TECHNISCHE UNIVERSITÄT DRESDEN

**André Berthold,**
Constantin Fürst, Antonia Obersteiner, Lennart Schmidt, Dirk Habich, Wolfgang Lehner, Horst Schirmeier

Dresden University of Technology, Faculty of Computer Science, Institute of Systems Architecture
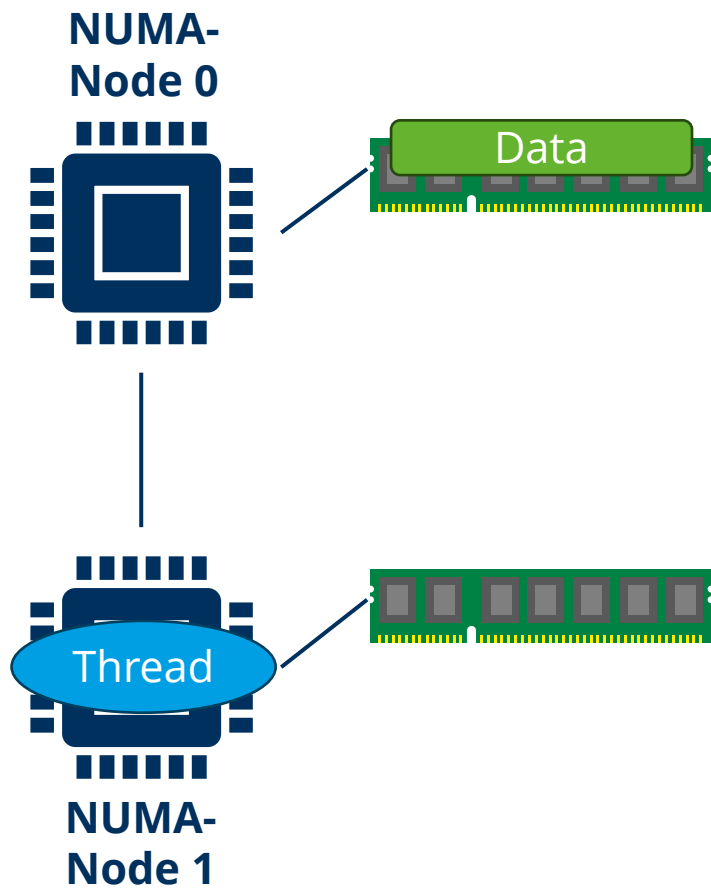
# Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing

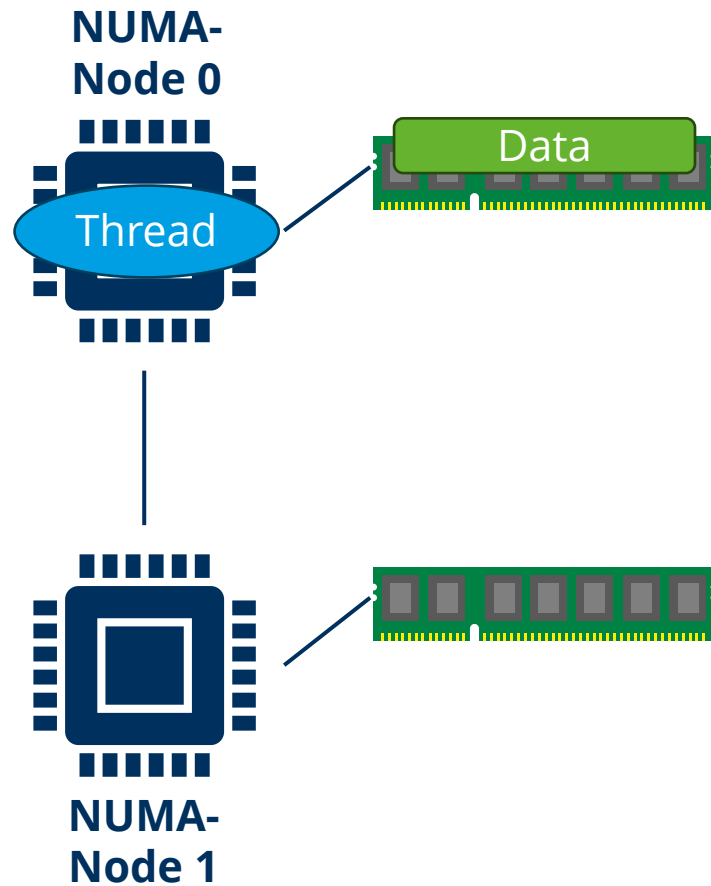Austin, Texas, DIMES 2nd Workshop on Disruptive Memory Systems // November 3, 2024
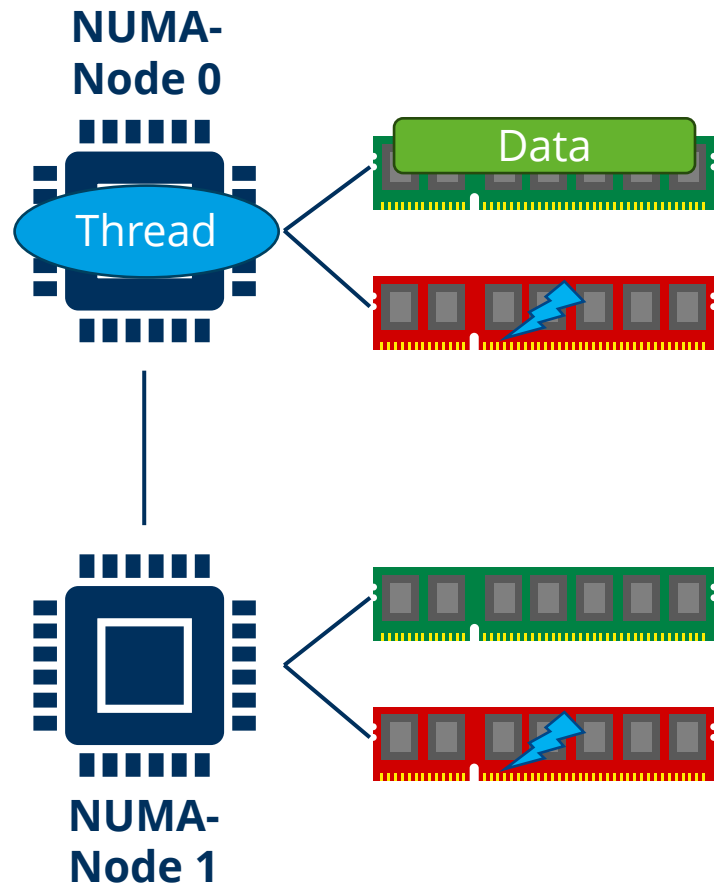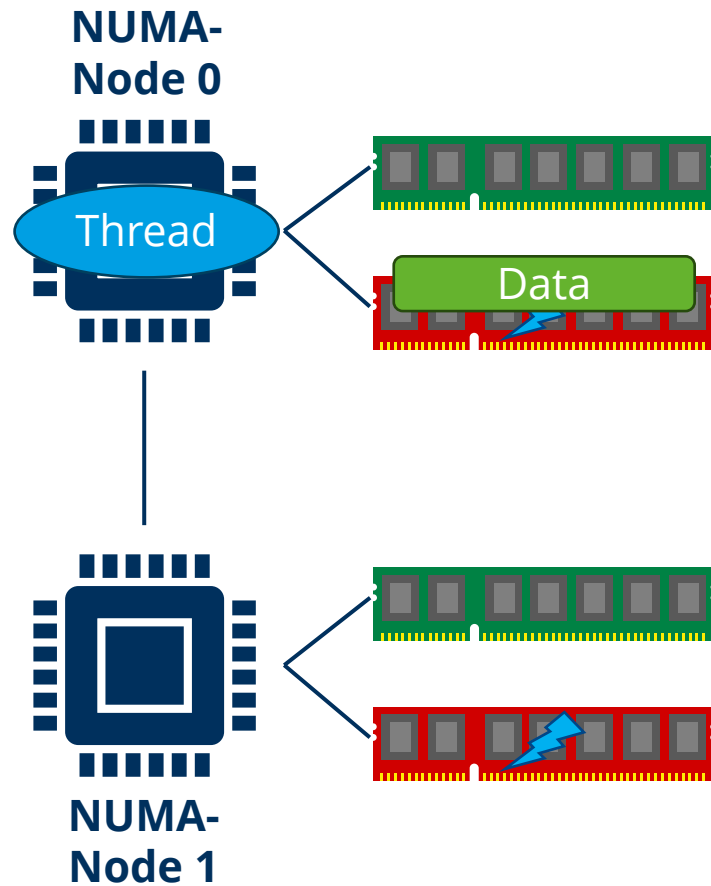
# Motivation

NUMA-
Node 0

NUMA-
Node 1

Funded by

# Motivation

# Motivation

# Motivation



Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 5

Funded by

# Motivation

Funded by

# Motivation



$$\texttt{std::memcpy(a, b, } 2^{30}\texttt{);}$$

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 7

Funded by

# Motivation

NUMA-Node 0

Thread

Data

$$\texttt{std::memcpy(a, b, } 2^{30}\texttt{);}$$

DMA-Engine

NUMA-Node 1

# Motivation



NUMA-Node 0

Thread

Data

NUMA-Node 1

```
std::memcpy(a, b, 2^30);
```

DMA-Engine

# Motivation

NUMA-Node 0

Thread

Data

NUMA-Node 1

$$\texttt{std::memcpy(a, b, } 2^{30}\texttt{);}$$

DMA-Engine

Sapphire Rapids Microarchitecture

Processor

TECHNISCHE UNIVERSITÄT DRESDEN

Funded by

DFG

DRESDEN concept

# Motivation



$$\texttt{std::memcpy(a, b, } 2^{30}\texttt{);}$$

NUMA-Node 0

Thread

Data

NUMA-Node 1

DMA-Engine

Sapphire Rapids Microarchitecture

DRAM

Processor

HBM

# Motivation



NUMA-Node 0

Thread

Data

NUMA-Node 1

```
std::memcpy(a, b, 2^30);
```

DMA-Engine

Sapphire Rapids Microarchitecture

DRAM

Processor

DSA

HBM

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 12

TECHNISCHE UNIVERSITÄT DRESDEN

Funded by

DFG

DRESDEN concept

# DSA - Intel Data Streaming Accelerator

**The DSA** supports:
- Memory Copy
- Memory Fill
- Memory Compare
  ⋮



Reese Kuper, et al. 2024. A Quantitative Analysis and Guidelines of Data Streaming Accelerator in Modern Intel Xeon Scalable Processors. In *ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*, 37–54.

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 13

Funded by
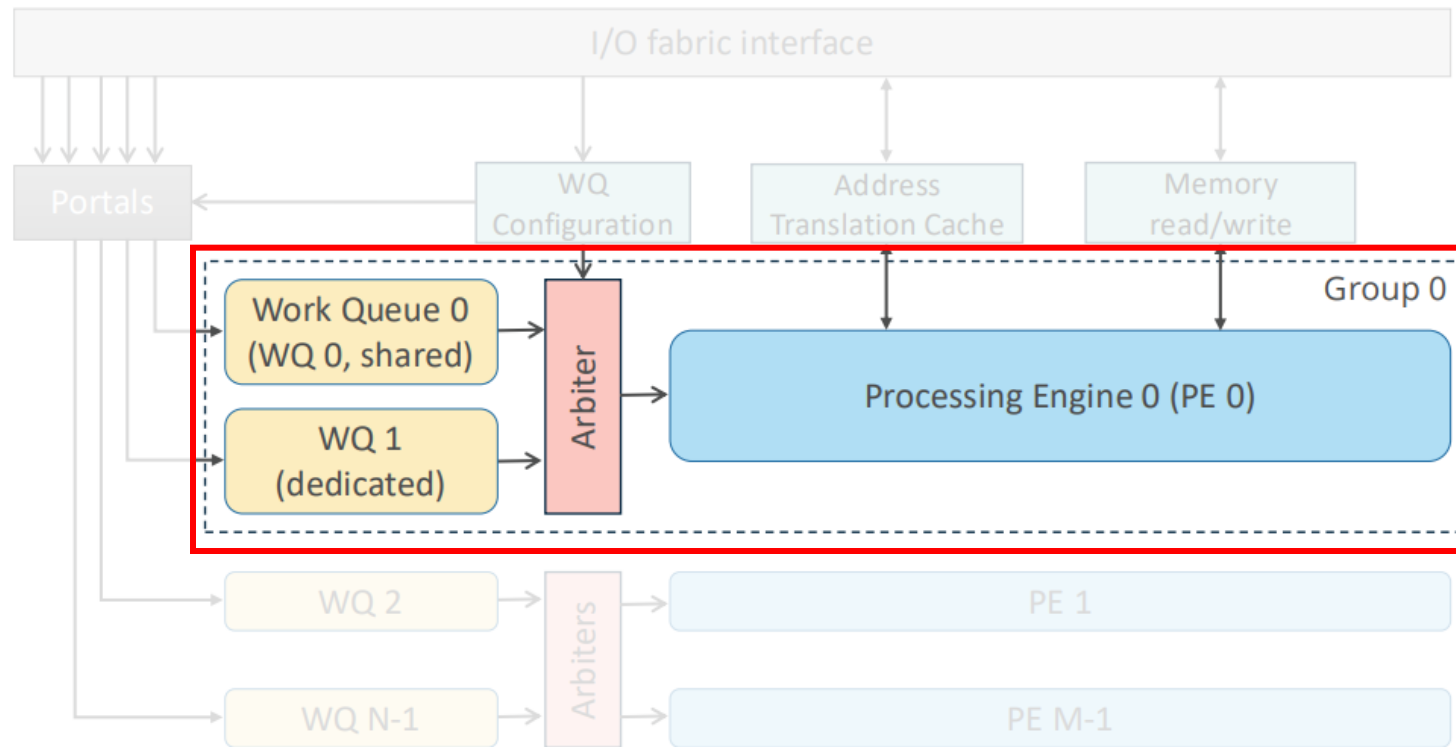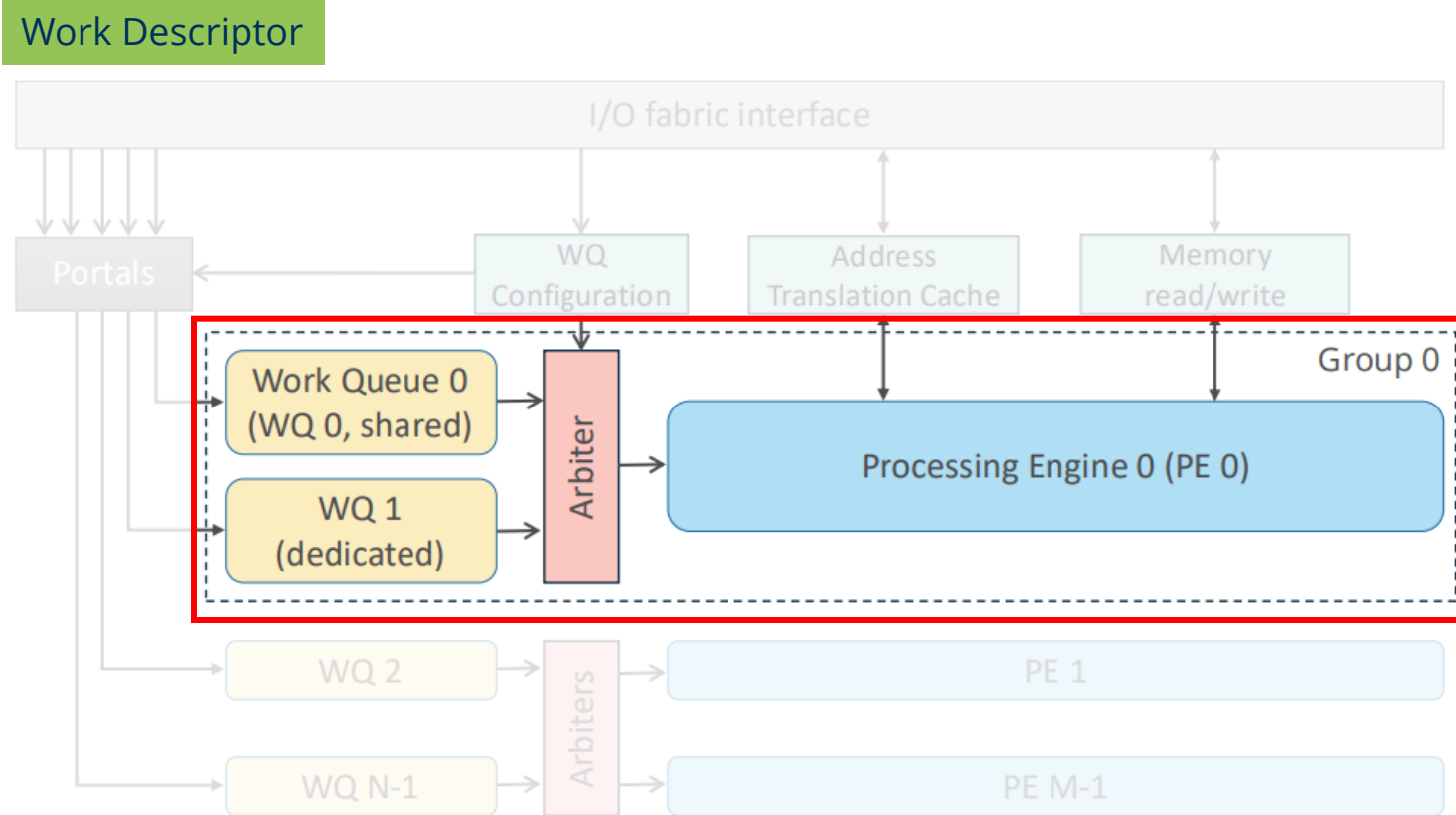
# DSA - Intel Data Streaming Accelerator

**The DSA** supports:
- Memory Copy
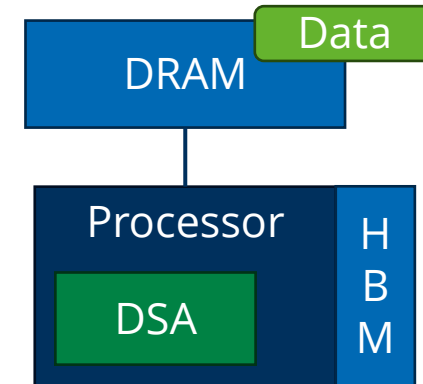- Memory Fill
- Memory Compare

⋮



Reese Kuper, et al. 2024. A Quantitative Analysis and Guidelines of Data Streaming Accelerator in Modern Intel Xeon Scalable Processors. In *ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*, 37–54.

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 14

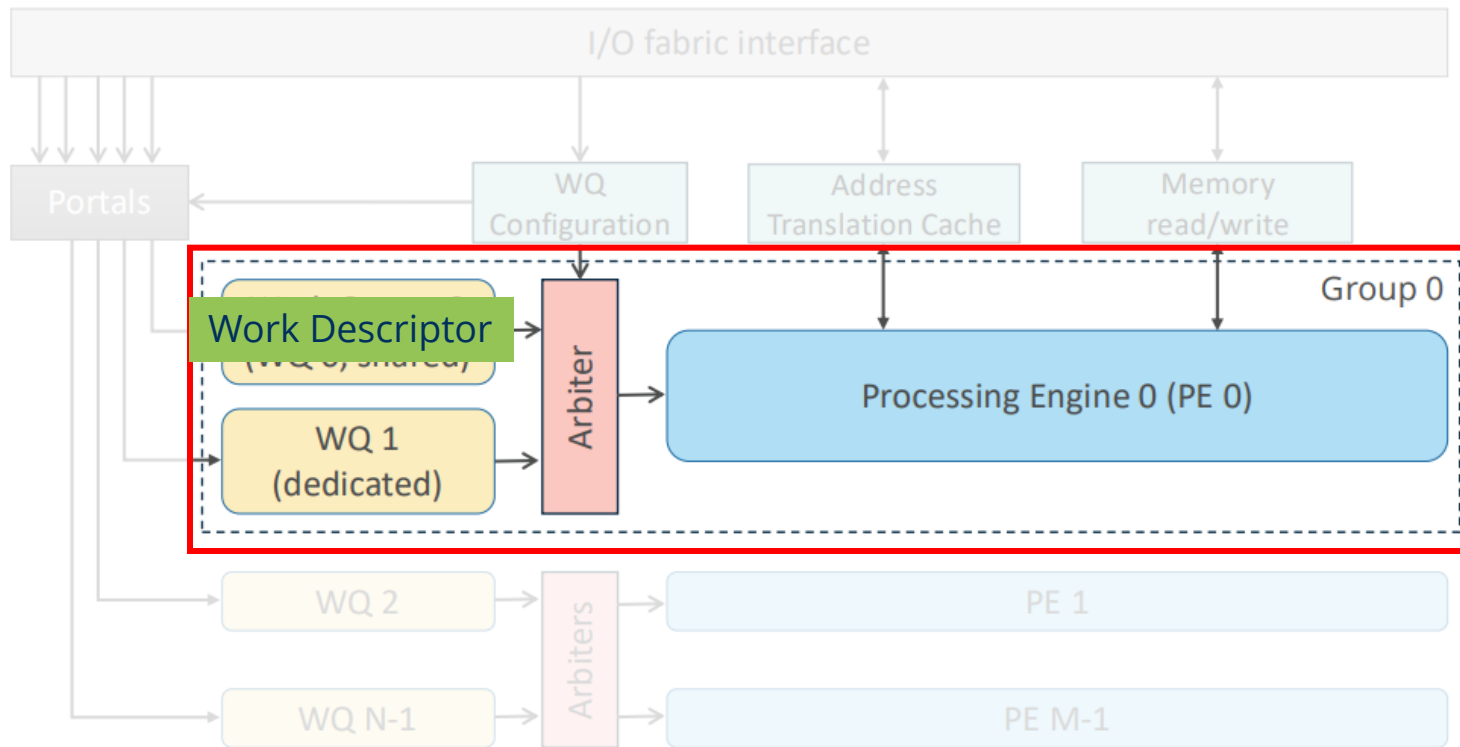Funded by

# DSA - Intel Data Streaming Accelerator

**The DSA** supports:
- Memory Copy
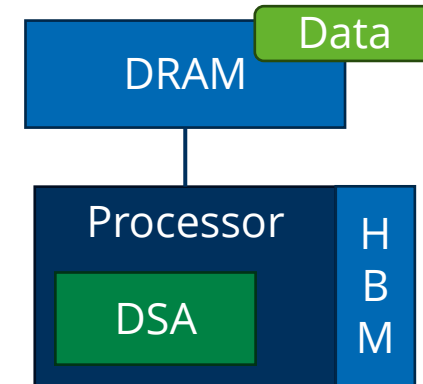- Memory Fill
- Memory Compare
⋮



Reese Kuper, et al. 2024. A Quantitative Analysis and Guidelines of Data Streaming Accelerator in Modern Intel Xeon Scalable Processors. In *ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*, 37–54.

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 15

Funded by

# DSA - Intel Data Streaming Accelerator

**The DSA** supports:
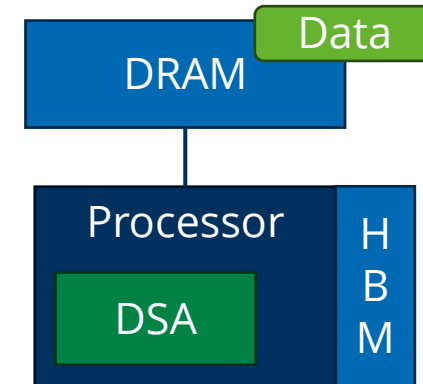- Memory Copy
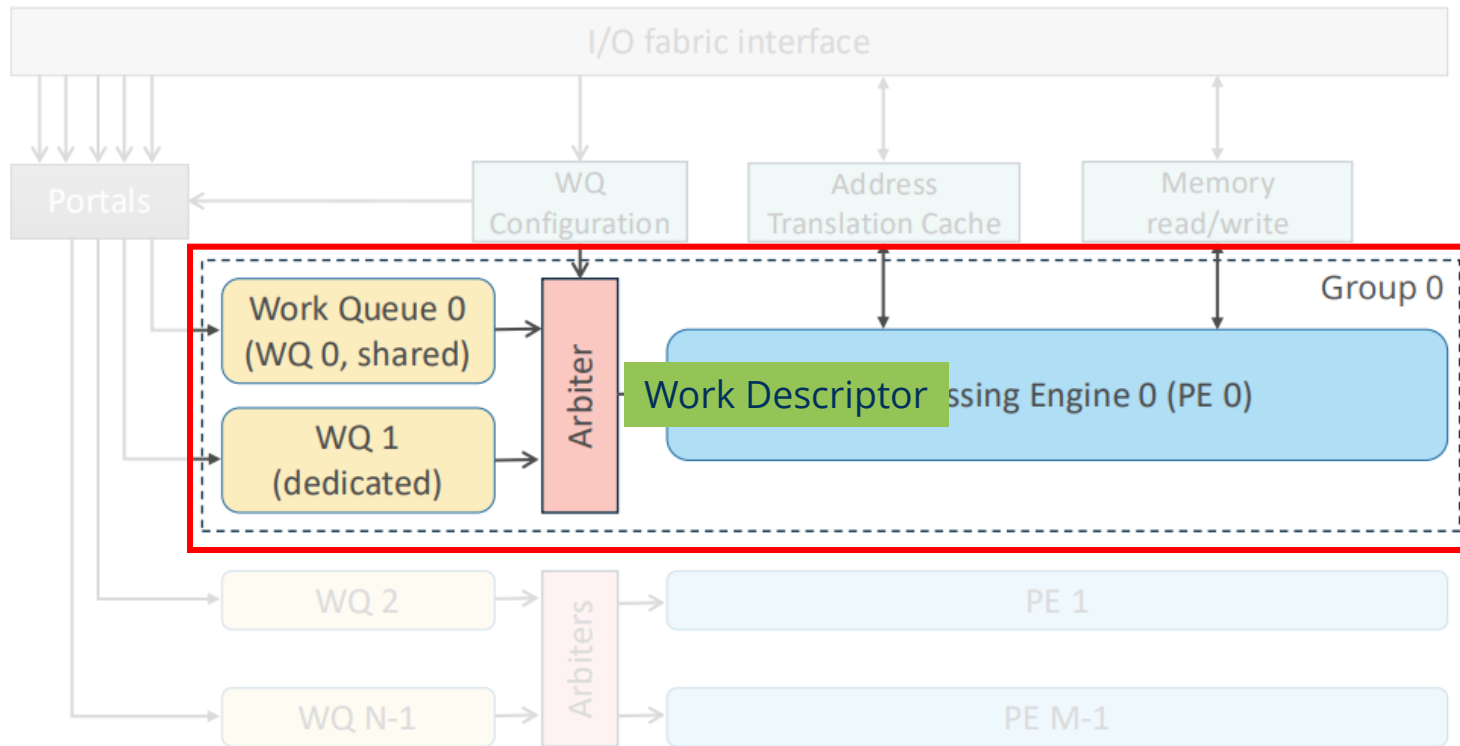- Memory Fill
- Memory Compare
  ⋮



Reese Kuper, et al. 2024. A Quantitative Analysis and Guidelines of Data Streaming Accelerator in Modern Intel Xeon Scalable Processors. In *ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*, 37–54.

