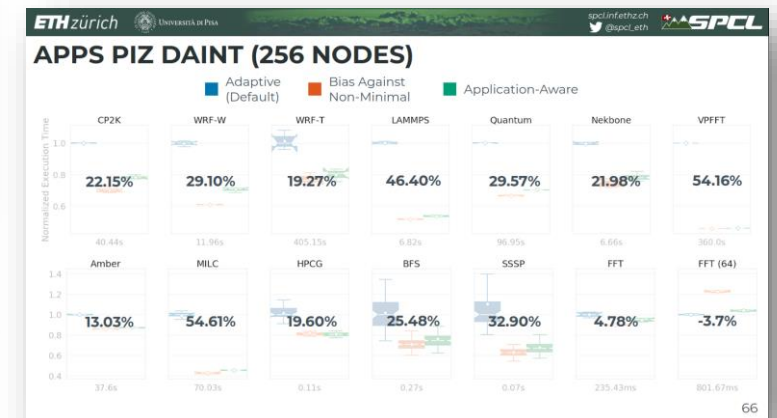
Salvatore Di Girolamo  
salvatore.digirolamo@inf.ethz.ch  
ETH Zurich  
Zurich, Switzerland

Torsten Hoefer  
htor@inf.ethz.ch  
ETH Zurich  
Zurich, Switzerland

## Design and implementation of a transparent solution



## Improvements up to 55% on real applications



# Slingshot the Exascale Interconnect

## Description of the main features of the interconnect

- All HPC traffic layered over **RoCEv2**
- Efficient **software stack**
- High-Radix **Switches**
- Low-Diameter **Topology**
- Congestion Control**
- Adaptive Routing**
- Quality of Service**

## An In-Depth Analysis of the Slingshot Interconnect

Daniele De Sensi  
Department of Computer Science  
ETH Zurich  
ddesensi@ethz.ch

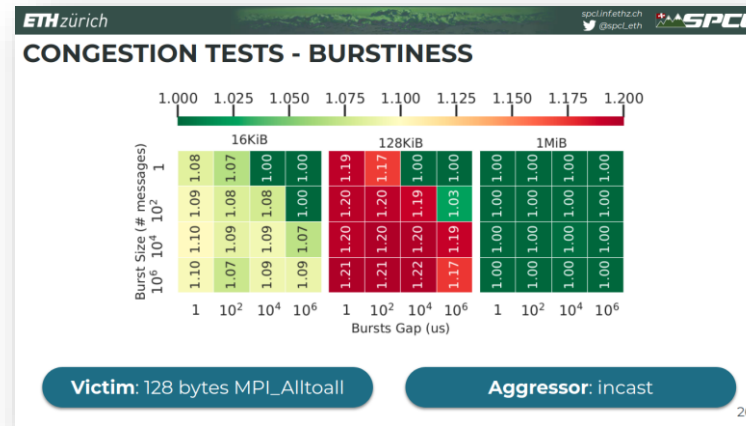
Salvatore Di Girolamo  
Department of Computer Science  
ETH Zurich  
salvatore.digirolamo@inf.ethz.ch

Kim H. McMahon  
Hewlett Packard Enterprise  
kim.mcmahon@hpe.com

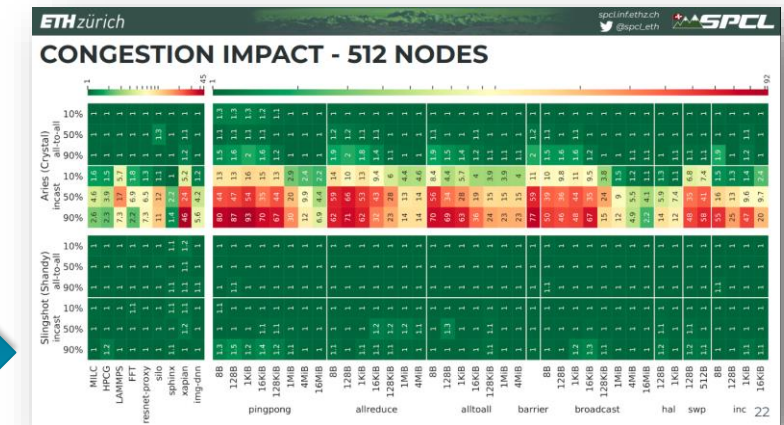
Duncan Roweth  
Hewlett Packard Enterprise  
duncan.roweth@hpe.com

Torsten Hoefler  
Department of Computer Science  
ETH Zurich  
torsten.hoefler@inf.ethz.ch

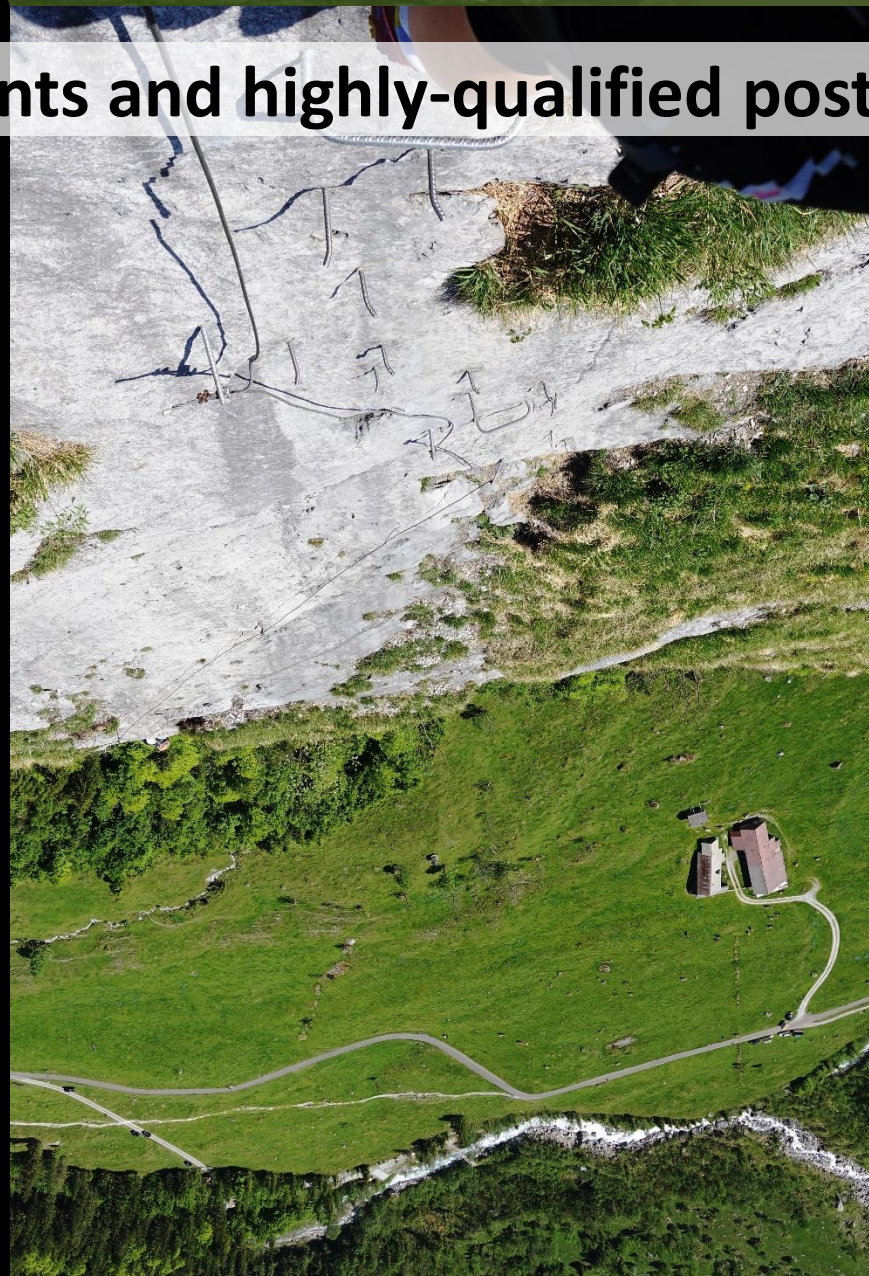
## In-depth benchmarking procedure that can be ported to other interconnect



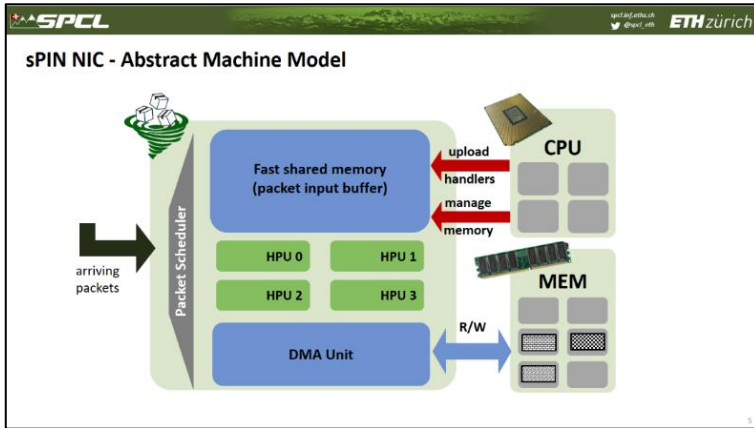
## Detailed results on performance, congestion control, and quality of service, on microbenchmarks, HPC, and DC applications