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 16

Funded by

# DSA - Intel Data Streaming Accelerator

**The DSA** supports:
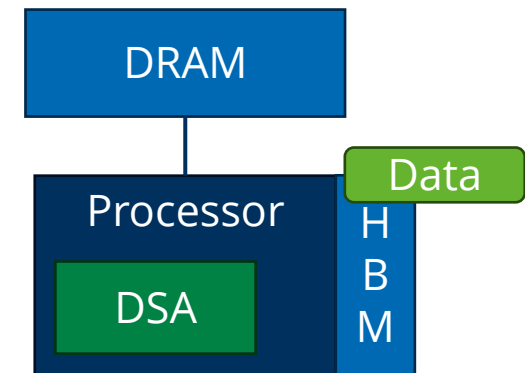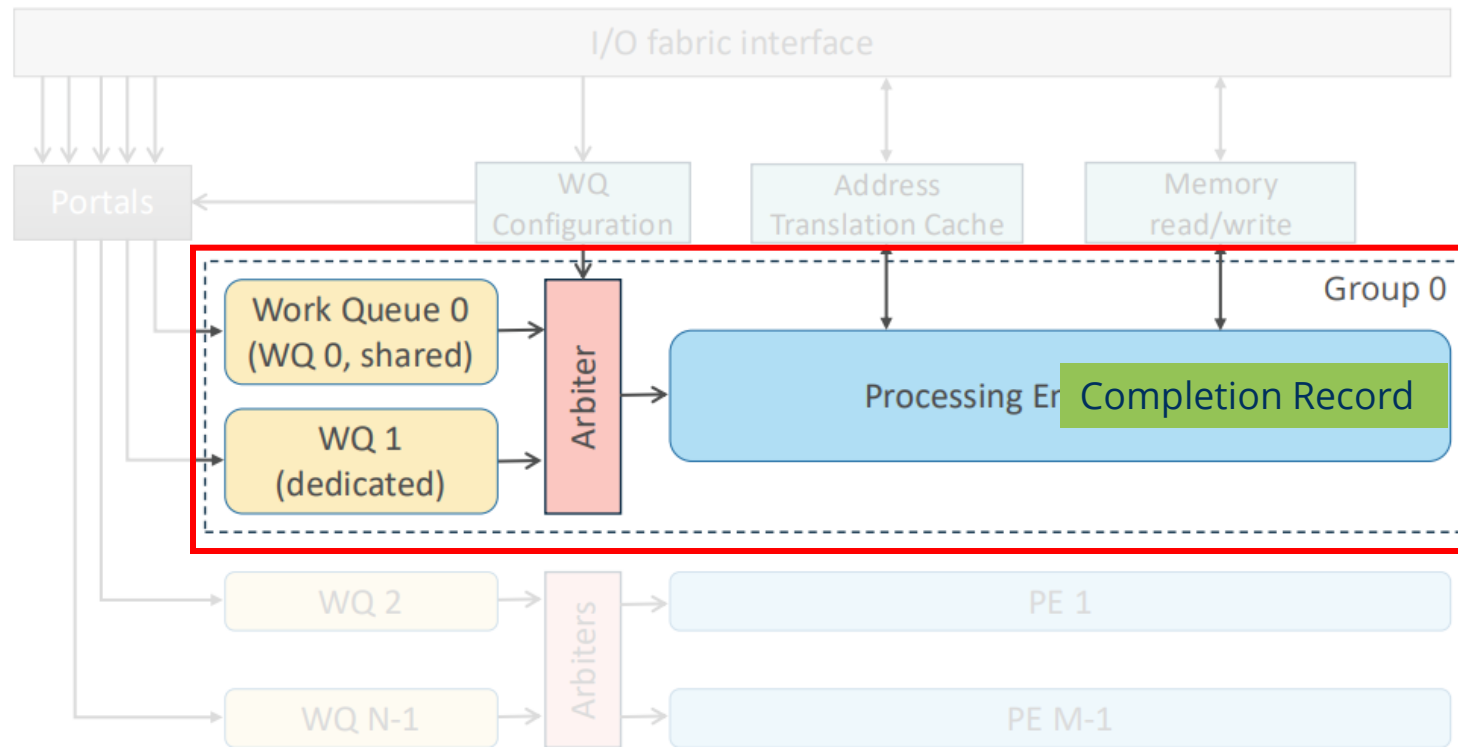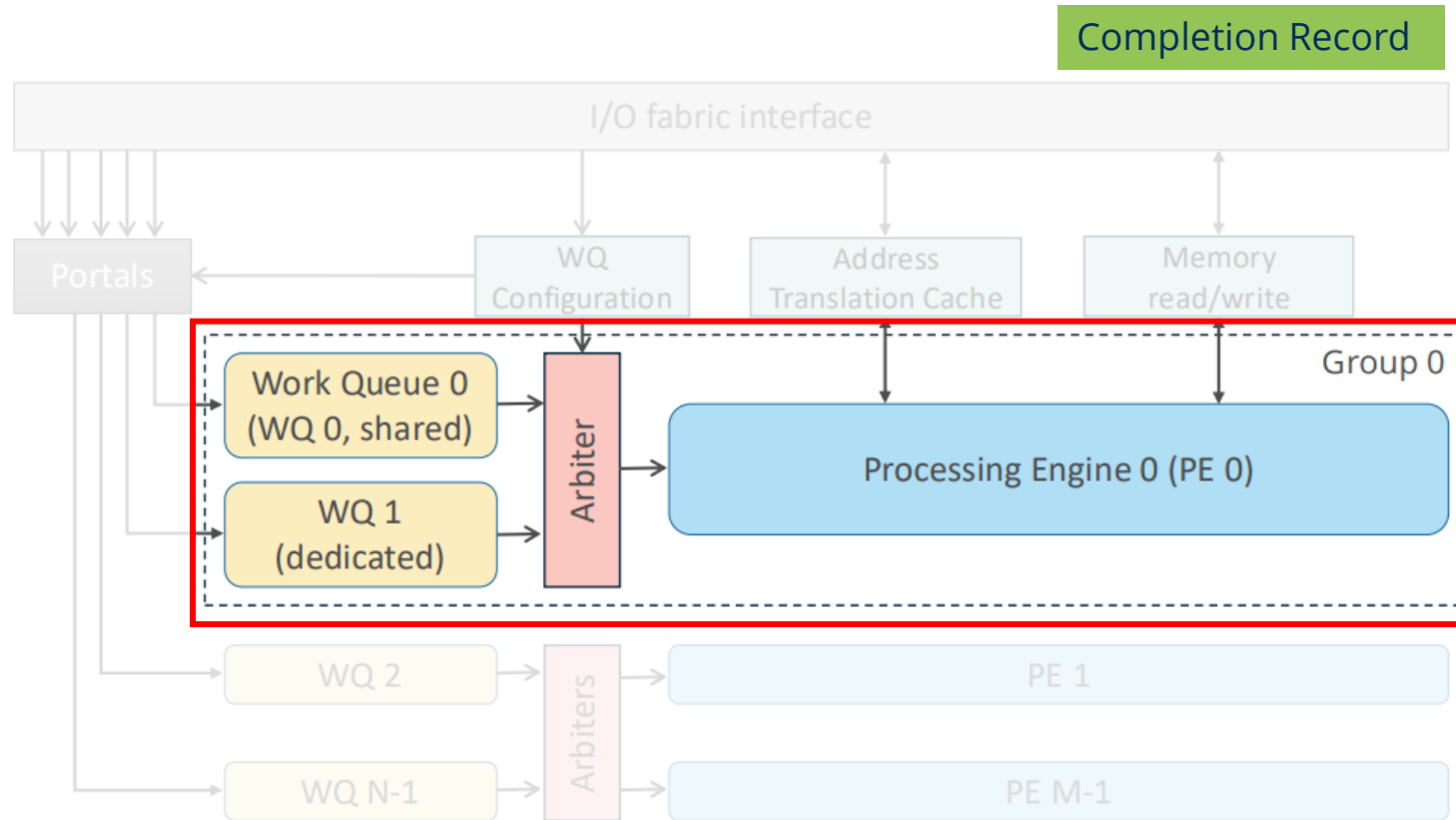- Memory Copy
- Memory Fill
- Memory Compare
  ⋮



Reese Kuper, et al. 2024. A Quantitative Analysis and Guidelines of Data Streaming Accelerator in Modern Intel Xeon Scalable Processors. In *ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*, 37–54.

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 17

Funded by

# DSA - Intel Data Streaming Accelerator



**The DSA** supports:
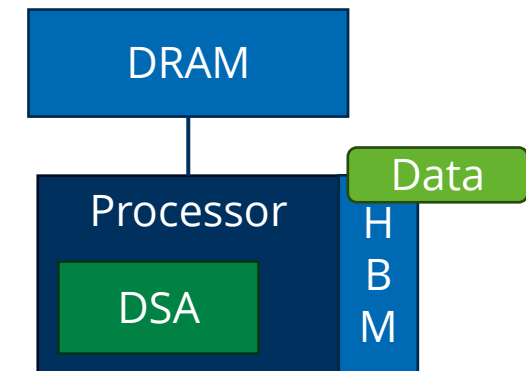- Memory Copy
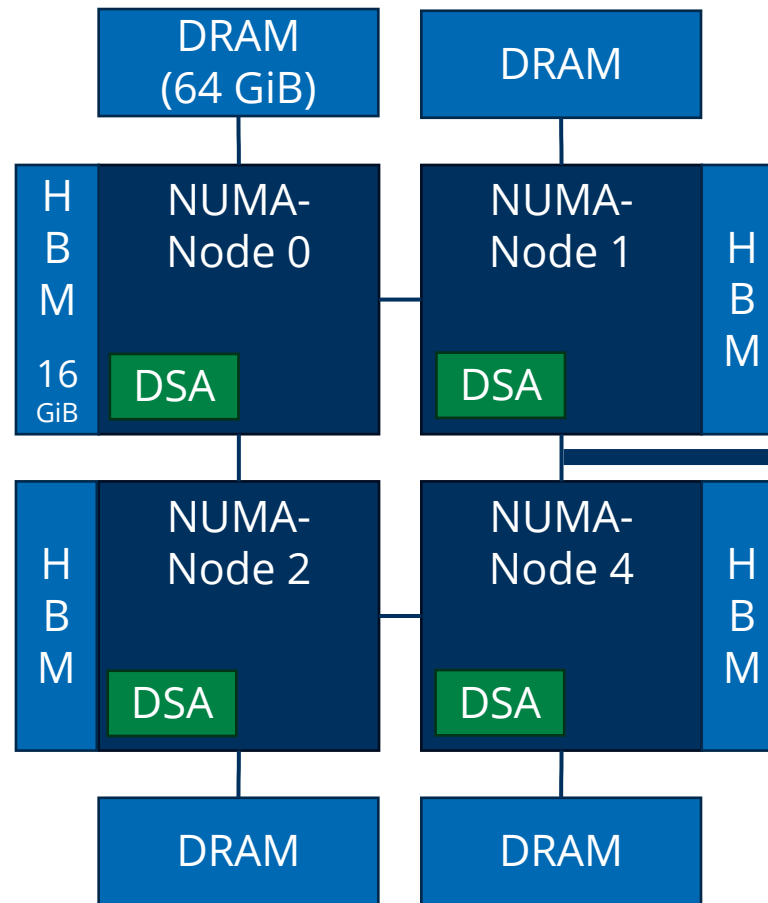- Memory Fill
- Memory Compare
  ⋮

Reese Kuper, et al. 2024. A Quantitative Analysis and Guidelines of Data Streaming Accelerator in Modern Intel Xeon Scalable Processors. In *ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS '24)*, 37–54.
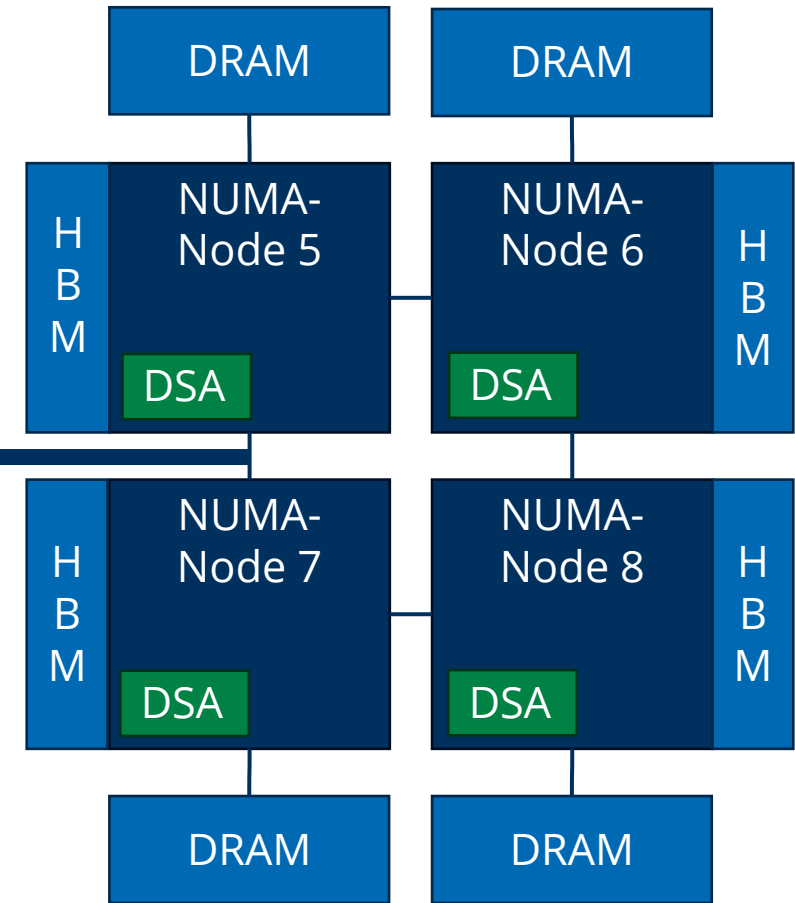
Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 18

Funded by

# Benchmarking the DSA

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 19

Funded by

# Benchmarking the DSA
## Setup

**System under Test:** Intel Xeon CPU Max 9468

# Benchmarking the DSA
## Setup

**System under Test:** Intel Xeon CPU Max 9468

1) CPU vs. DSA

# Benchmarking the DSA
## Setup

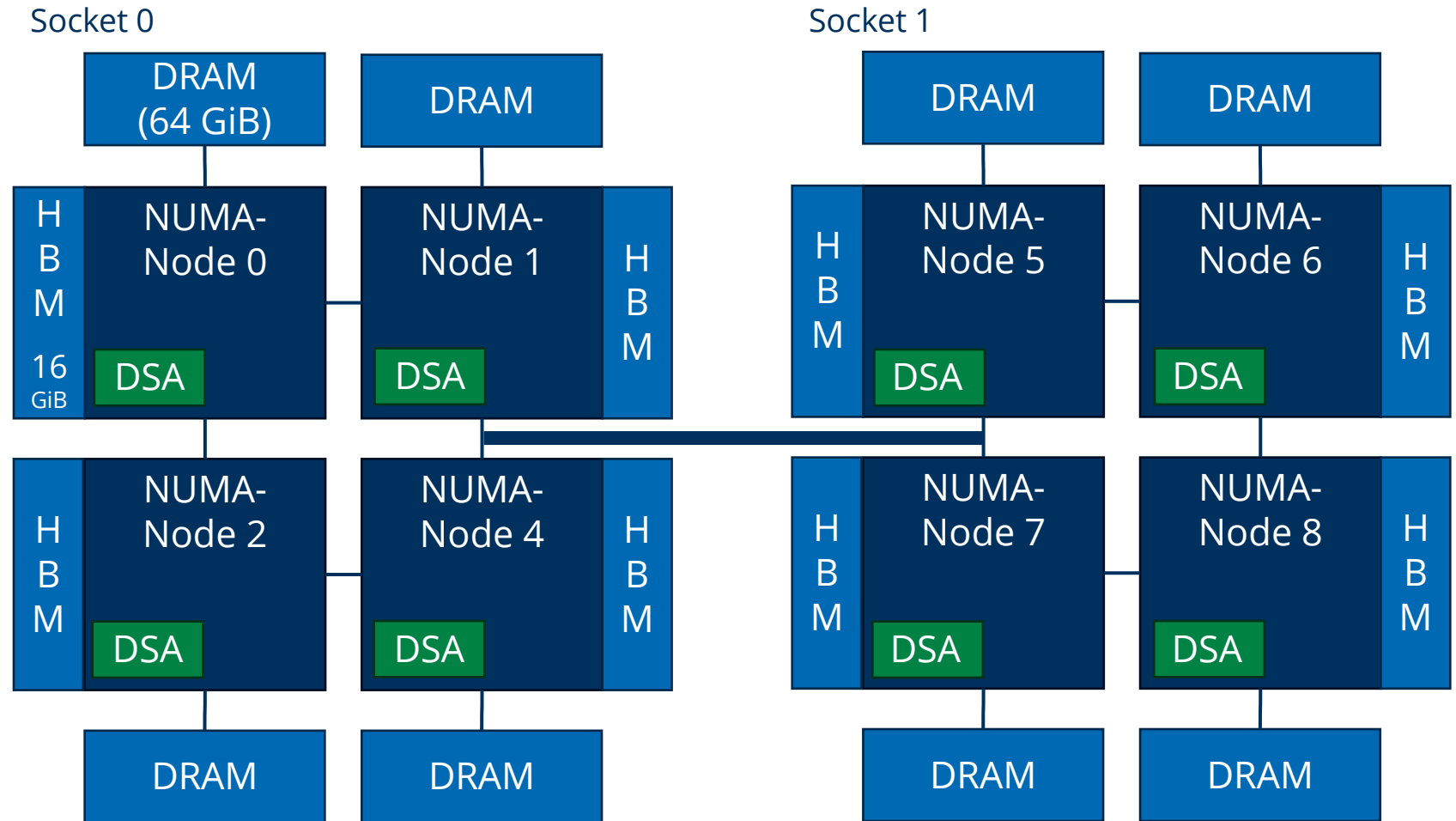**System under Test:** Intel Xeon CPU Max 9468

1) CPU vs. DSA
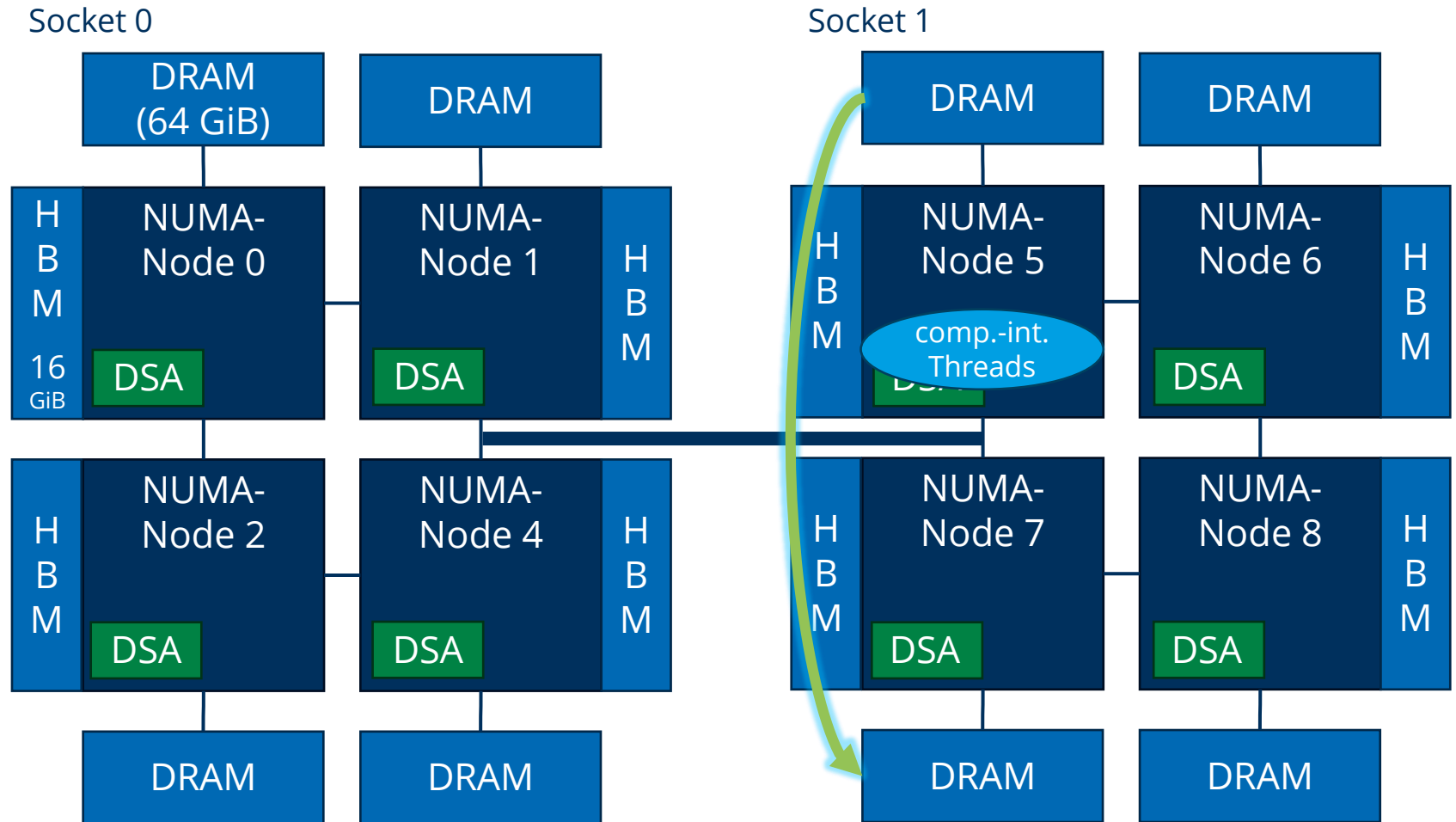
2) Interference between DSA and CPU

Funded by

# Benchmarking the DSA
## Setup

**System under Test:** Intel Xeon CPU Max 9468

1) CPU vs. DSA
2) Interference between DSA and CPU
   1) Compute-Intensive CPU Threads

Socket 0

Socket 1

DRAM (64 GiB)

DRAM

DRAM

DRAM

HBM 16 GiB

NUMA-Node 0

DSA

NUMA-Node 1

DSA

HBM

HBM

NUMA-Node 5

DSA

comp.-int. Threads

NUMA-Node 6

DSA

HBM

HBM

NUMA-Node 2

DSA

NUMA-Node 4

DSA

HBM

HBM

NUMA-Node 7

DSA

NUMA-Node 8

DSA

HBM

DRAM

DRAM

DRAM

DRAM

# Benchmarking the DSA
## Setup

1) CPU vs. DSA

2) Interference between
   DSA and CPU
   1) Compute-Intensive
      CPU Threads
   2) Data-Intensive
      CPU Threads

**System under Test:** Intel Xeon CPU Max 9468

# Benchmarking the DSA
## Setup

1) CPU vs. DSA
2) Interference between
   DSA and CPU
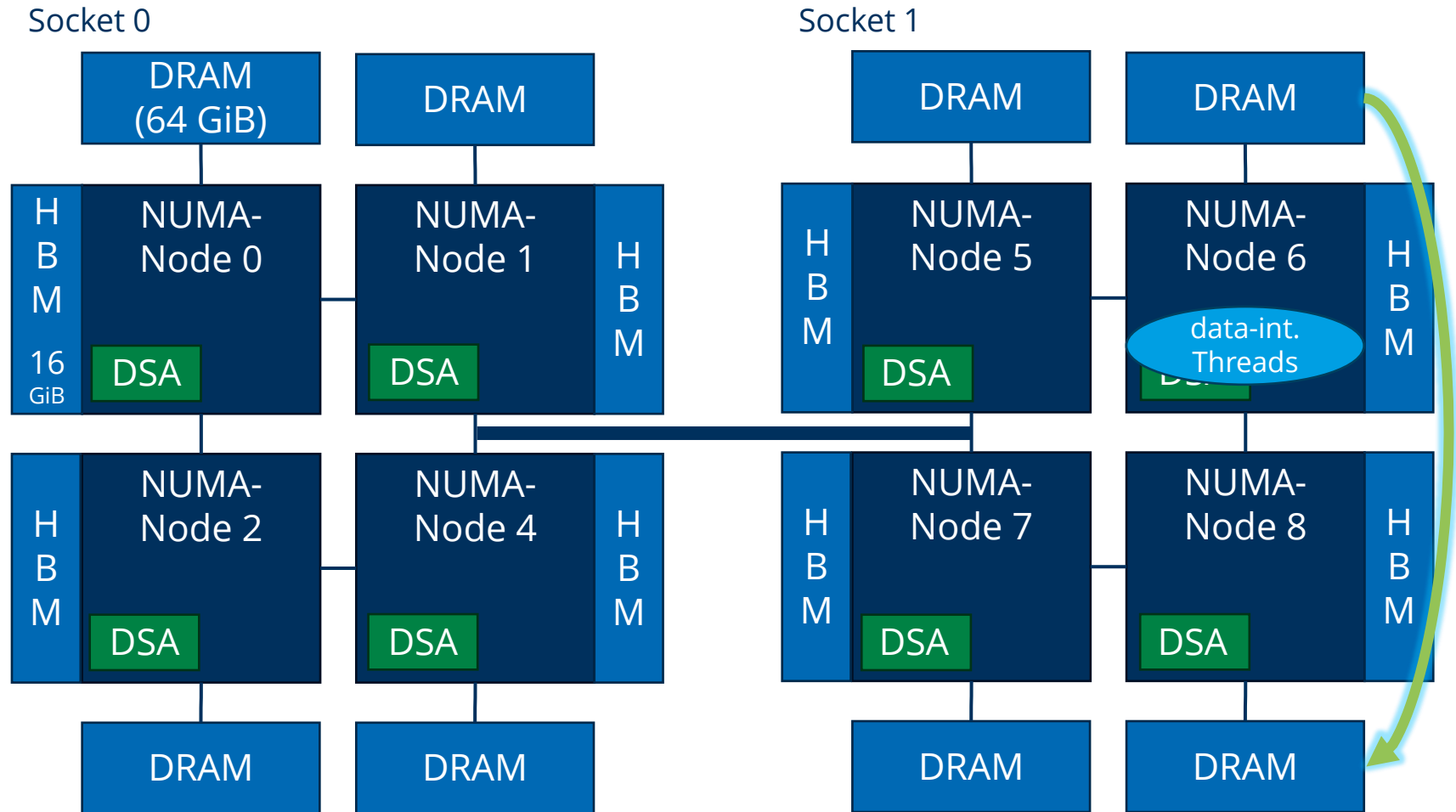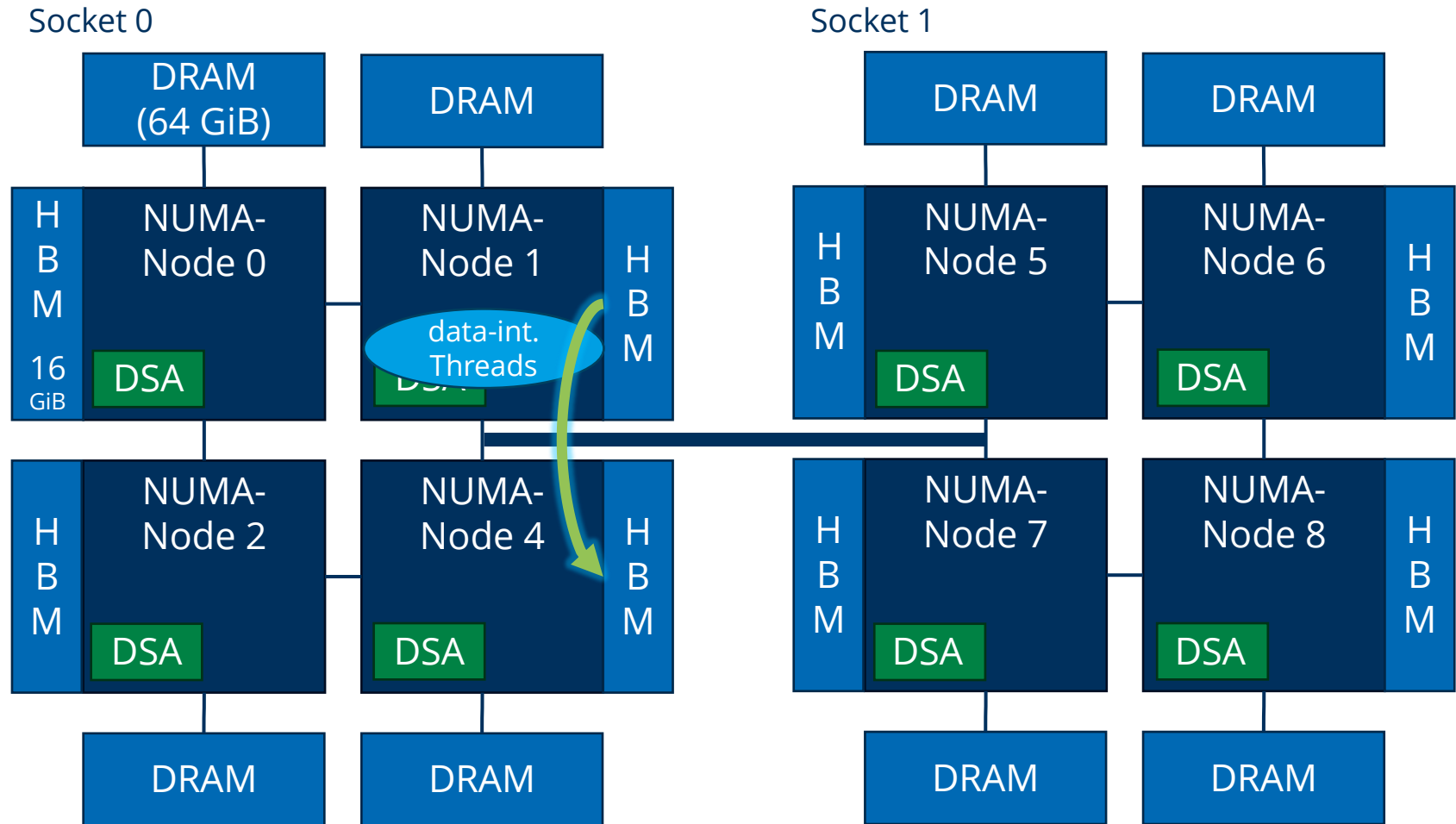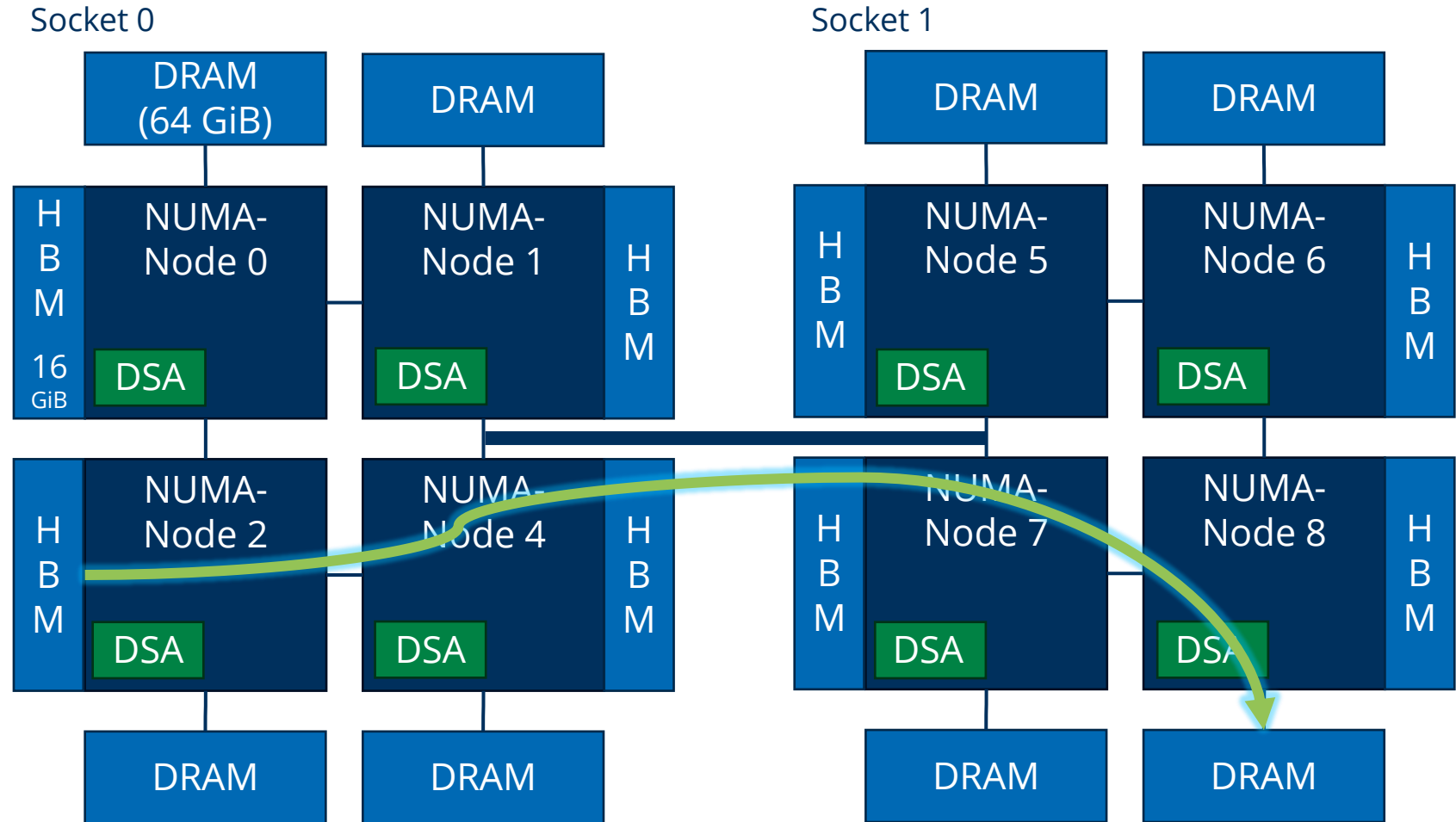   1) Compute-Intensive
      CPU Threads
   2) Data-Intensive
      CPU Threads
   3) Data-Intensive
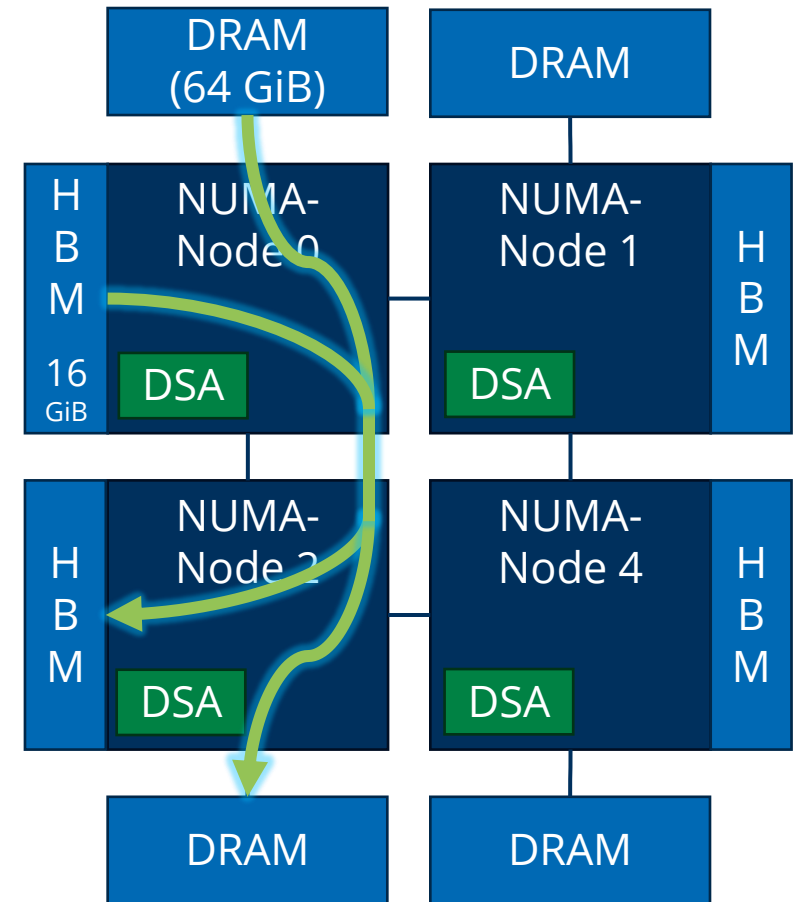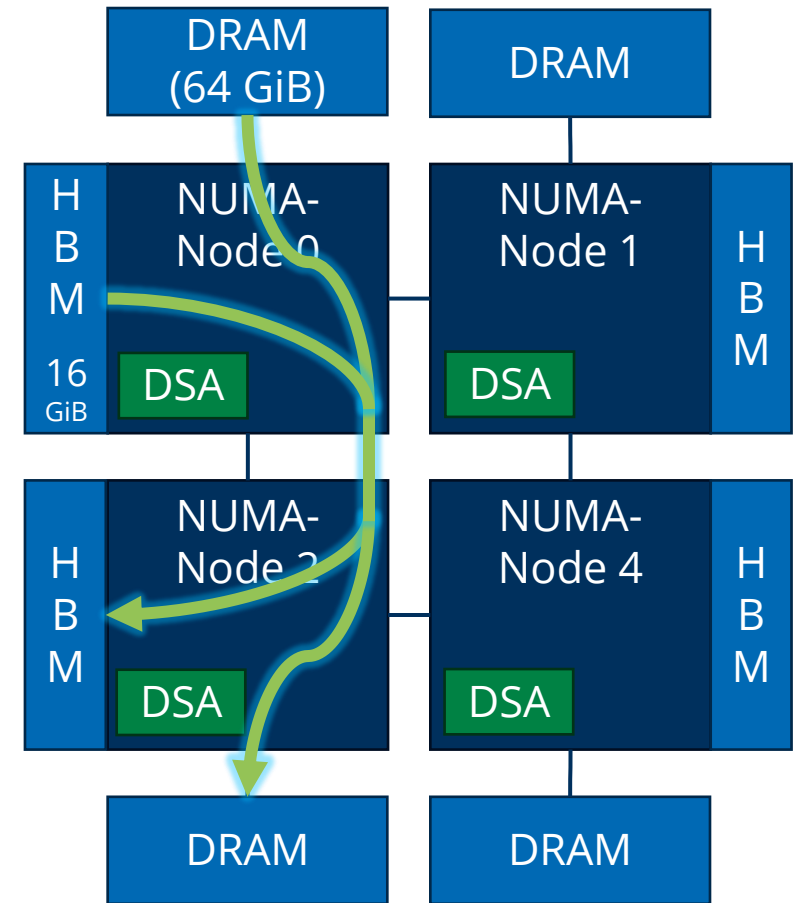      CPU Threads
      on HBM

**System under Test:** Intel Xeon CPU Max 9468

# Benchmarking the DSA
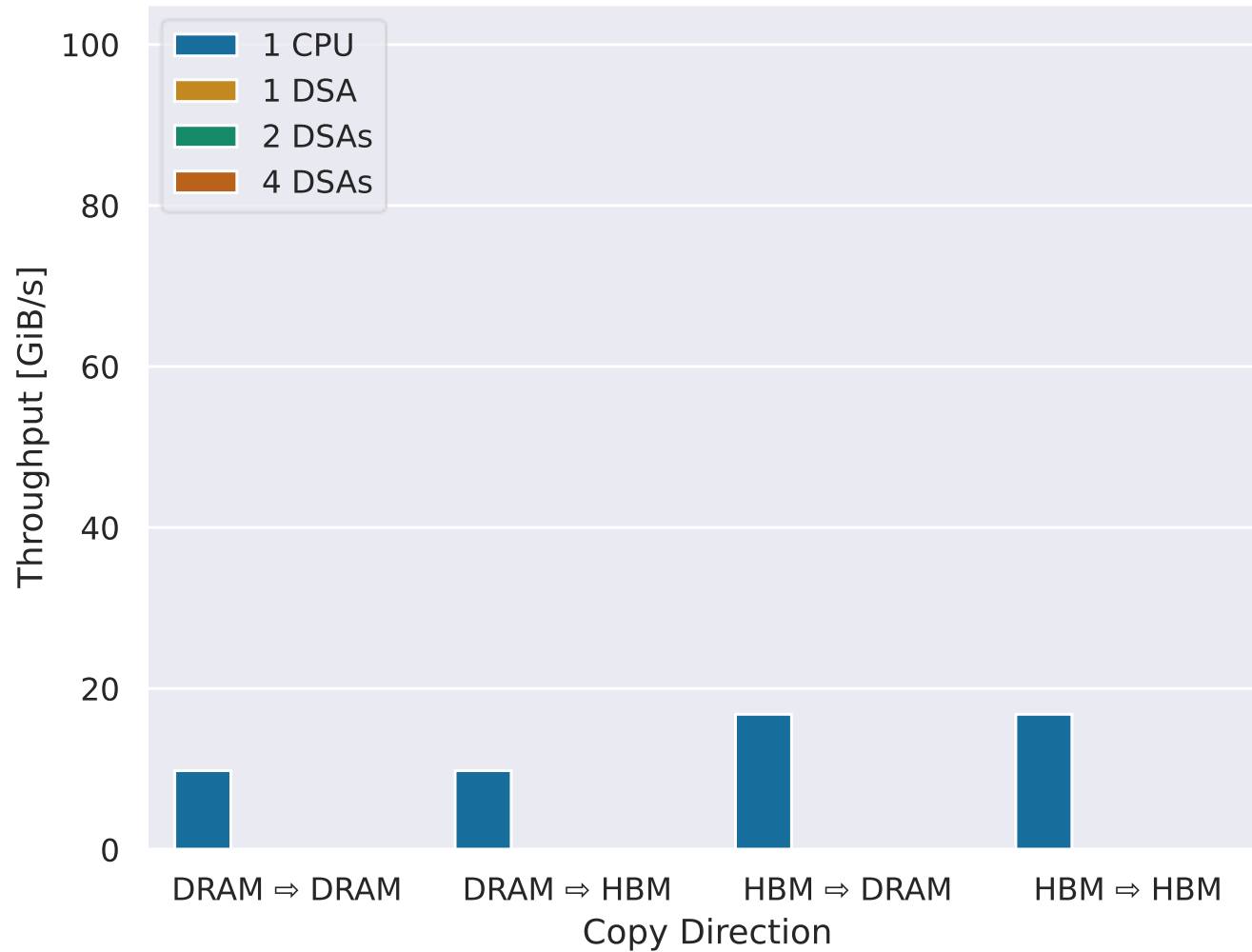## Setup

**System under Test:** Intel Xeon CPU Max 9468

1) CPU vs. DSA
2) Interference between DSA and CPU
   1) Compute-Intensive CPU Threads
   2) Data-Intensive CPU Threads
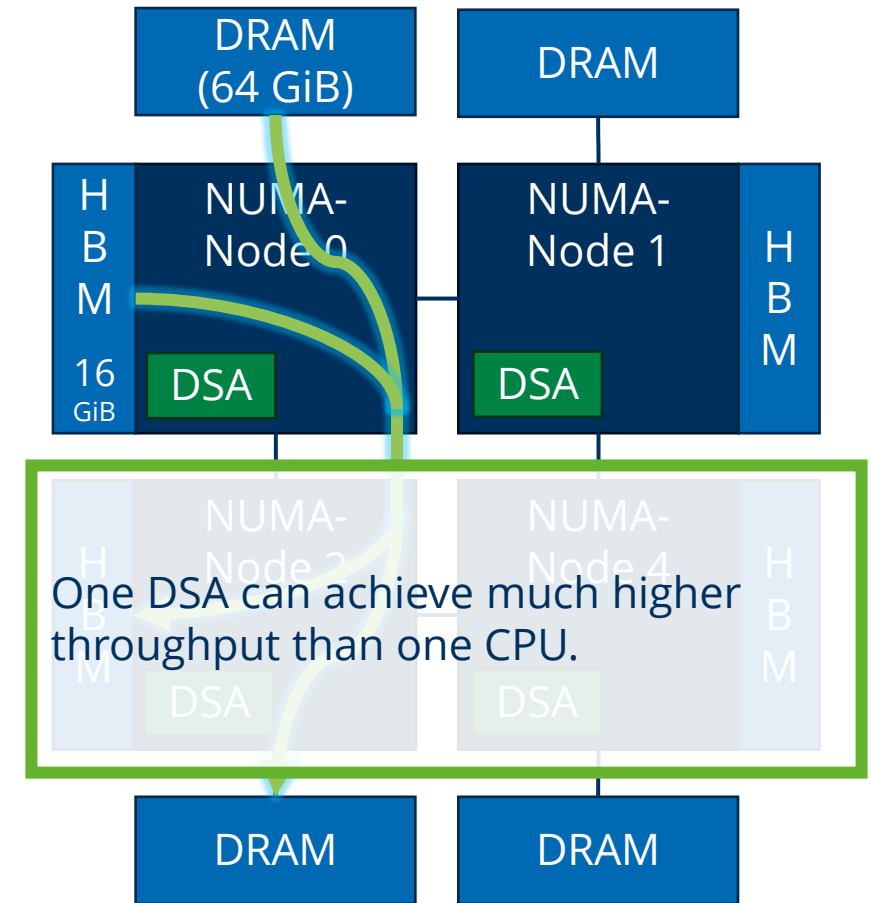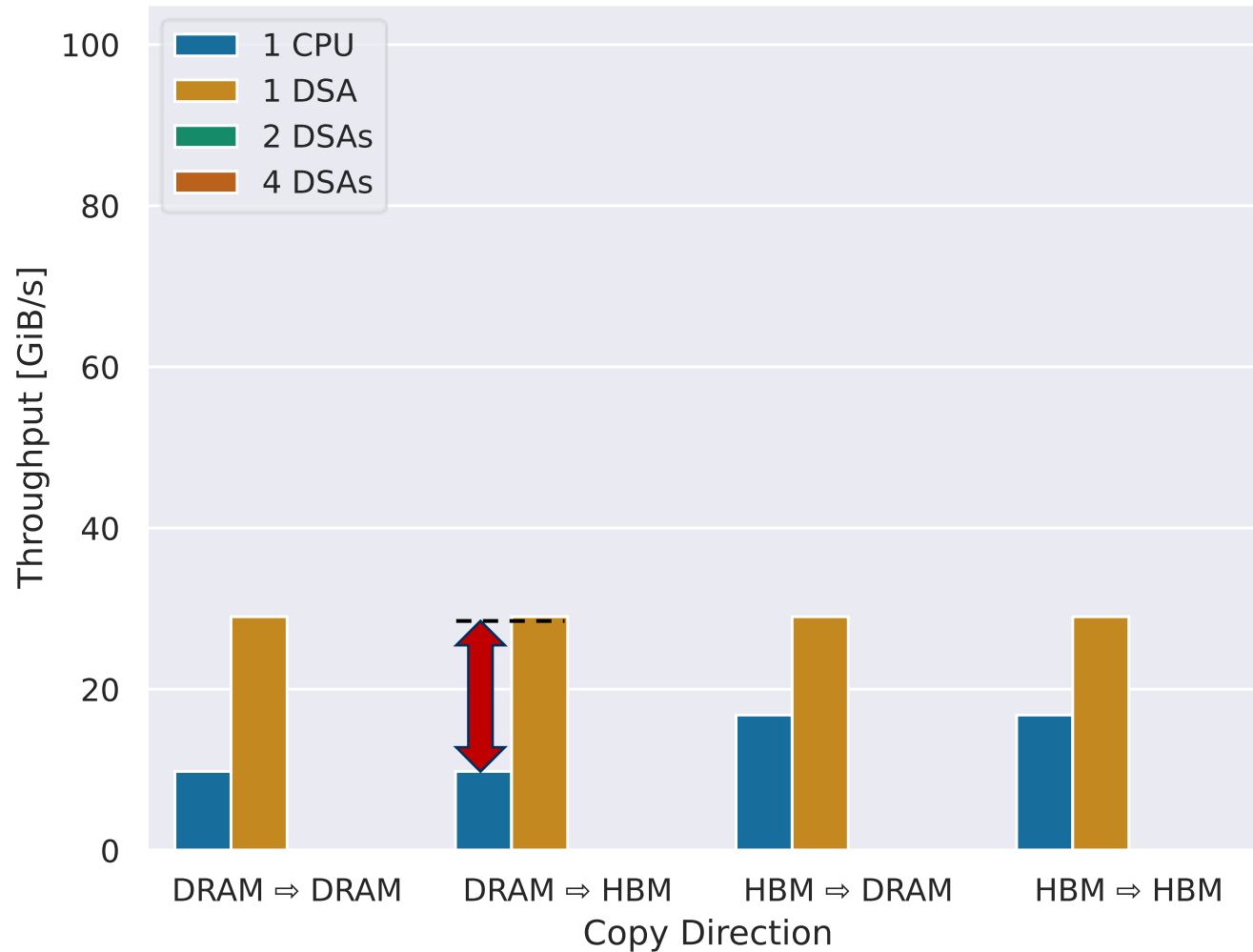   3) Data-Intensive CPU Threads on HBM
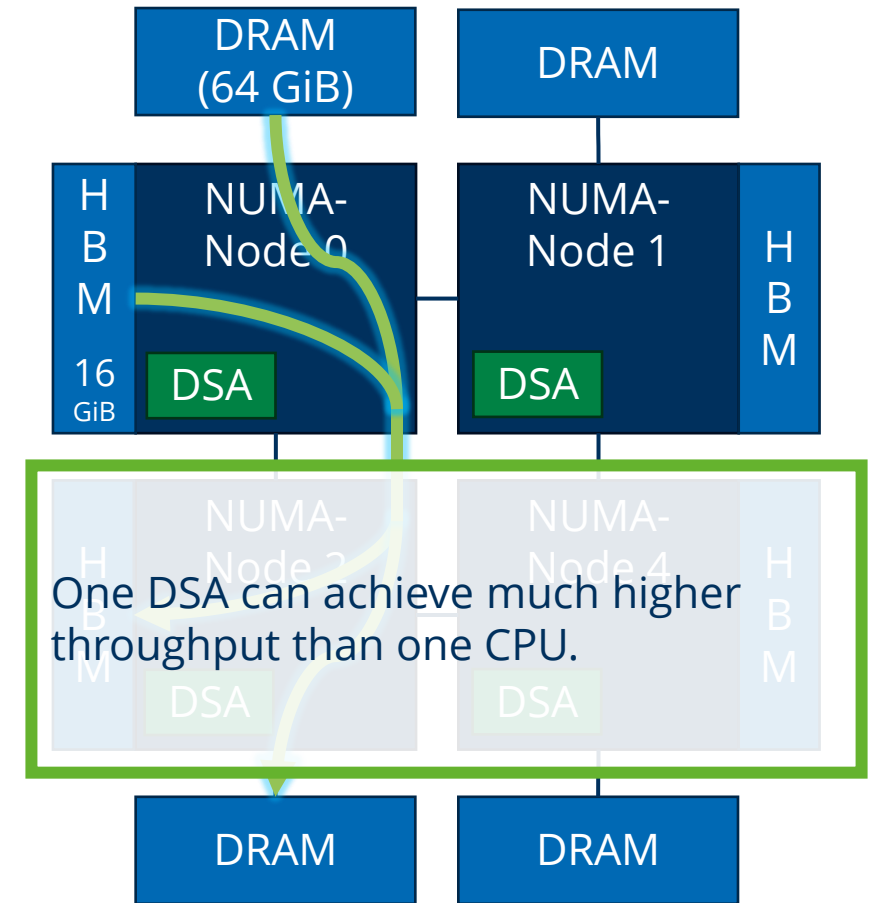3) Inter-Socket Data Transfer with DSA