**SPCL is hiring PhD students and highly-qualified postdocs to reach new heights!**



# sPIN Streaming Processing in the Network for Network Acceleration



```

sPIN - Programming Interface

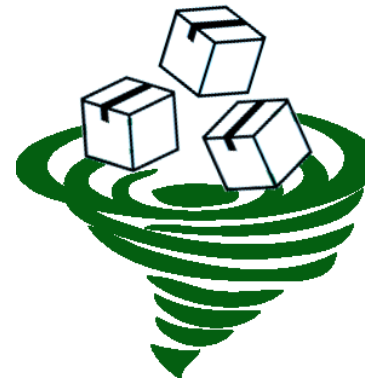
Header handler
_handler int pp_header_handler(const pti_header_t h, void *state) {
    pingpong_info_t *i = >state;
    i->source = h.source_id;
    return PROCESS_DATA; // execute payload handler to put from device
}

Payload handler
_handler int pp_payload_handler(const pti_payload_t p, void *state) {
    pingpong_info_t *i = >state;
    PtiHandlerPutFromDevice(p.base, p.length, 1, 0, i->source, 10, 0, NULL, 0);
    return SUCCESS;
}

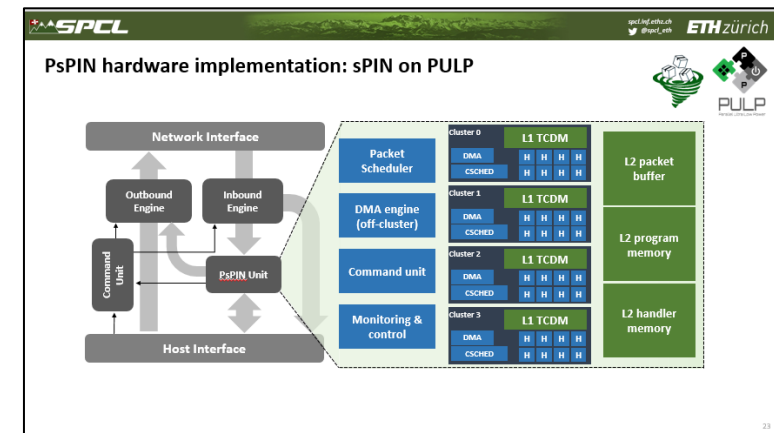
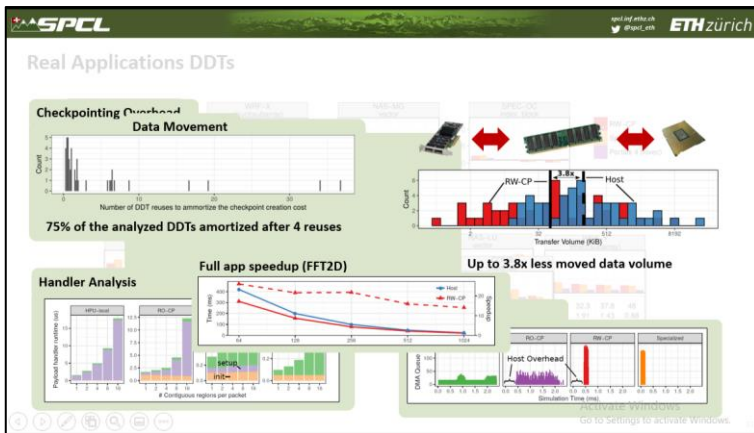
Completion handler
_handler int pp_completion_handler(int dropped_bytes,
    bool flow_control_triggered, void *state) {
    return SUCCESS;
}

connect(peer, /* ... */, &pp_header_handler, &pp_payload_handler, &pp_completion_handler);
    
```

sPIN



beyond RDMA





# Backup Slides

## But why PULP/RISC-V?

- RISC-V is an open source ISA
  - Allows and supports extensions
    - Doing this in ARM may be complex and expensive*
- PULP aims to provide high performance per Watt
  - Energy efficient
  - Provides tight control over compute and data movement schedule
  - Fits well the sPIN abstract machine model (e.g., removing cache coherency on ARM could be painful)
  - PULP is actively researched + we can leverage ISS group expertise at ETH