TECHNISCHE UNIVERSITÄT DRESDEN
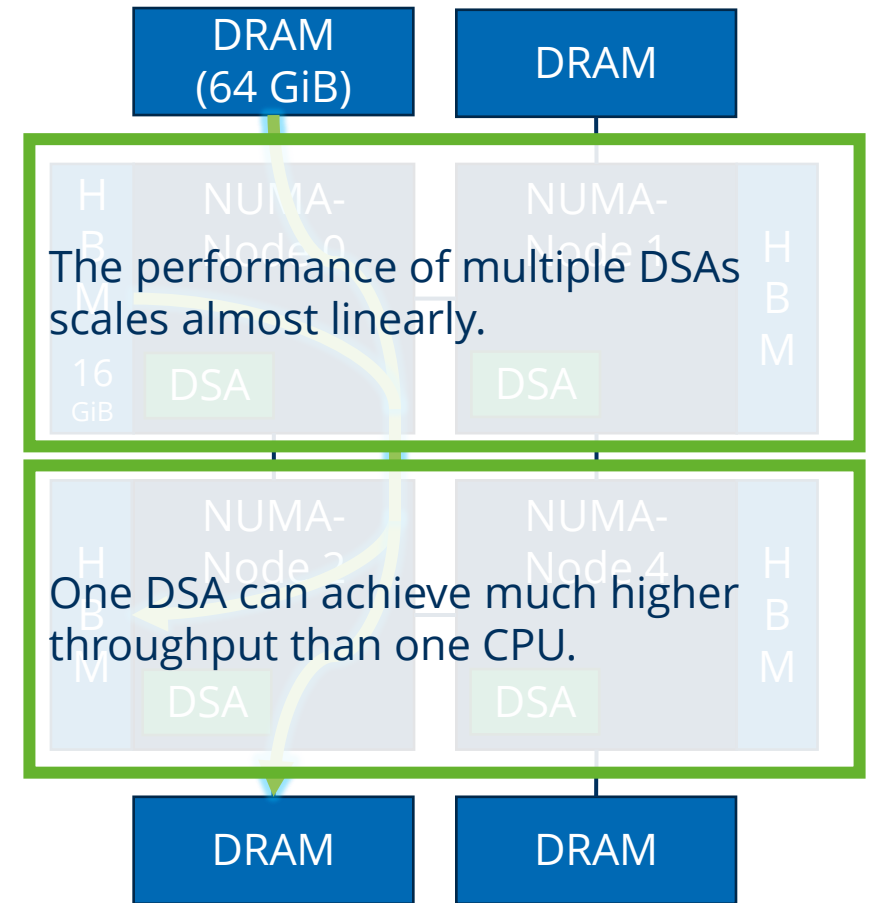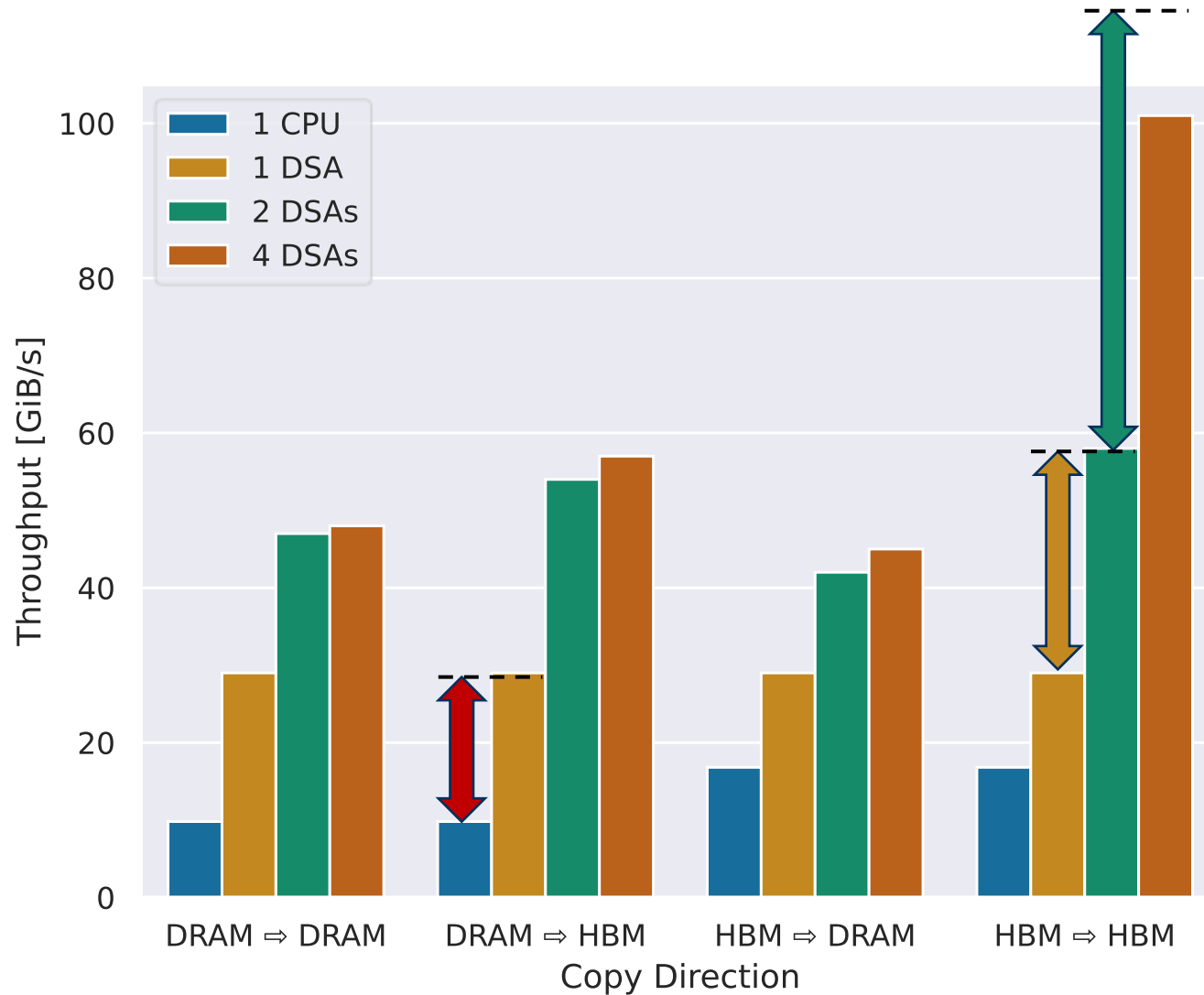
Funded by

DFG

DRESDEN concept

# Benchmark – CPU vs. DSA

# Benchmark – CPU vs. DSA

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Funded by

# Benchmark – CPU vs. DSA



One DSA can achieve much higher throughput than one CPU.

# Benchmark – CPU vs. DSA



One DSA can achieve much higher throughput than one CPU.

# Benchmark – CPU vs. DSA



The performance of multiple DSAs scales almost linearly.

One DSA can achieve much higher throughput than one CPU.

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 31

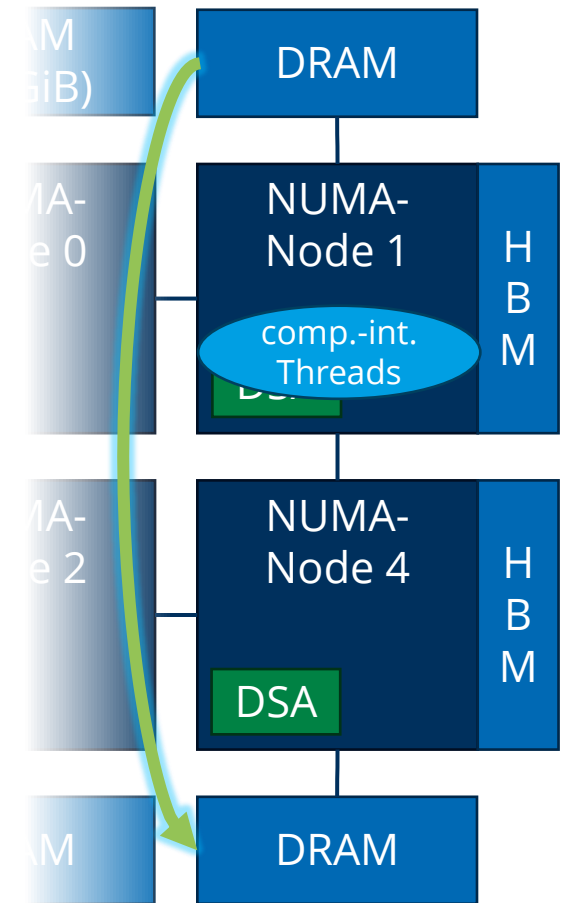# Benchmark – Interference between DSA and CPU
## Compute-Intensive CPU Threads

# Benchmark – Interference between DSA and CPU
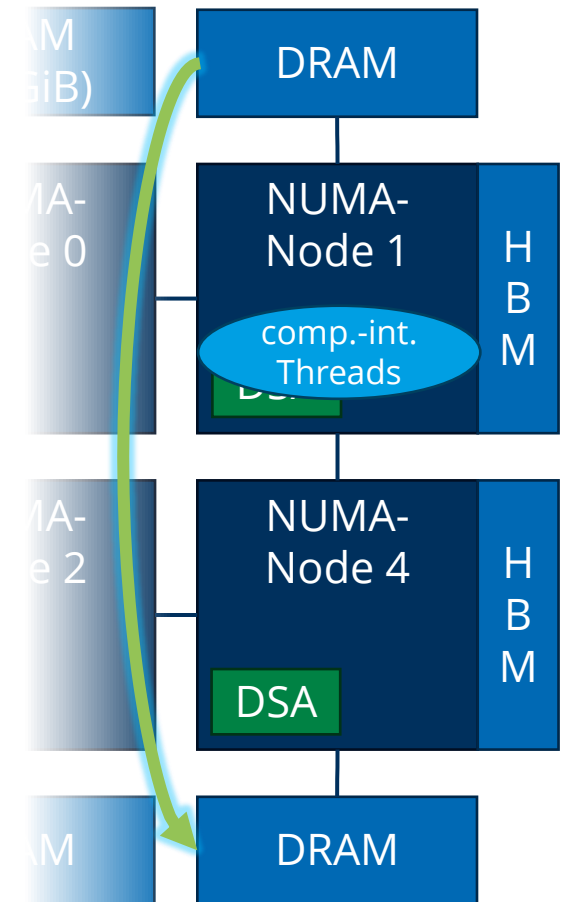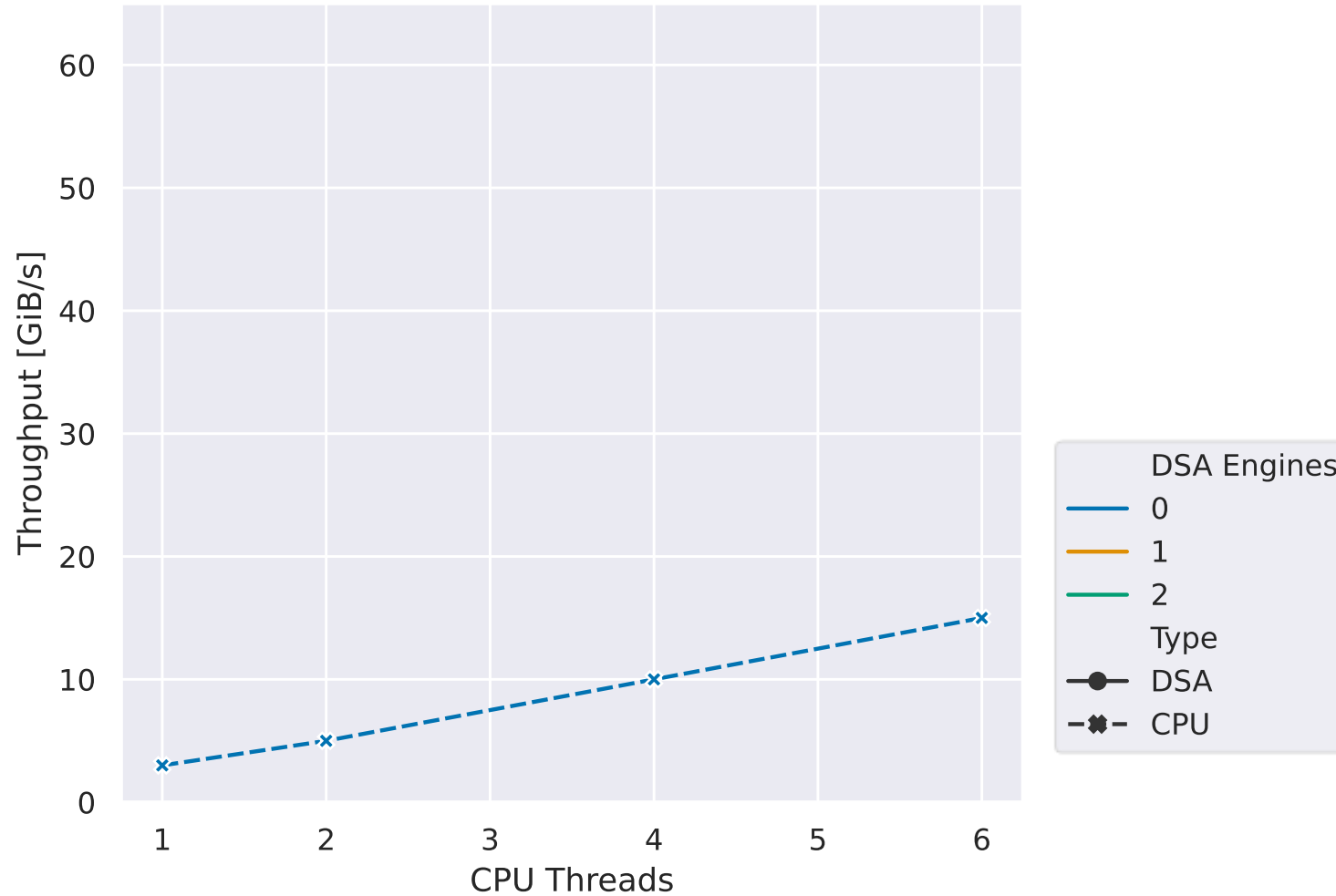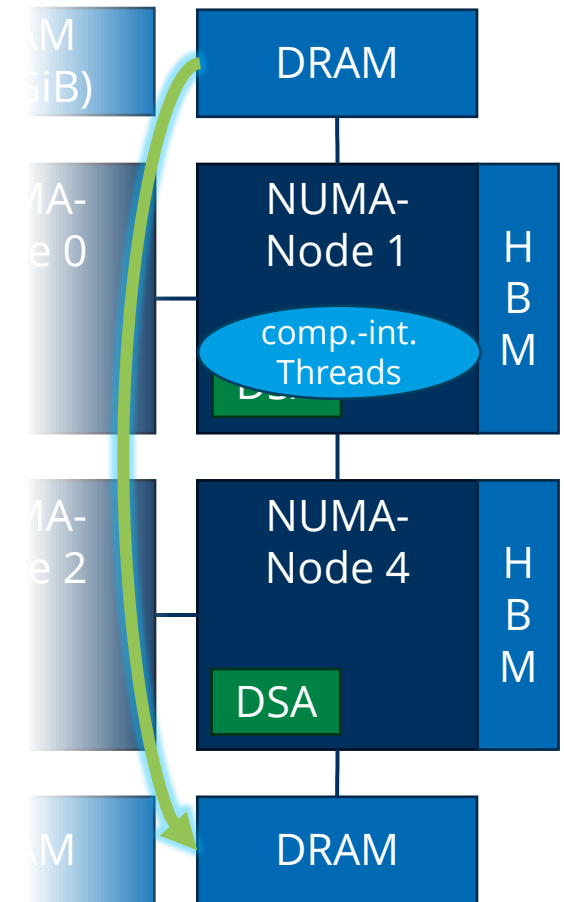## Compute-Intensive CPU Threads

# Benchmark – Interference between DSA and CPU
## Compute-Intensive CPU Threads

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
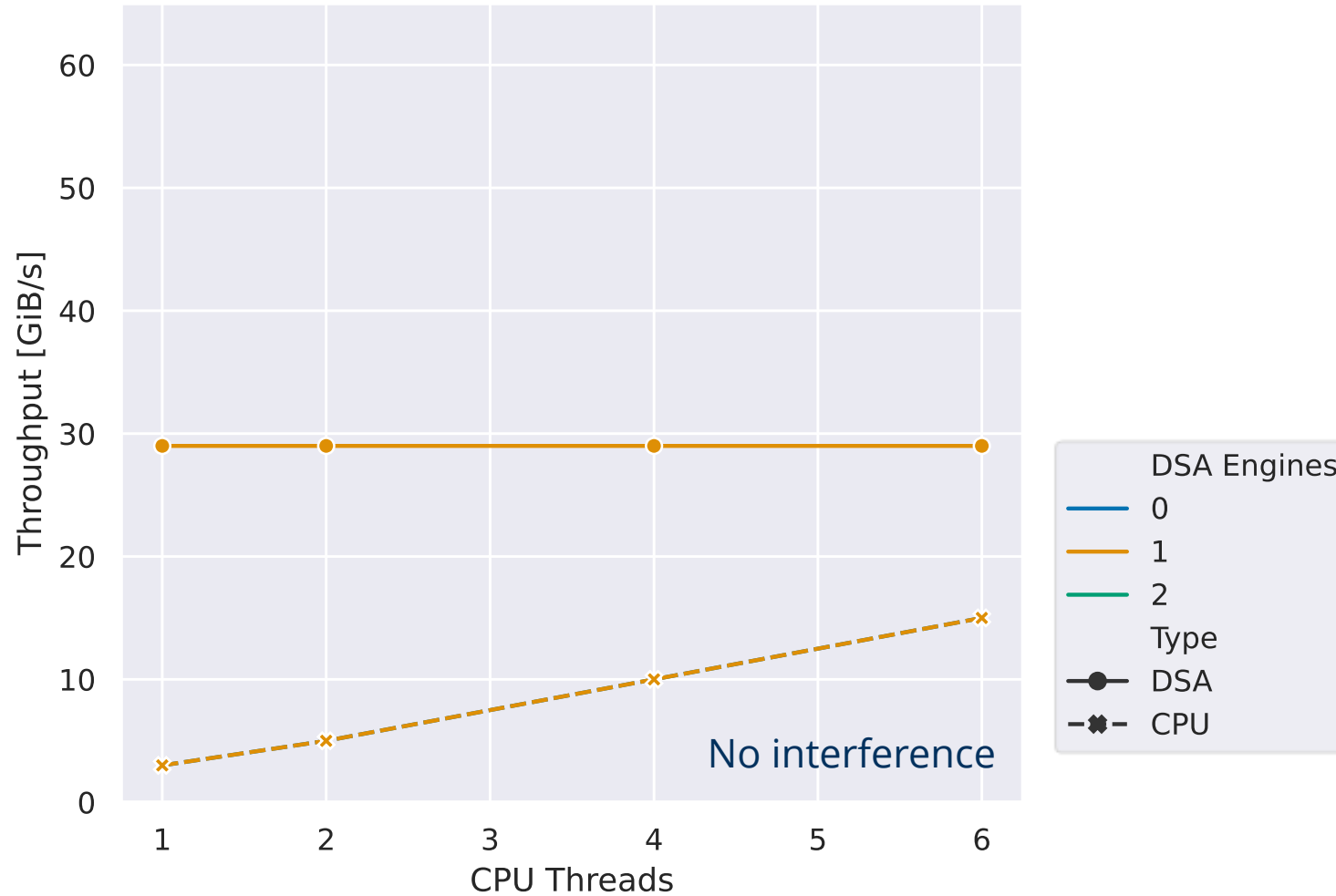Dresden University of Technology / André Berthold

Funded by

# Benchmark – Interference between DSA and CPU
## Compute-Intensive CPU Threads

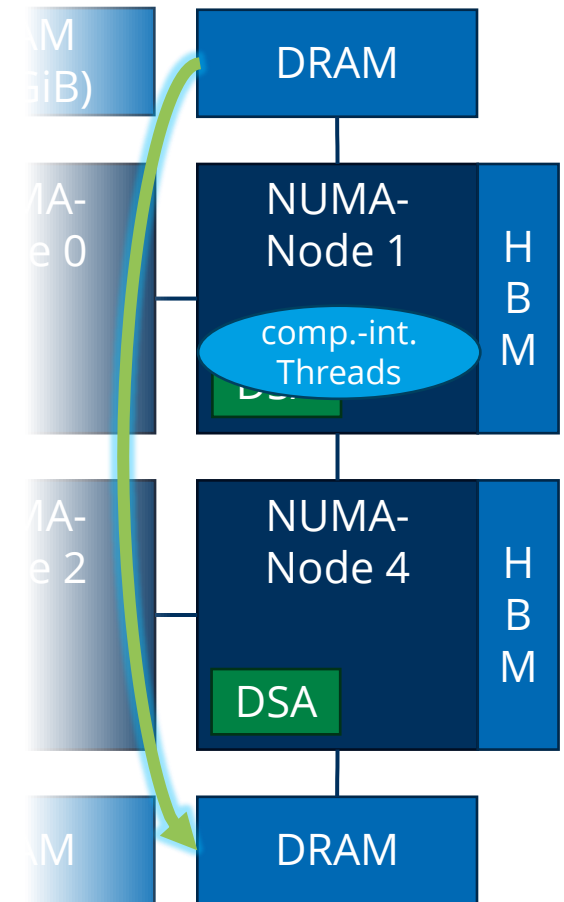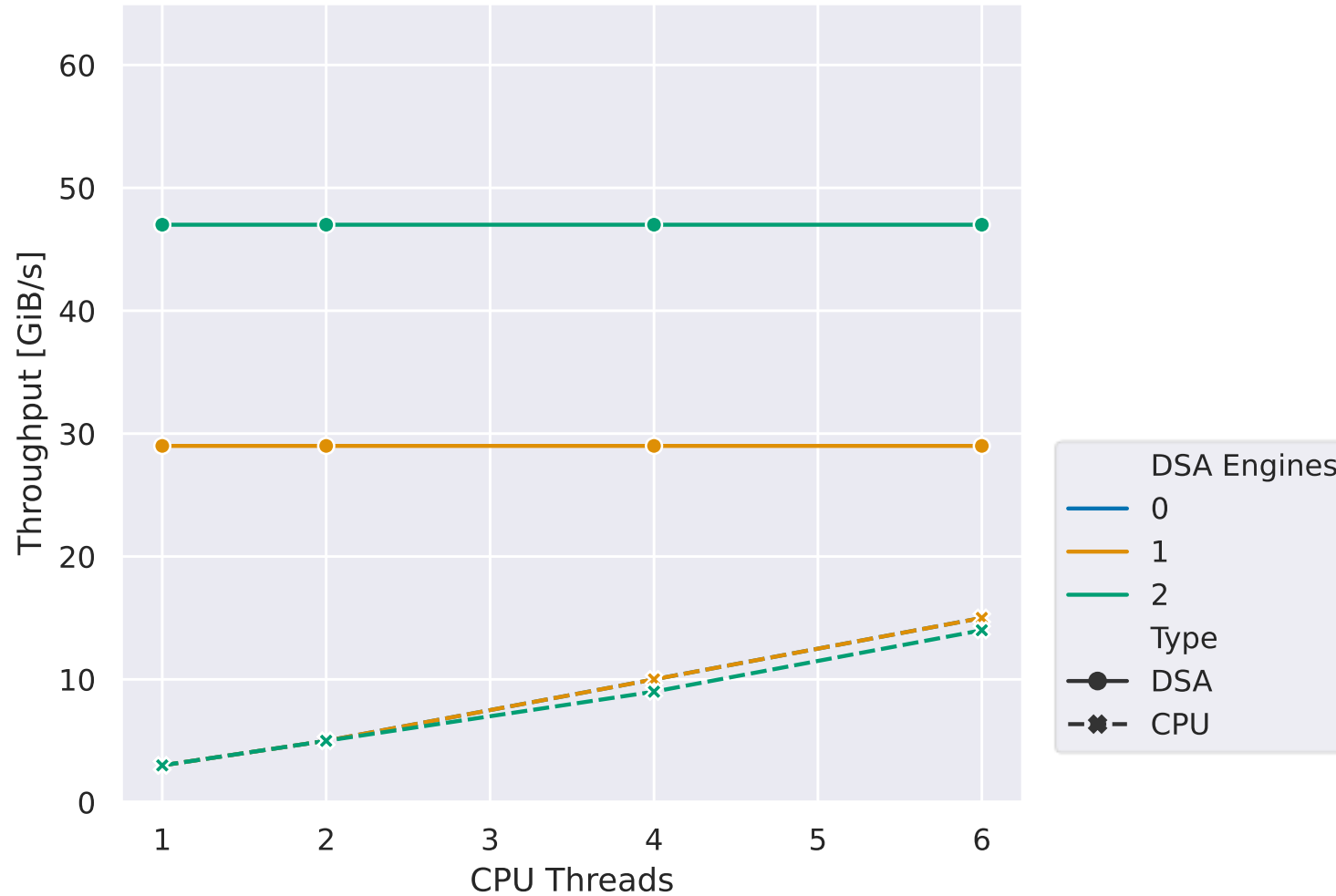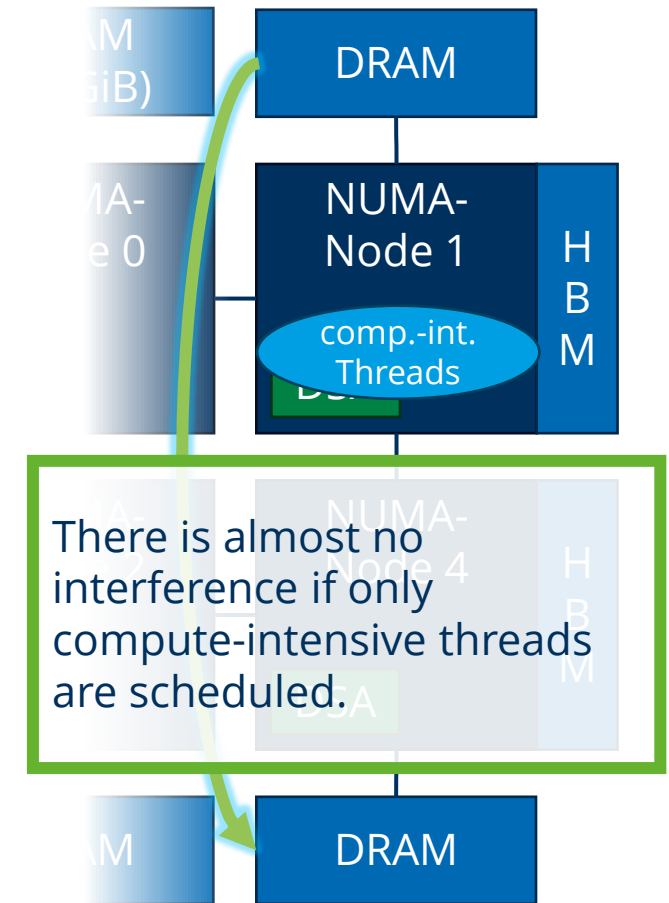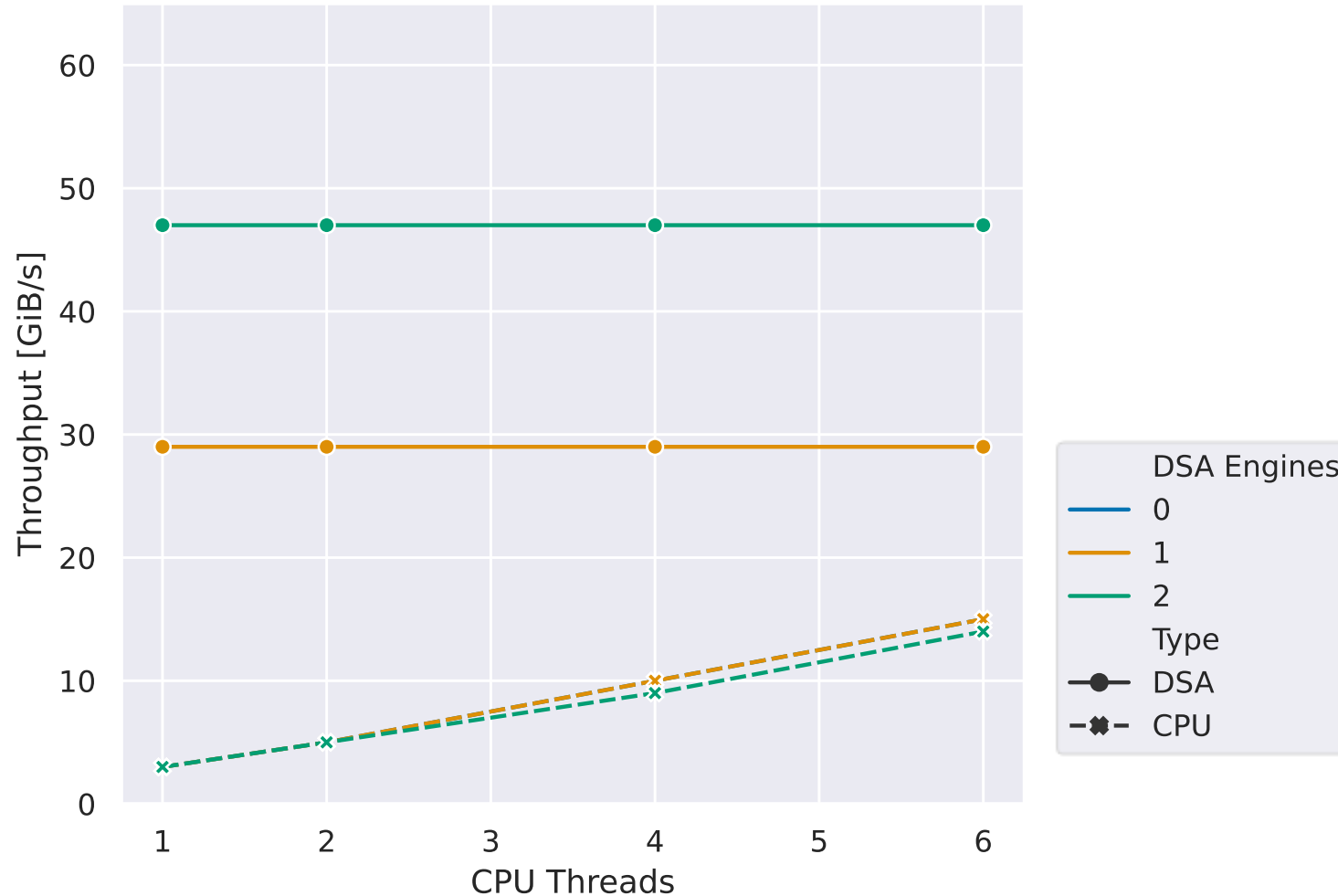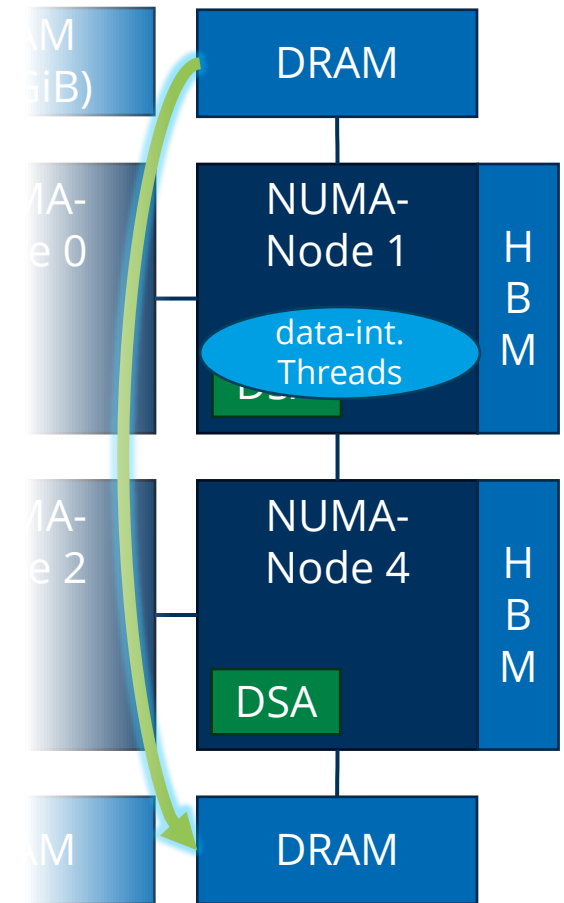# Benchmark – Interference between DSA and CPU
## Compute-Intensive CPU Threads



There is almost no interference if only compute-intensive threads are scheduled.

# Benchmark – Interference between DSA and CPU
## Data-Intensive CPU Threads (Vectorized, AVX-512)

# Benchmark – Interference between DSA and CPU
## Data-Intensive CPU Threads (Vectorized, AVX-512)

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
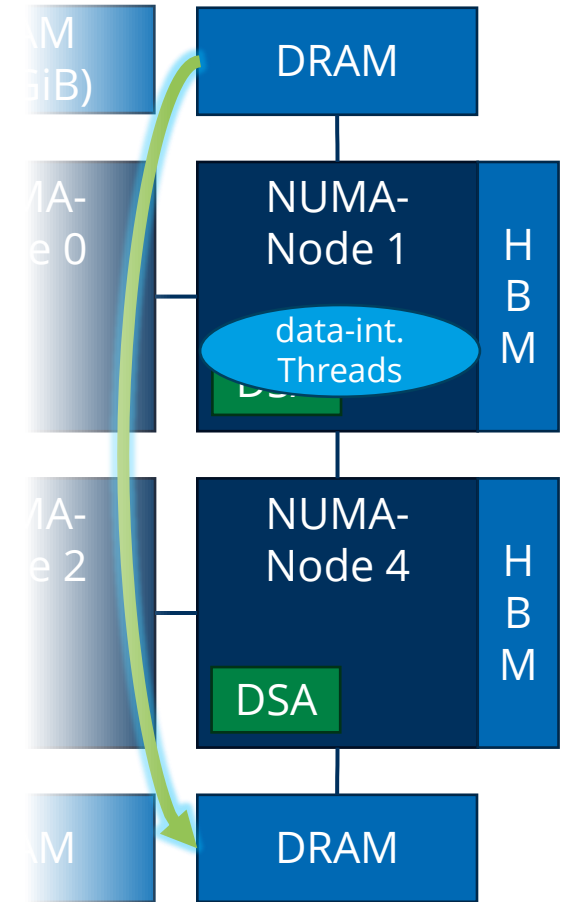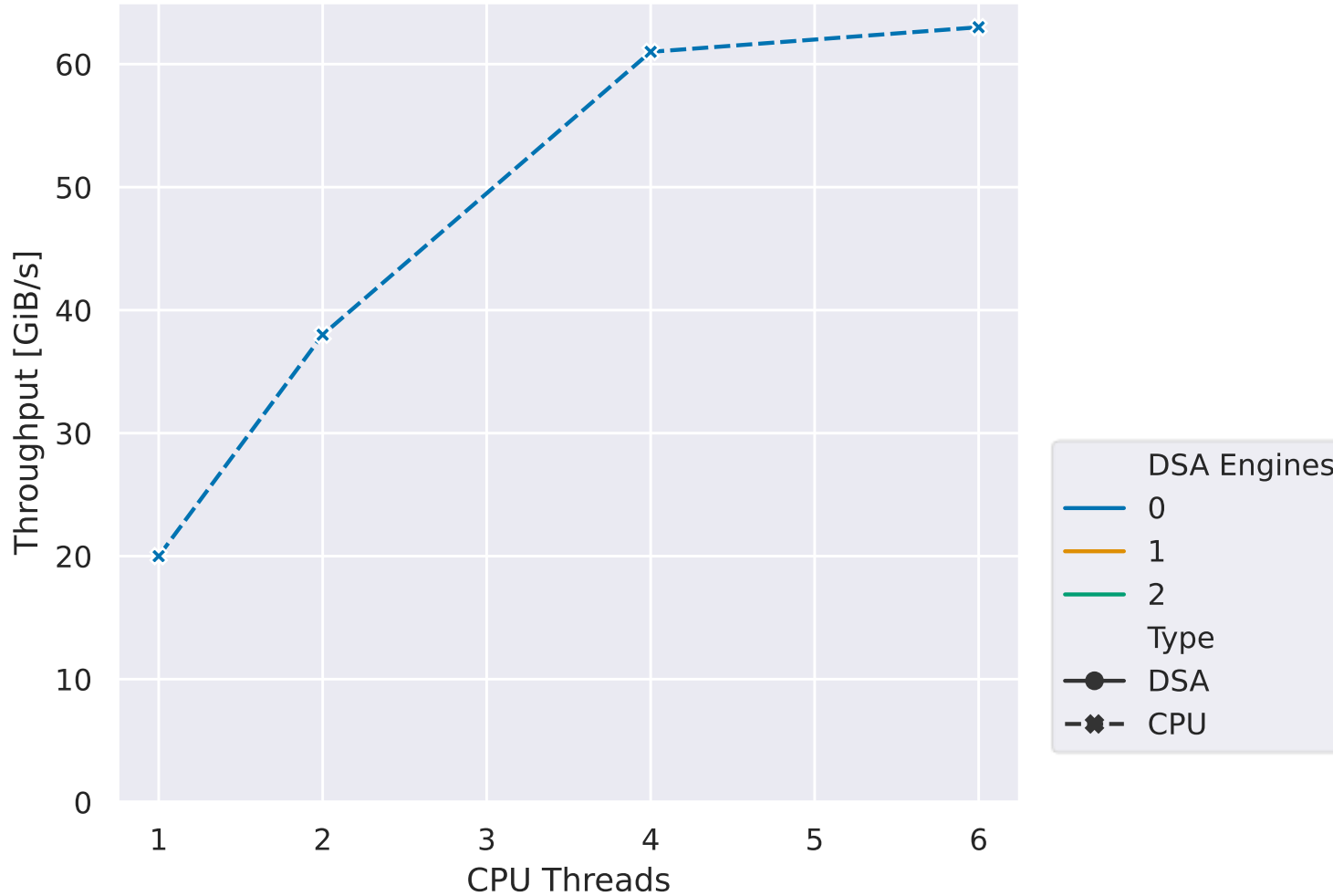Dresden University of Technology / André Berthold

Funded by

# Benchmark – Interference between DSA and CPU
## Data-Intensive CPU Threads (Vectorized, AVX-512)

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold
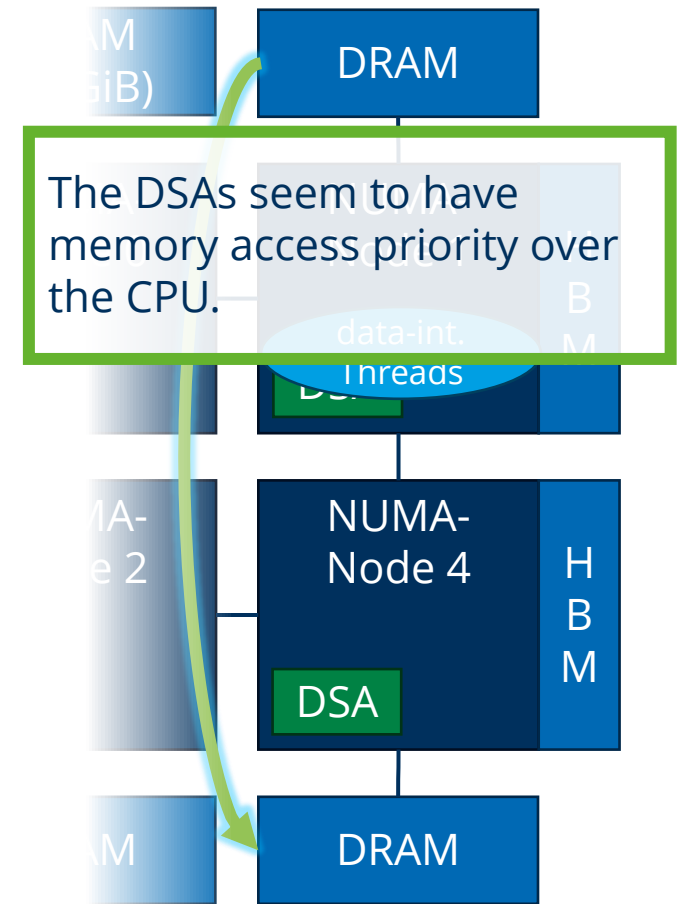
Slide 39

Funded by

# Benchmark – Interference between DSA and CPU
## Data-Intensive CPU Threads (Vectorized, AVX-512)



The DSAs seem to have memory access priority over the CPU.

Concurrent data-intensive threads result in reduced throughput for **multiple** DSA engines.

# Benchmark – Interference between DSA and CPU
## Data-Intensive CPU Threads on HBM

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Funded by

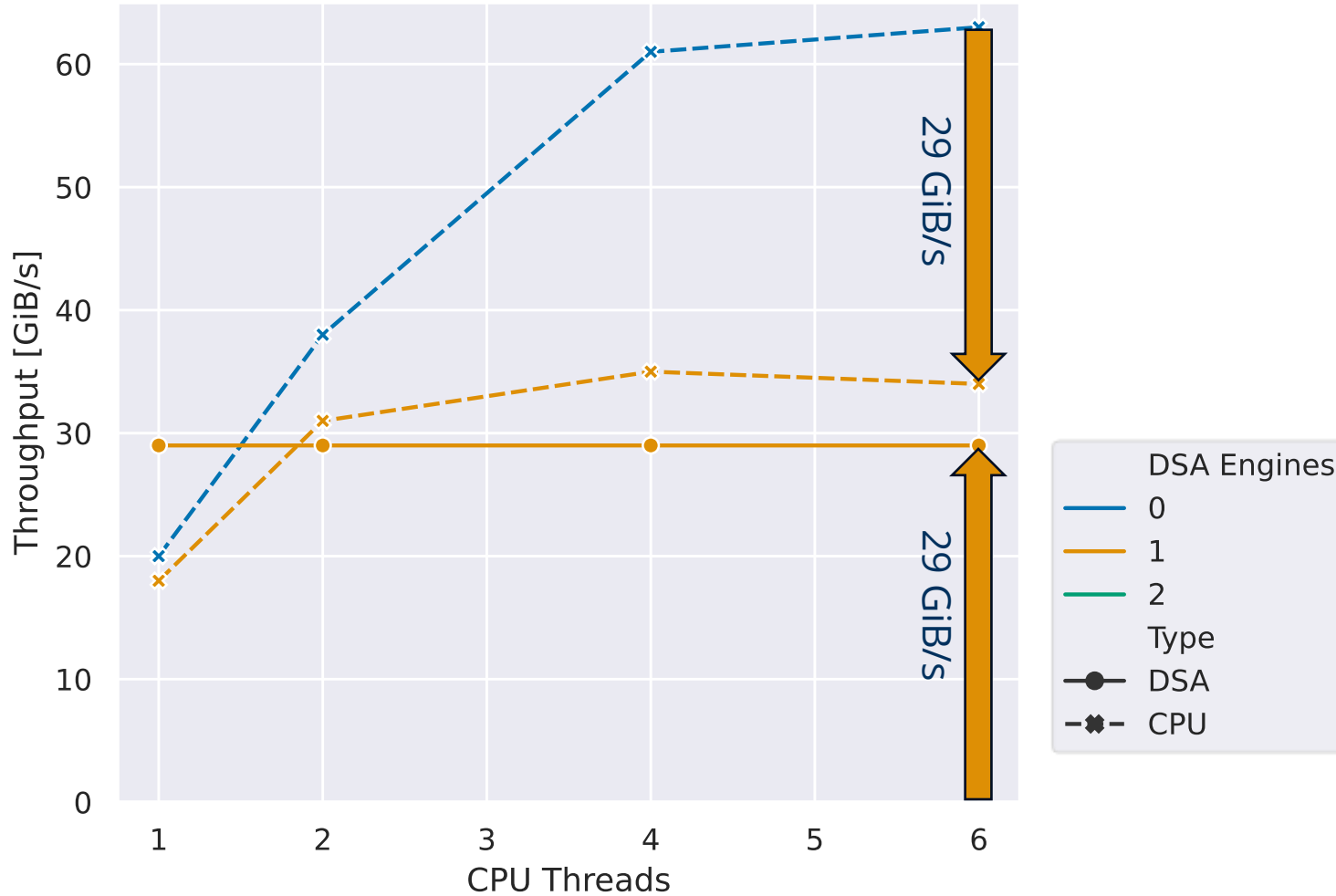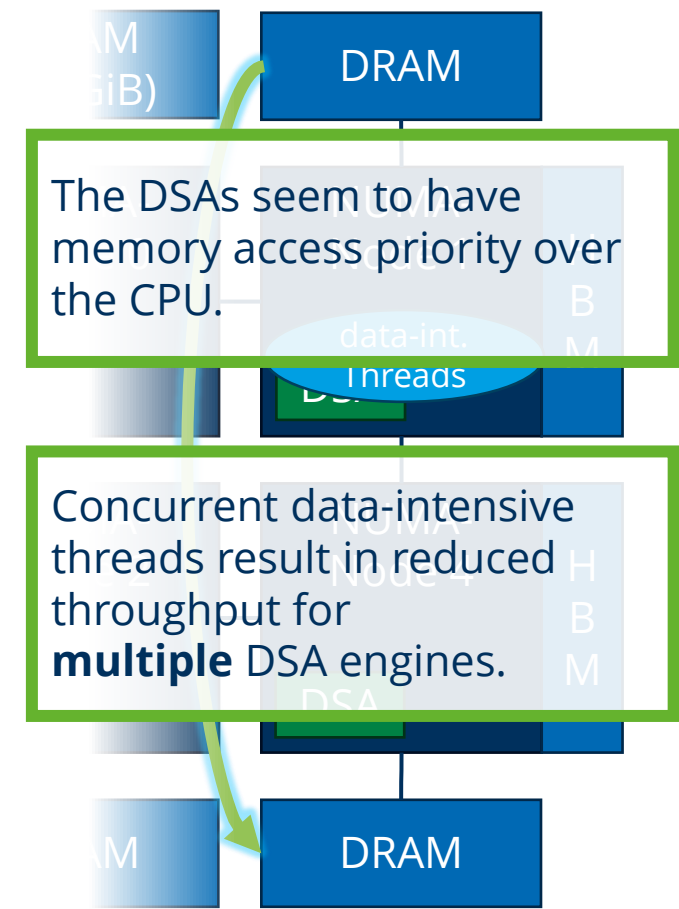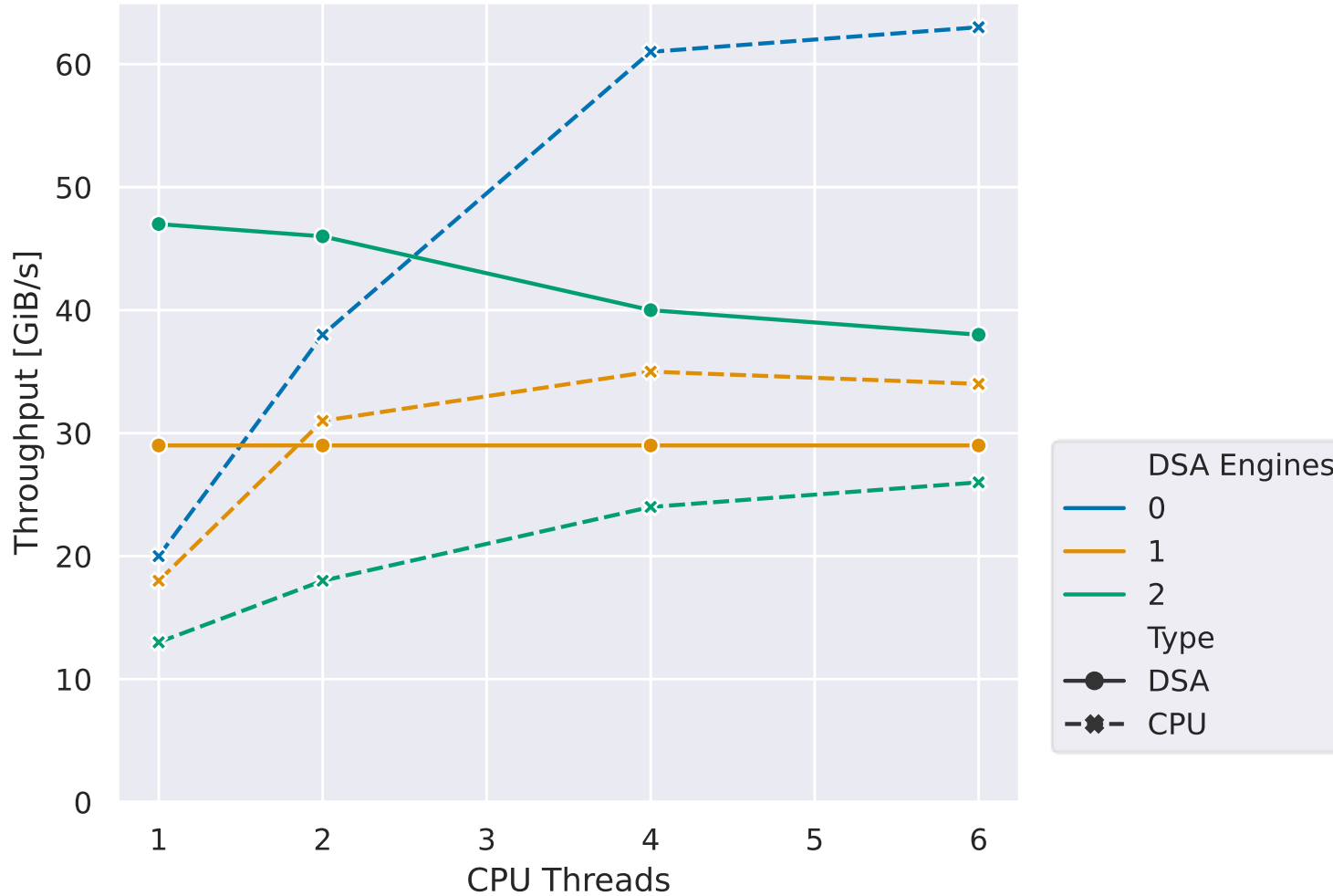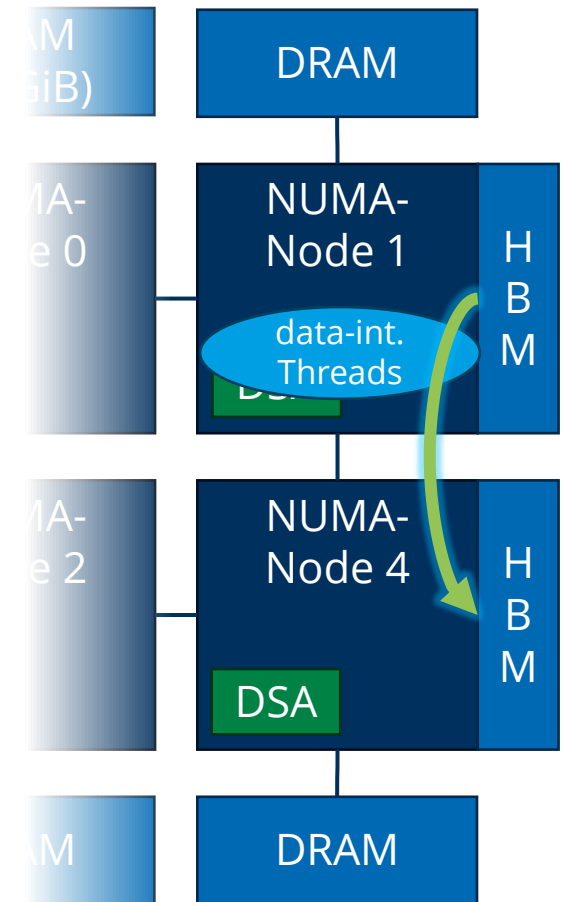# Benchmark – Interference between DSA and CPU
## Data-Intensive CPU Threads on HBM

# Benchmark – Interference between DSA and CPU
## Data-Intensive CPU Threads on HBM

# Benchmark – Interference between DSA and CPU
## Data-Intensive CPU Threads on HBM

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
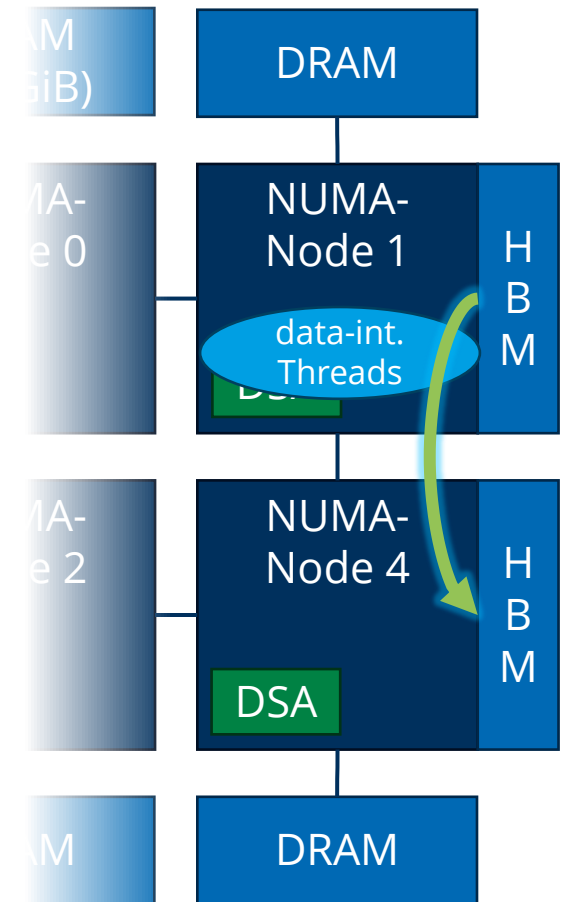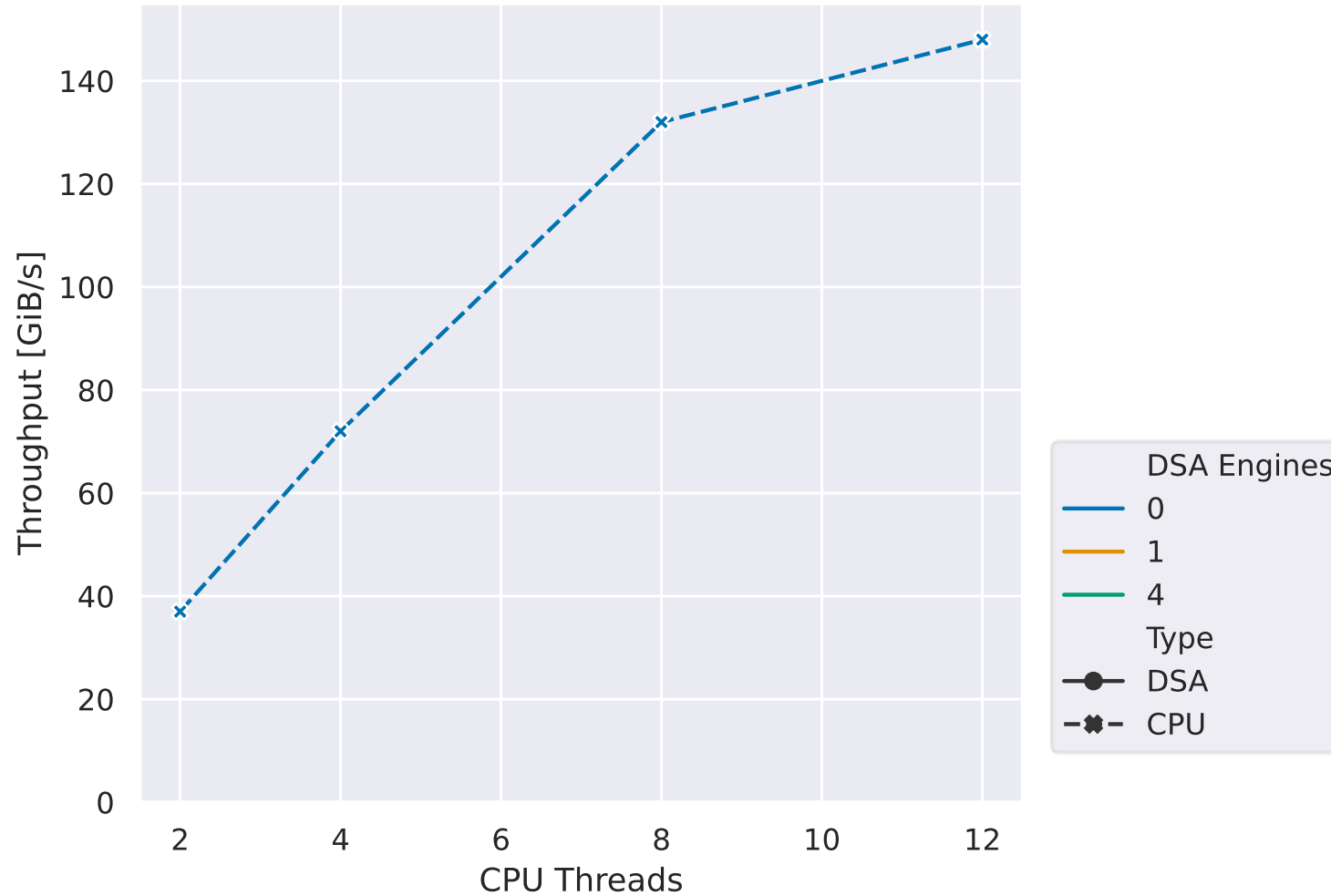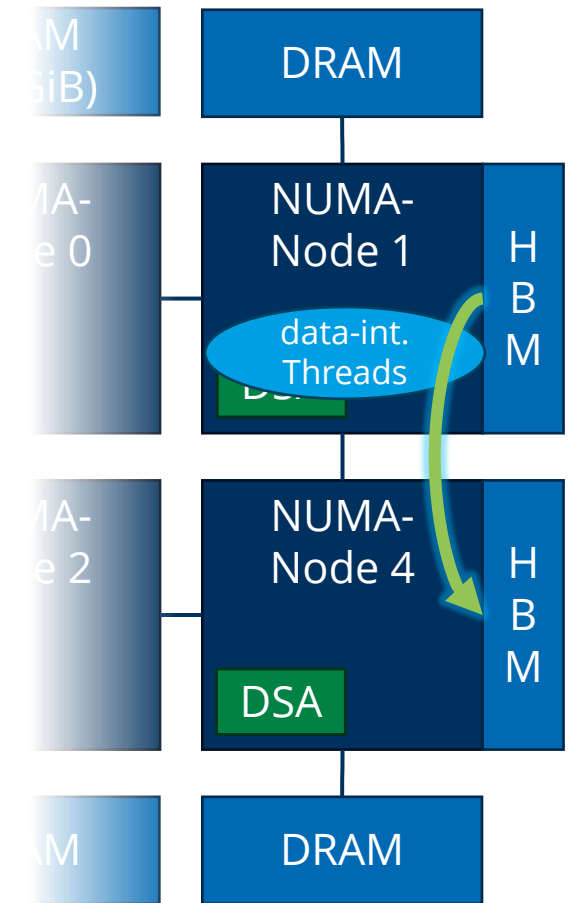Dresden University of Technology / André Berthold

Funded by

# Benchmark – Interference between DSA and CPU
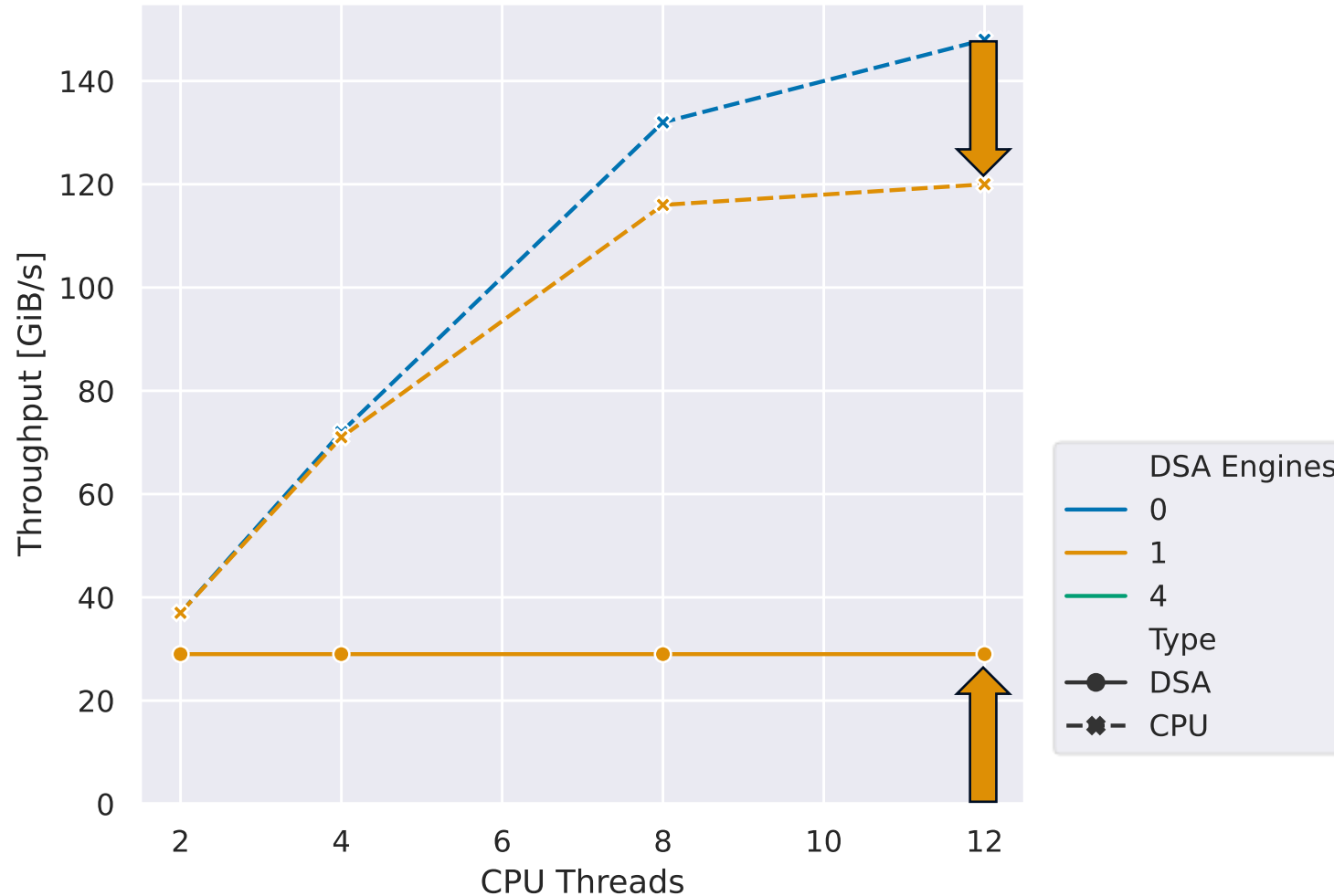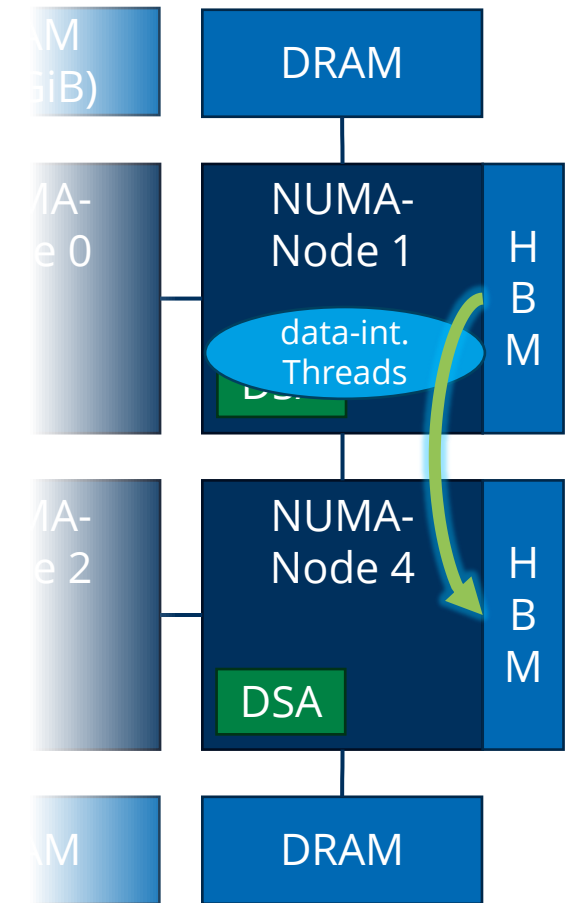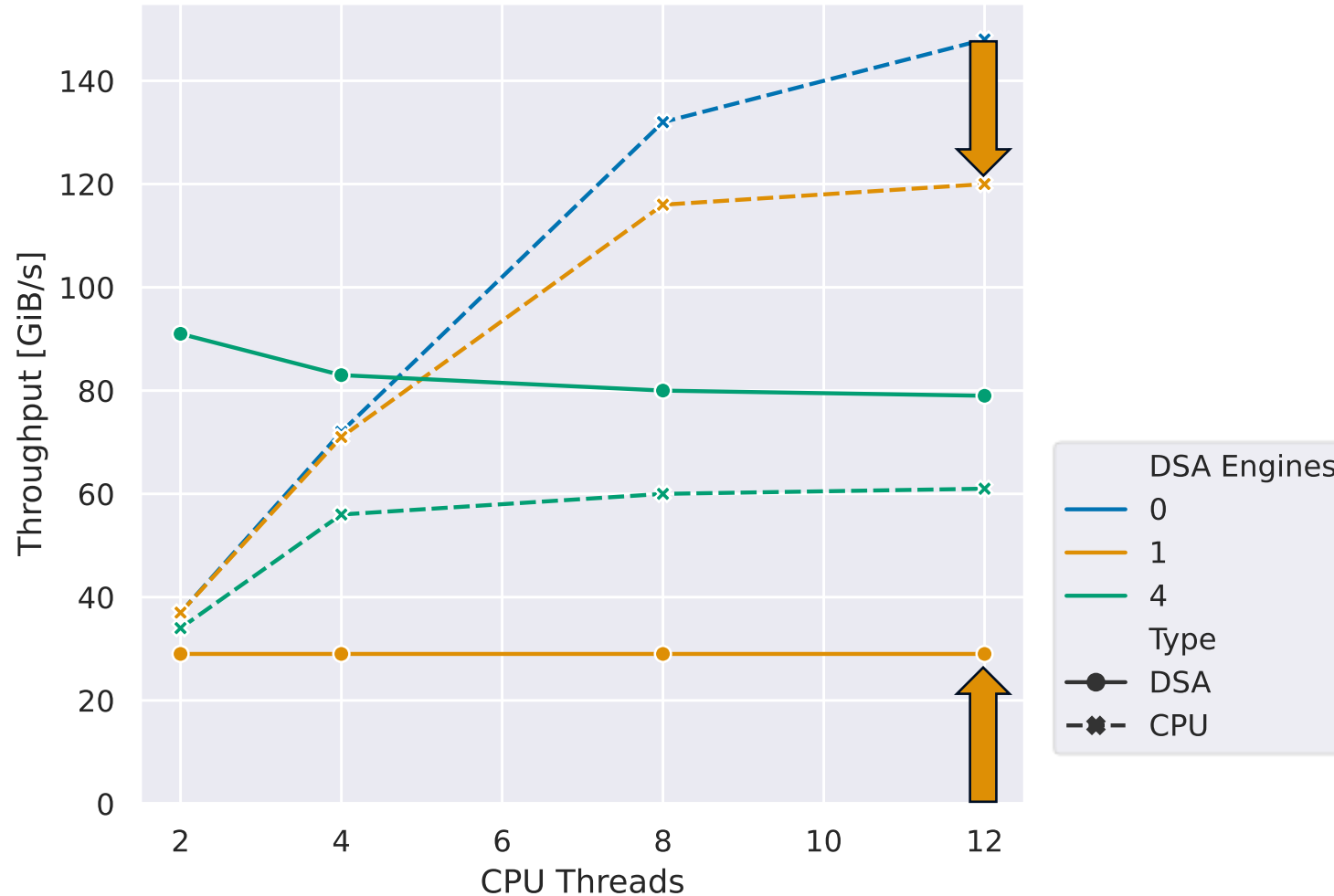## Data-Intensive CPU Threads on HBM



In HBM we see similar effects as in DRAM

# Benchmark – Inter-Socket Data Transfer with DSA

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 46

Funded by

# Benchmark – Inter-Socket Data Transfer with DSA

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 47

Funded by

# Benchmark – Inter-Socket Data Transfer with DSA

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Funded by

# Benchmark – Inter-Socket Data Transfer with DSA



One DSA can achieve much higher throughput than one CPU.

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 49

# Benchmark – Inter-Socket Data Transfer with DSA



No throughput scaling for multiple DSA engines.
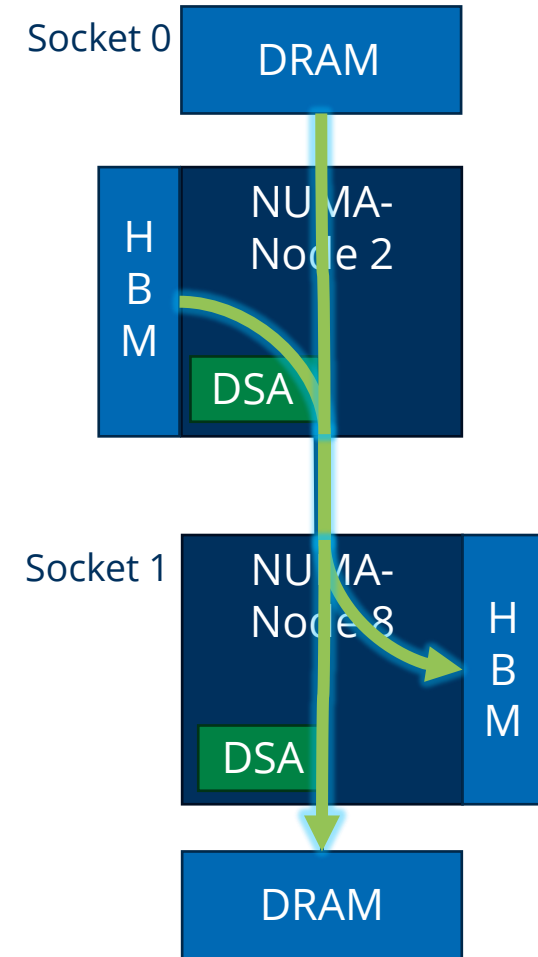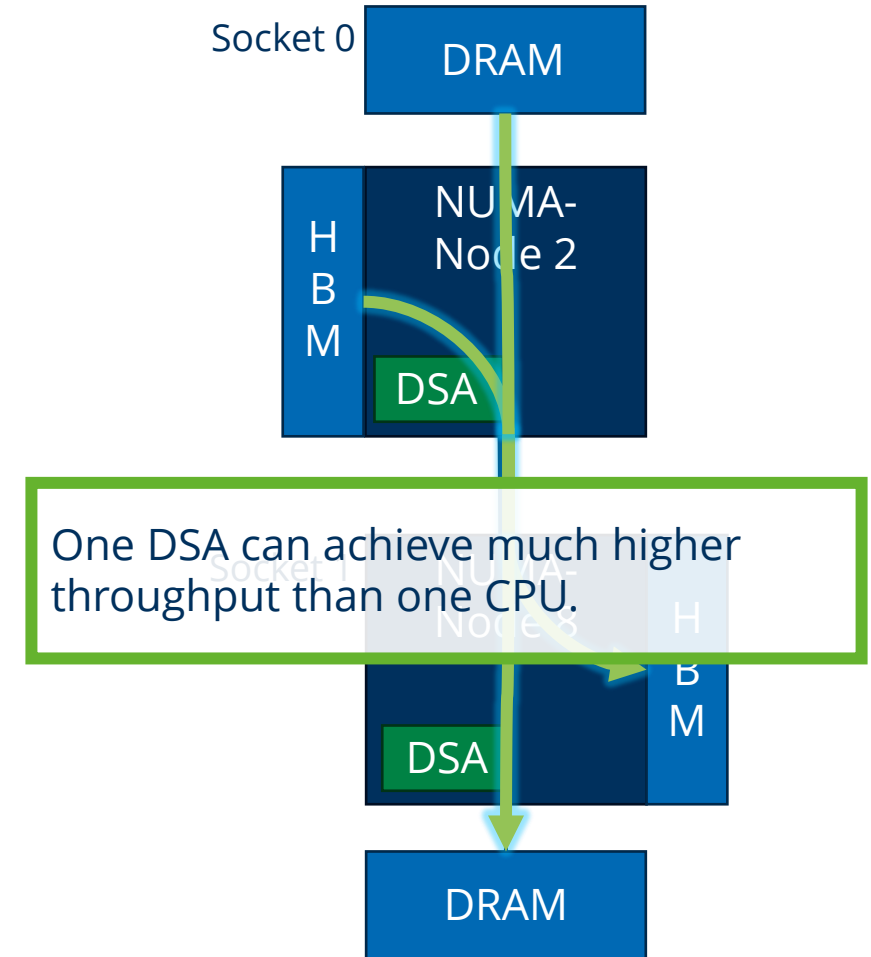
One DSA can achieve much higher throughput than one CPU.
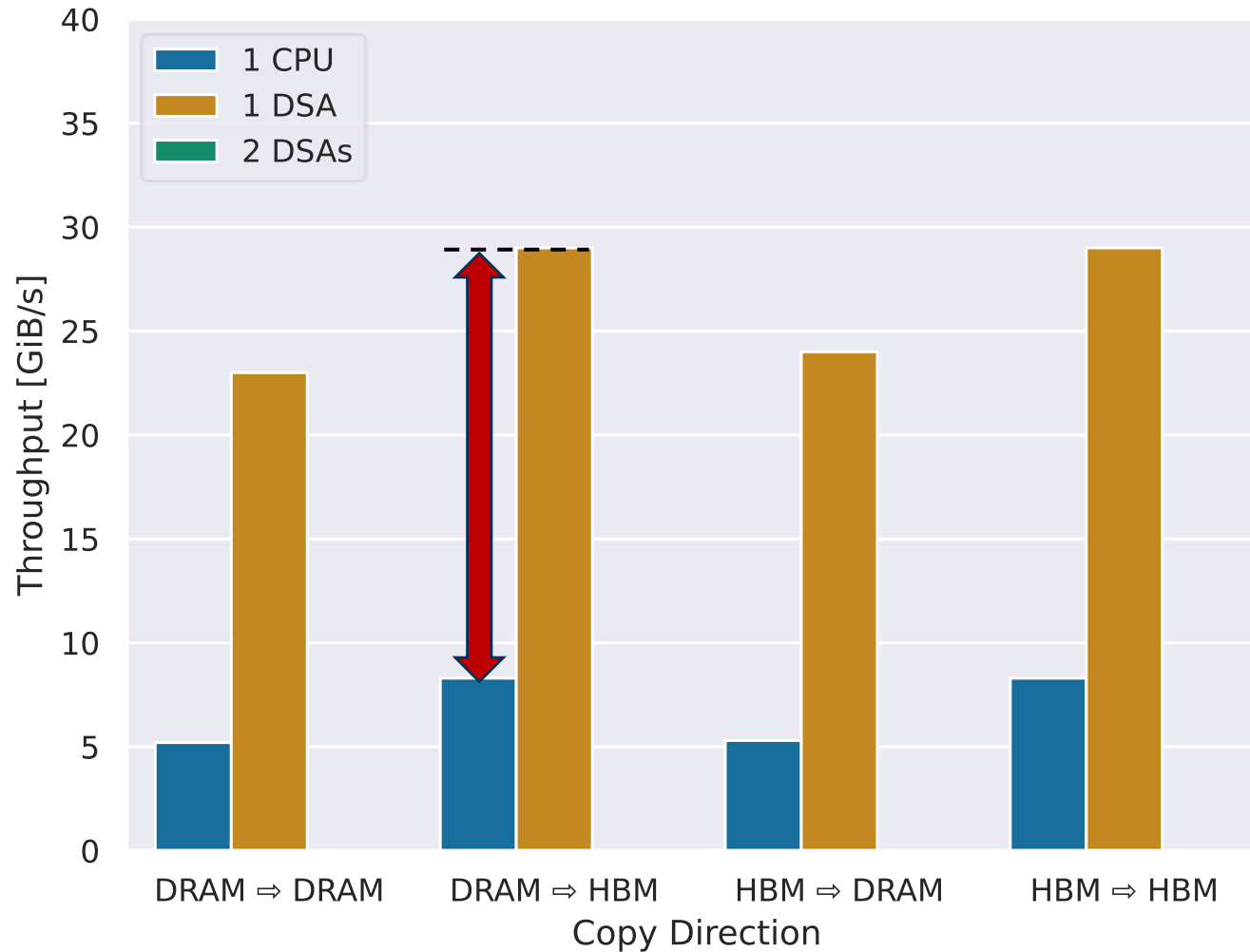
# Benchmark – Inter-Socket Data Transfer with DSA



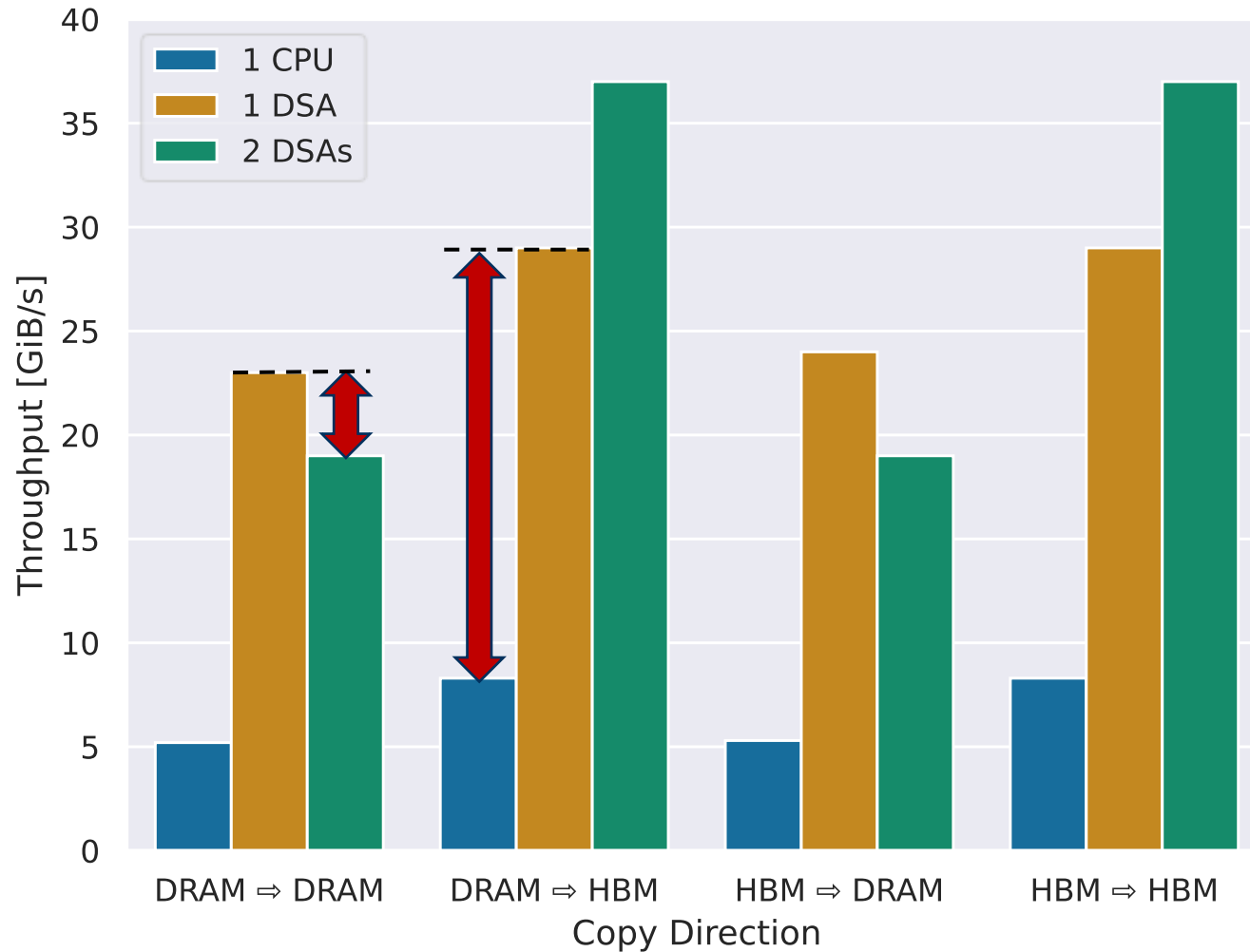Inter-socket transfers incur huge throughput penalties.

No throughput scaling for multiple DSA engines.

One DSA can achieve much higher throughput than one CPU.

# On-the-fly Data Distribution

André Berthold, Lennart Schmidt, Anton Obersteiner, Dirk Habich, Wolfgang Lehner, Horst Schirmeier. 2024. On-The-Fly Data Distribution to Accelerate Query Processing in Heterogeneous Memory Systems. In *28th European Conference on Advances in Databases and Information Systems (ADBIS '24).* Springer Nature Switzerland, Cham, 170–183.

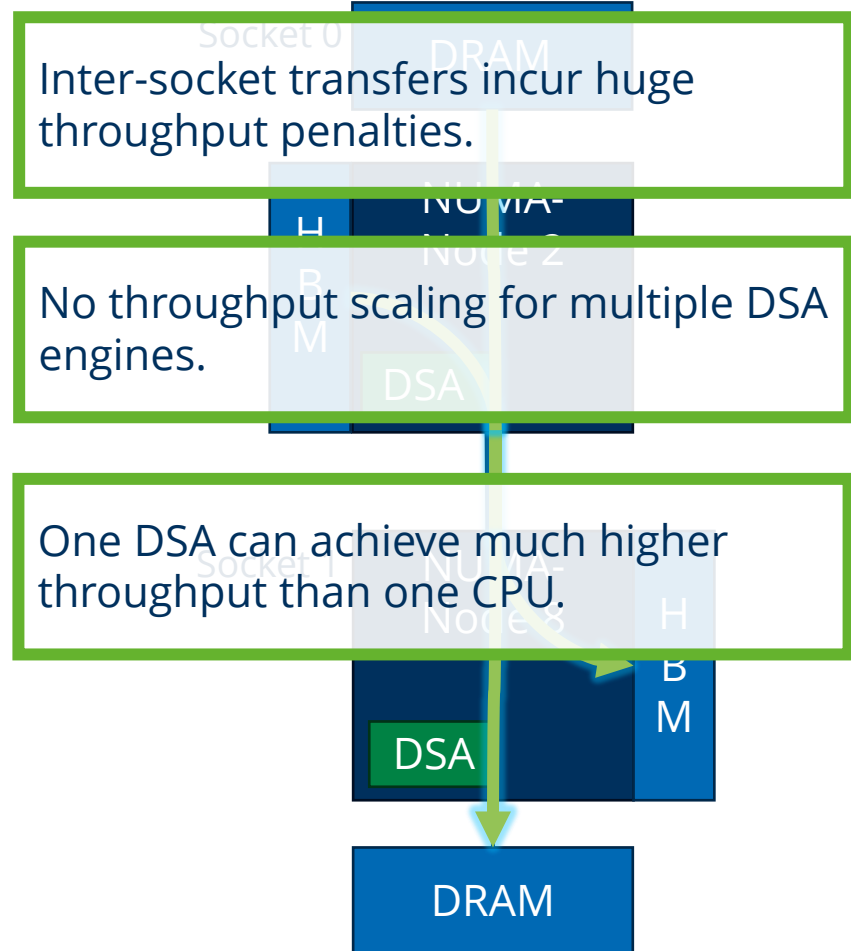Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 52

Funded by

TECHNISCHE UNIVERSITÄT DRESDEN

DFG

DRESDEN concept

# On-the-fly Data Distribution

**Query**

Filter a $\longrightarrow$ Filter b $\longrightarrow$ Sum b

DRAM

Column a    Column b

NUMA-Node    HBM

# On-the-fly Data Distribution
## Baseline Execution Time

**Query**

Filter a ⟶ Filter b ⟶ Sum b

Throughput | #Threads

0     1tu     2tu     3tu    Execution Time

DRAM

Column a     Column b

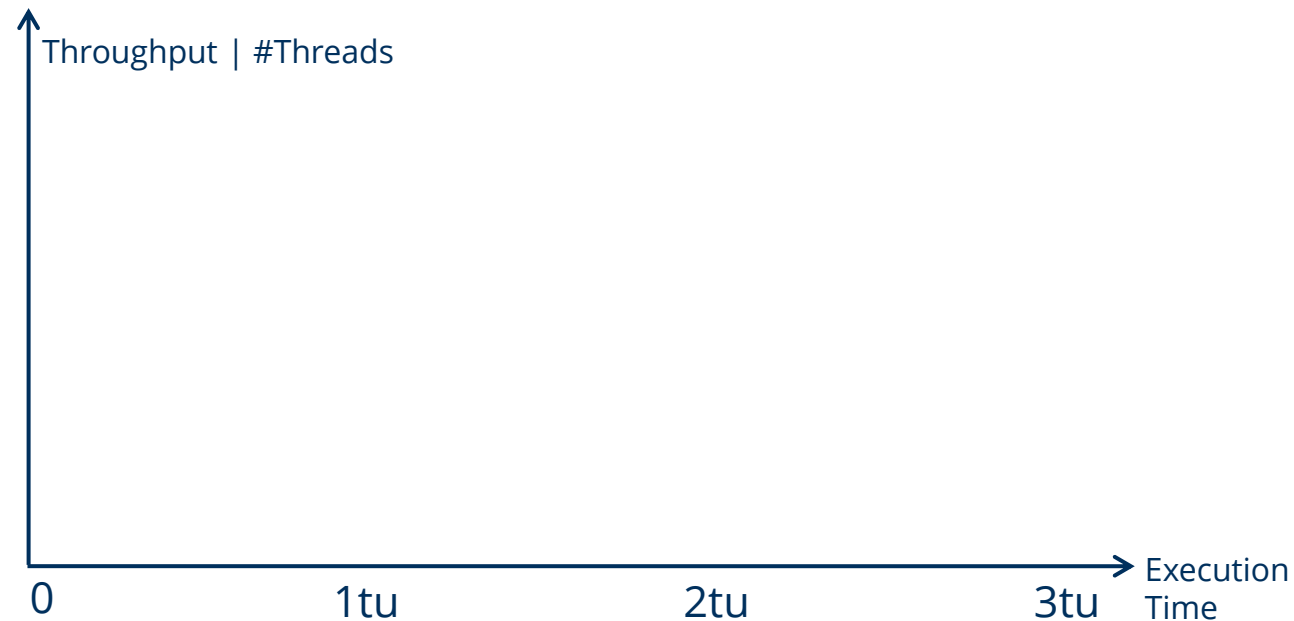NUMA-Node     HBM

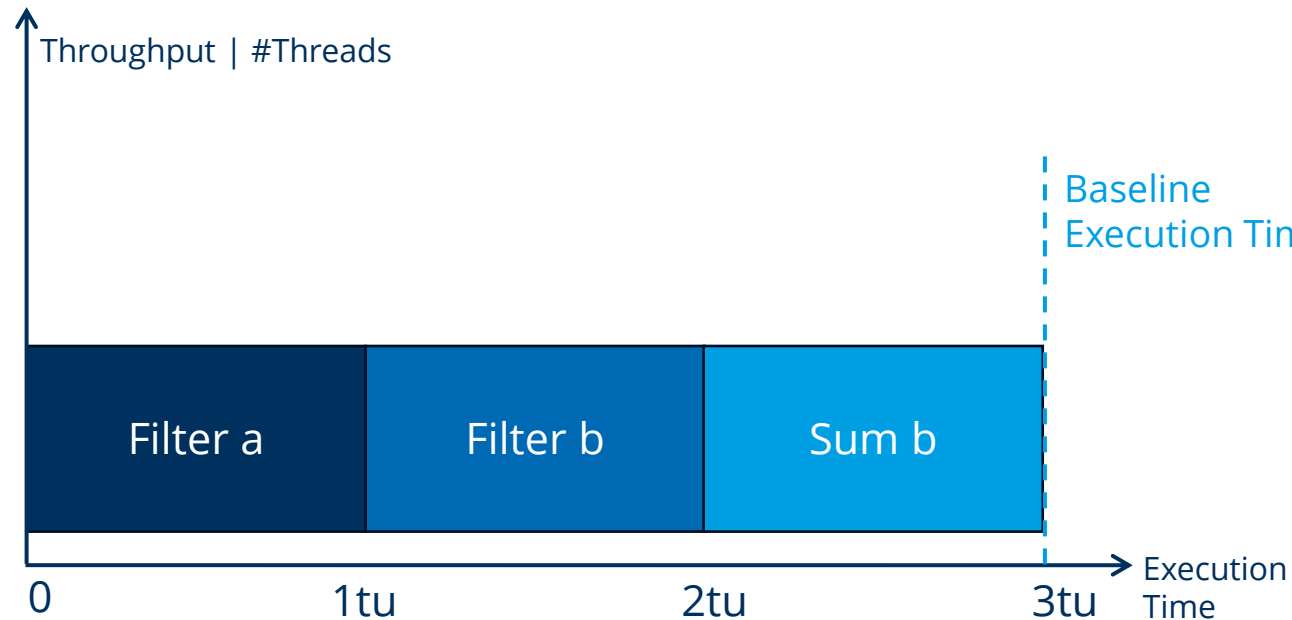TECHNISCHE UNIVERSITÄT DRESDEN

DFG

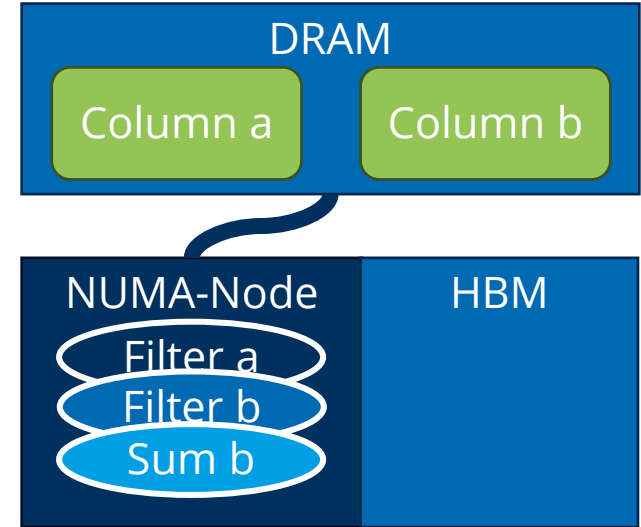DRESDEN concept

# On-the-fly Data Distribution
## Baseline Execution Time

**Query**

Filter a → Filter b → Sum b

Throughput | #Threads

Baseline Execution Time

| Filter a | Filter b | Sum b |

0   1tu   2tu   3tu   Execution Time

DRAM

Column a      Column b

NUMA-Node      HBM

Filter a
Filter b
Sum b

# On-the-fly Data Distribution
## Optimized Execution Time

**Query**

Filter a ⟶ Filter b ⟶ Sum b

Throughput | #Threads

Baseline Execution Time

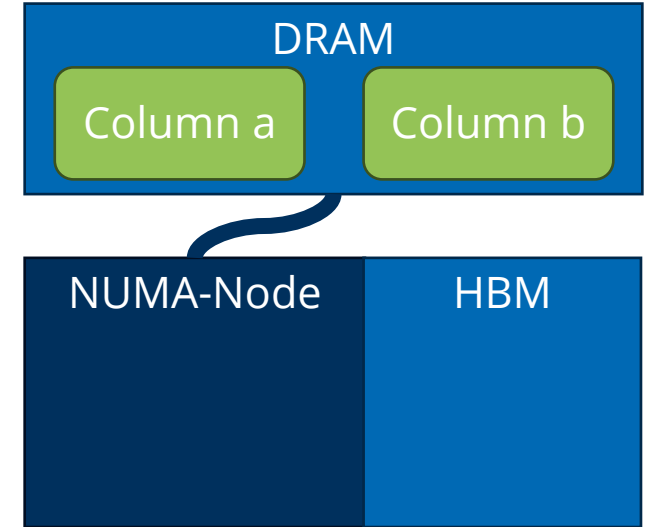0   1tu   2tu   3tu   Execution Time

DRAM

Column a   Column b

NUMA-Node   HBM

# On-the-fly Data Distribution
## Optimized Execution Time

**Query**

Look into the **future**

Filter a ⟶ Filter b ⟶ Sum b

DRAM

Column a    Column b

NUMA-Node    HBM

Throughput | #Threads

Baseline
Execution Time

0    1tu    2tu    3tu    Execution Time
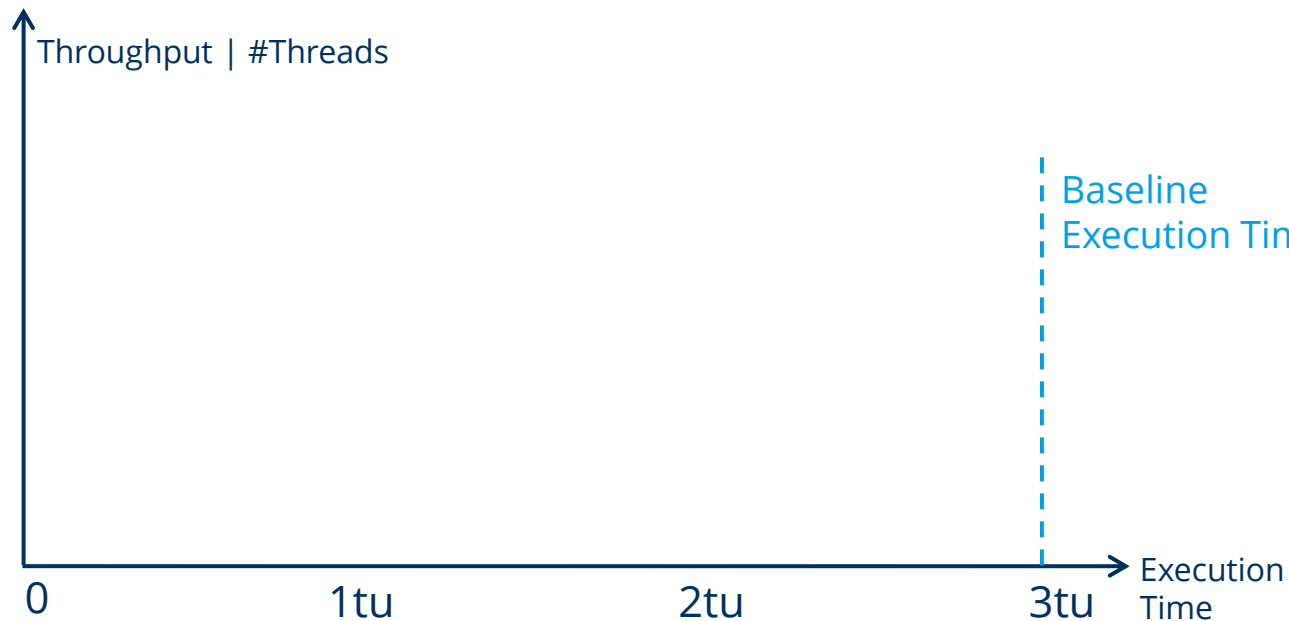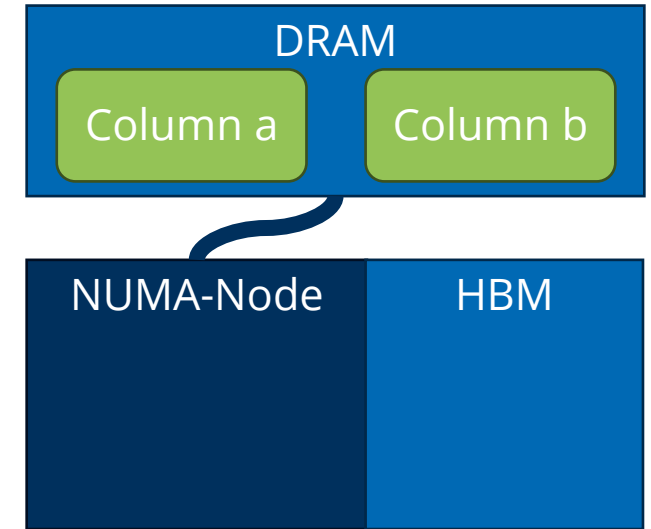
# On-the-fly Data Distribution
## Optimized Execution Time

**Query**

Look into the **future**

Filter a → Filter b → Sum b

Throughput | #Threads

Baseline
Execution Time

0    1tu    2tu    3tu    Execution Time

DRAM

Column a    Column b

Copy

NUMA-Node    HBM

Filter a

Column b

TECHNISCHE UNIVERSITÄT DRESDEN
Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
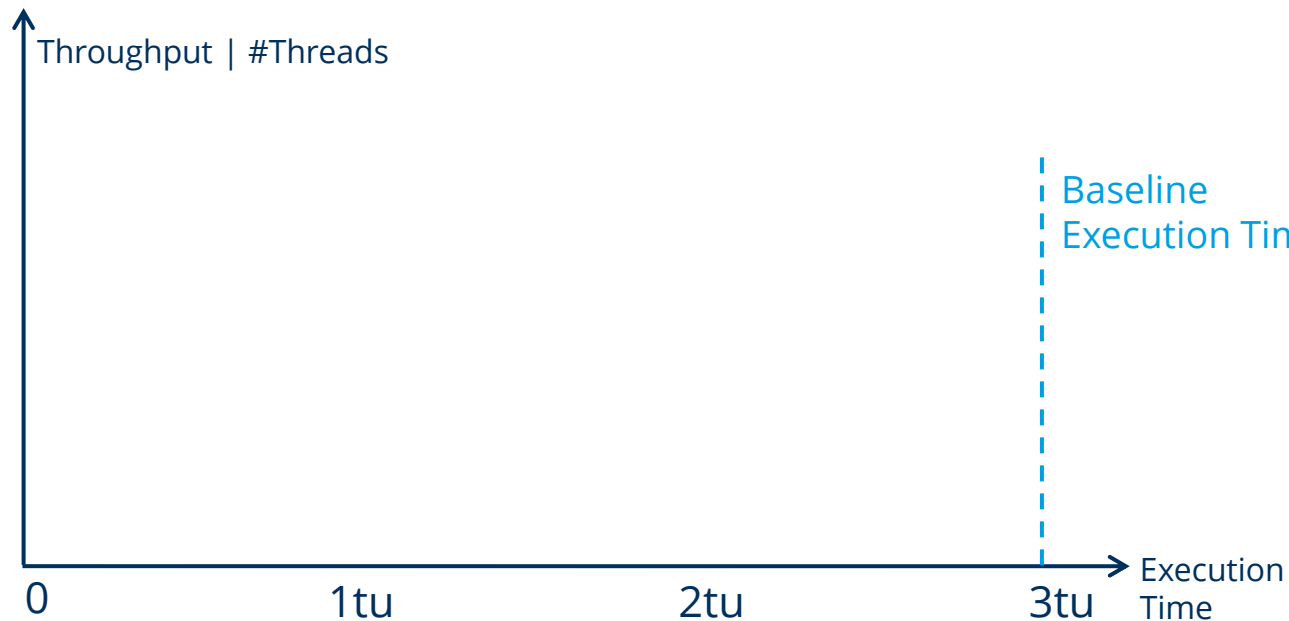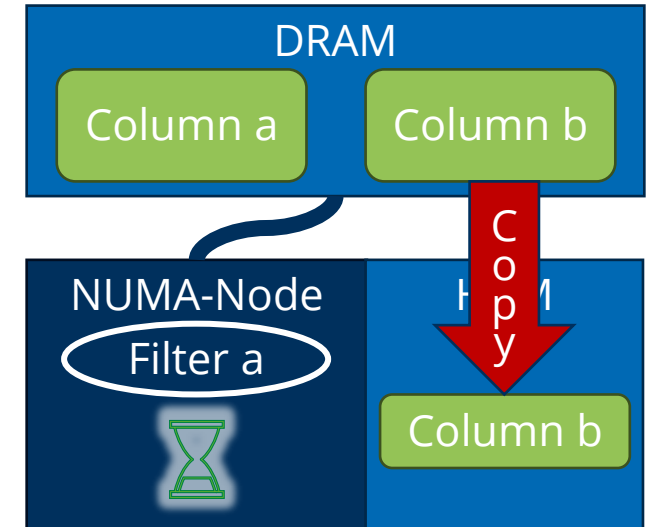Dresden University of Technology / André Berthold

DFG    DRESDEN concept

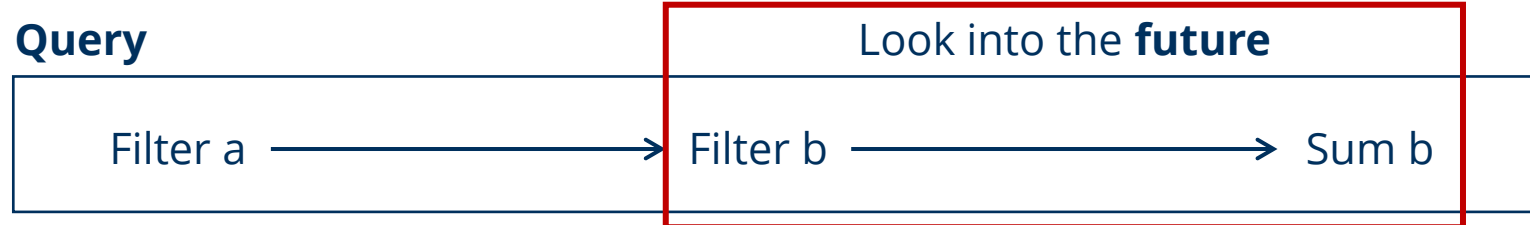# On-the-fly Data Distribution
## Optimized Execution Time

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Funded by

# On-the-fly Data Distribution
## Optimized Execution Time

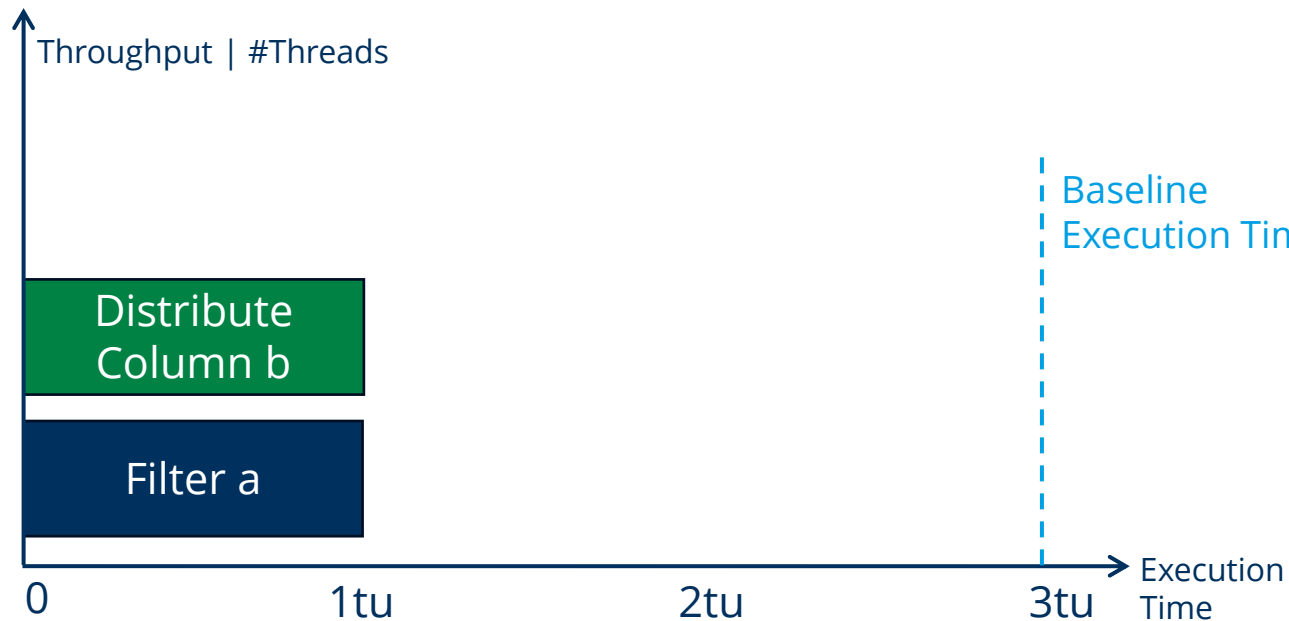**Query**
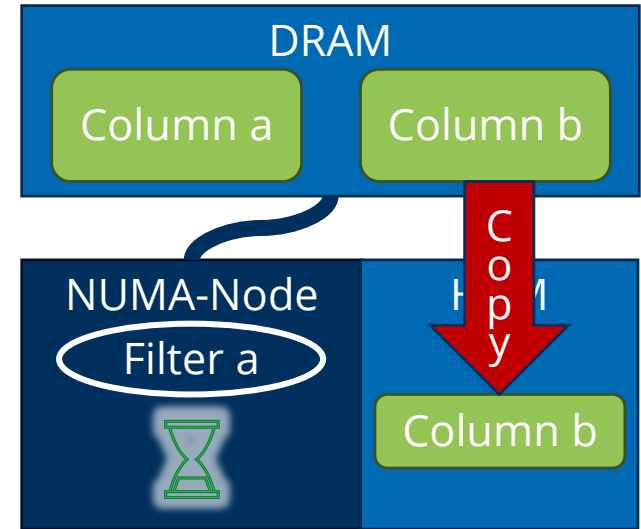
Filter a → Filter b → Sum b

DRAM
Column a   Column b

NUMA-Node   HBM
Filter b
Sum b   Column b

Throughput | #Threads

Baseline
Execution Time

Distribute Column b

Filter a

0      1tu      2tu      3tu      Execution Time

# On-the-fly Data Distribution
## Optimized Execution Time

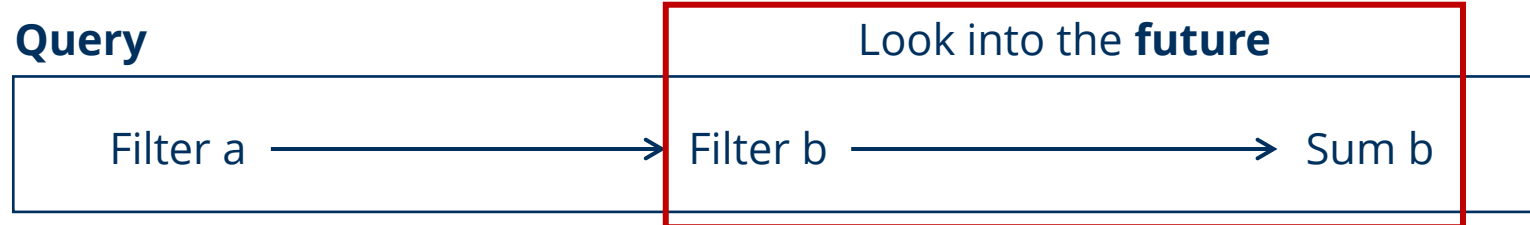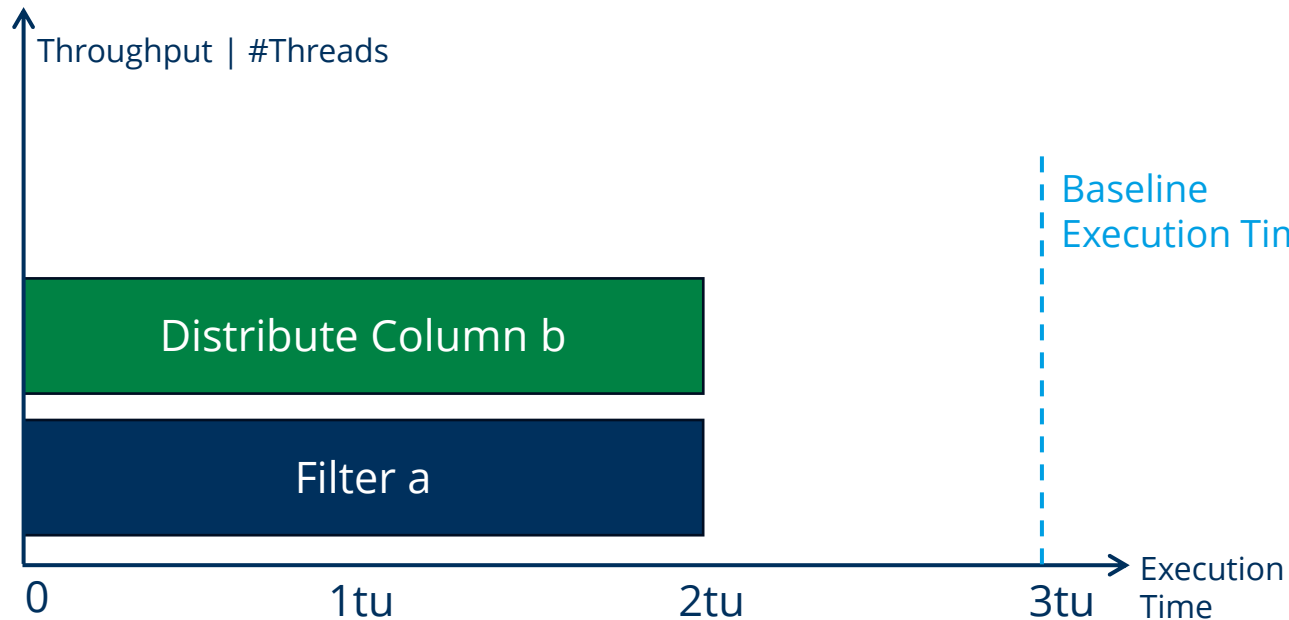**Query**

Filter a → Filter b → Sum b

Throughput | #Threads

Baseline Execution Time

Distribute Column b

Filter a

0          1tu          2tu          3tu          Execution Time

DRAM

Column a          Column b

NUMA-Node          HBM

Filter b

Sum b          Column b

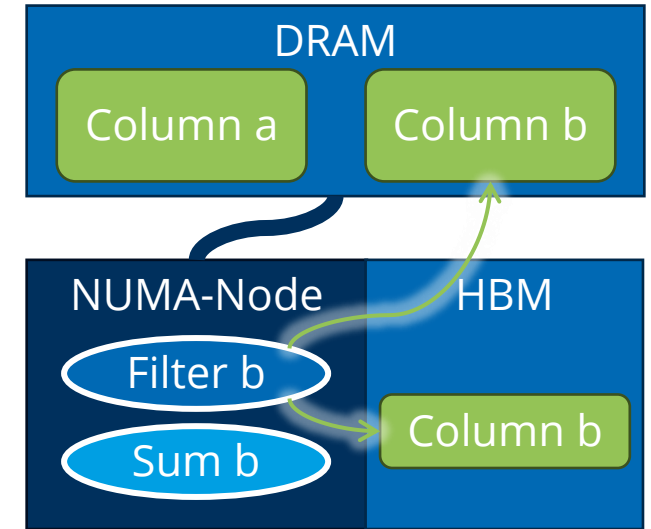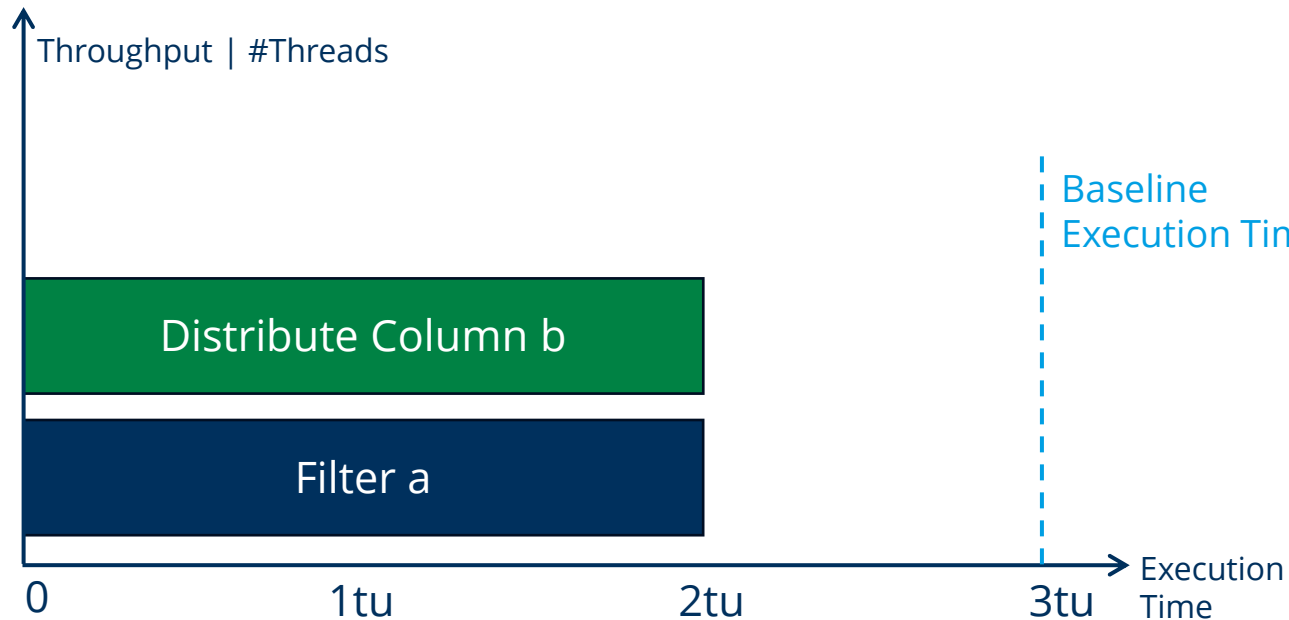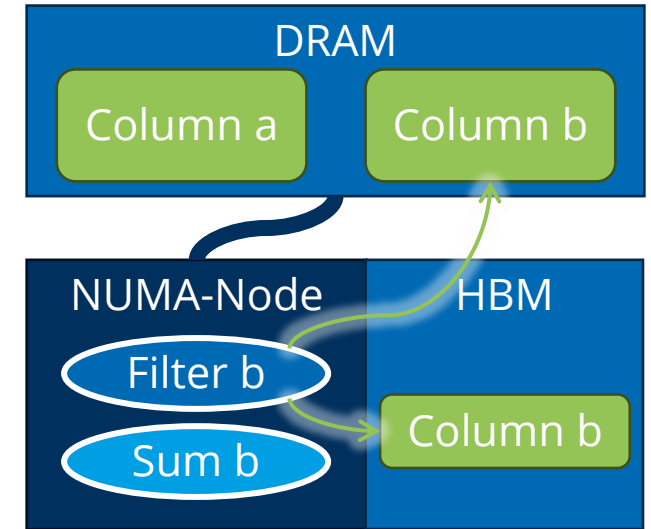TECHNISCHE UNIVERSITÄT DRESDEN

DFG          DRESDEN concept

# On-the-fly Data Distribution
## Optimized Execution Time

**Query**

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold
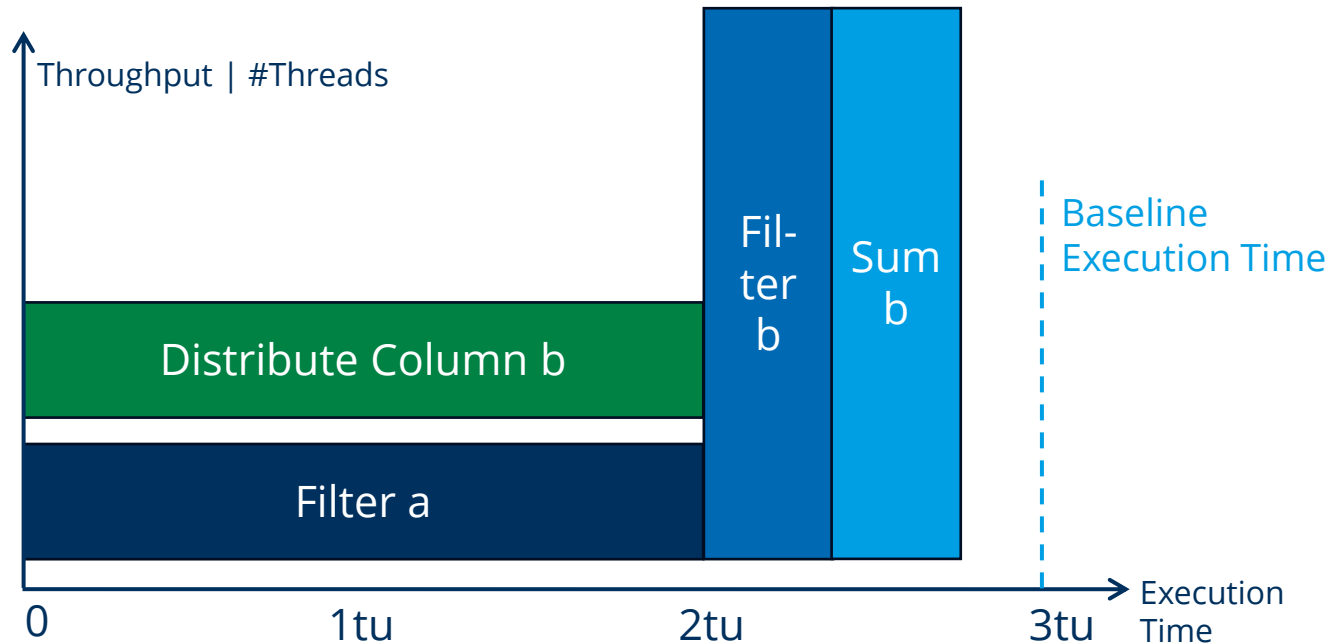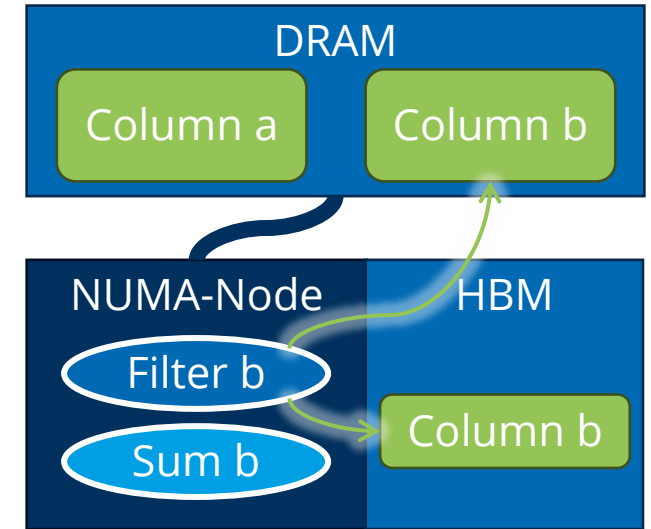
Funded by

# On-the-fly Data Distribution
## Optimized Execution Time



**Query**

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 63
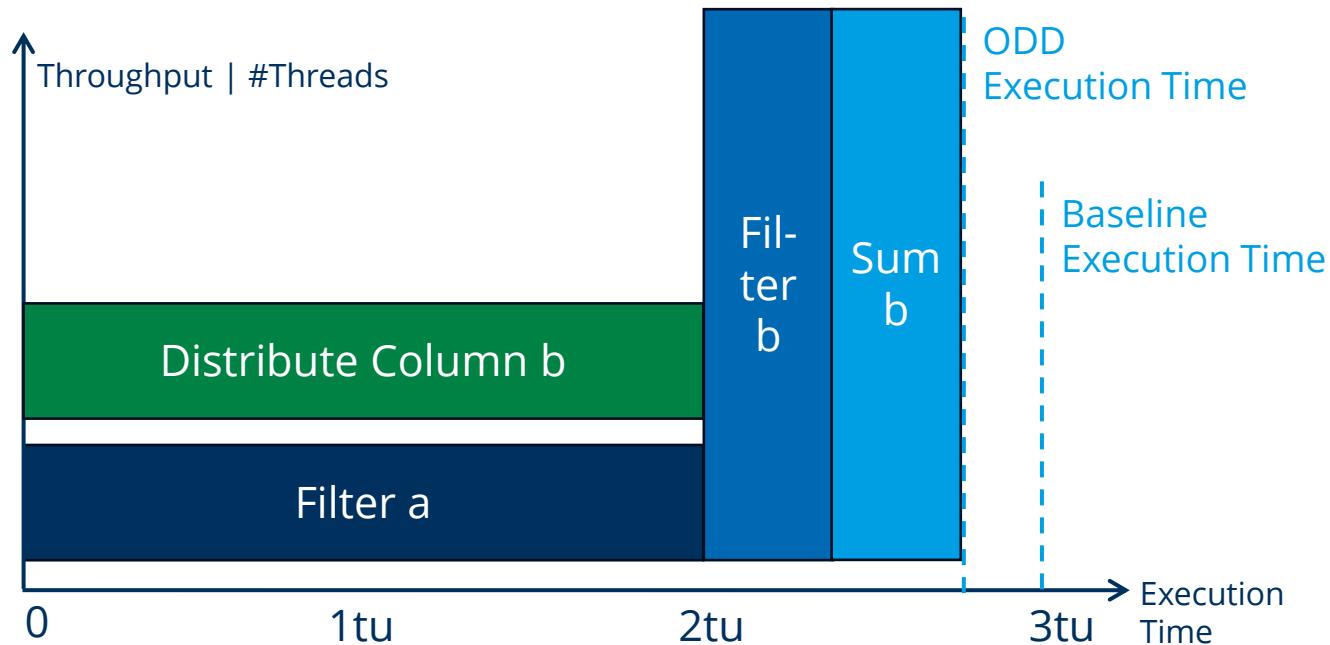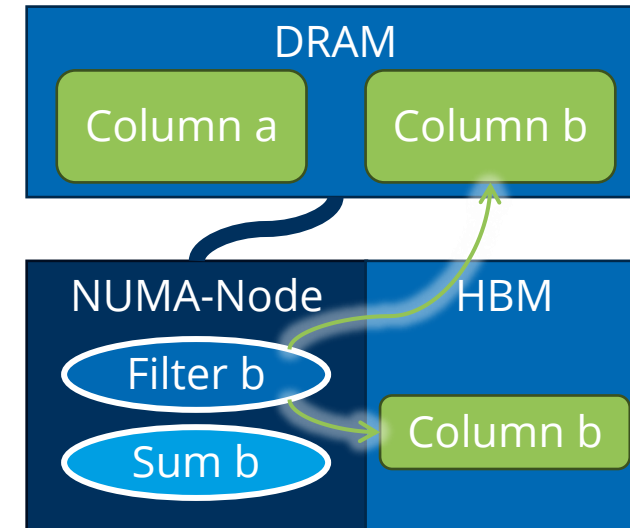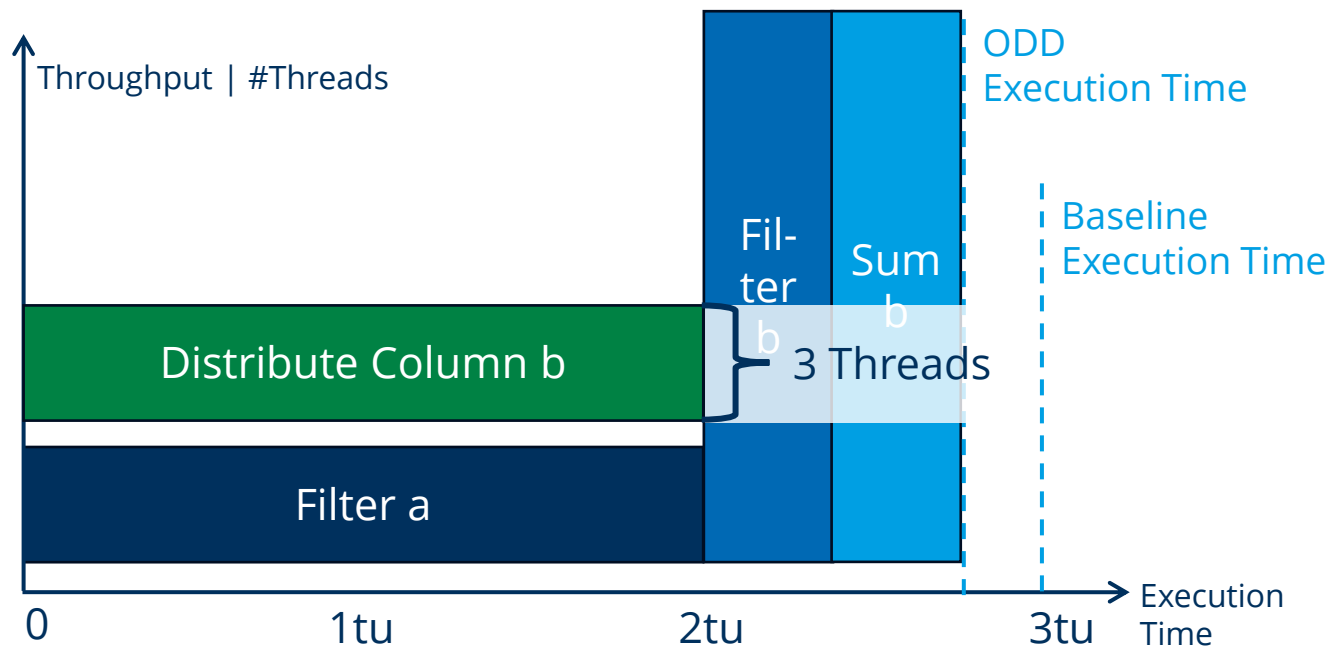
Funded by

# On-the-fly Data Distribution
## Optimized Execution Time

**Query**

Filter a → Filter b → Sum b



Throughput | #Threads

Distribute Column b

Filter a

Fil-ter b

Sum b

3 Threads

ODD Execution Time

Baseline Execution Time

0  1tu  2tu  3tu  Execution Time

DRAM

Column a   Column b

NUMA-Node   HBM

Filter b

Sum b

Column b

# On-the-fly Data Distribution
## Optimized using DSA

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
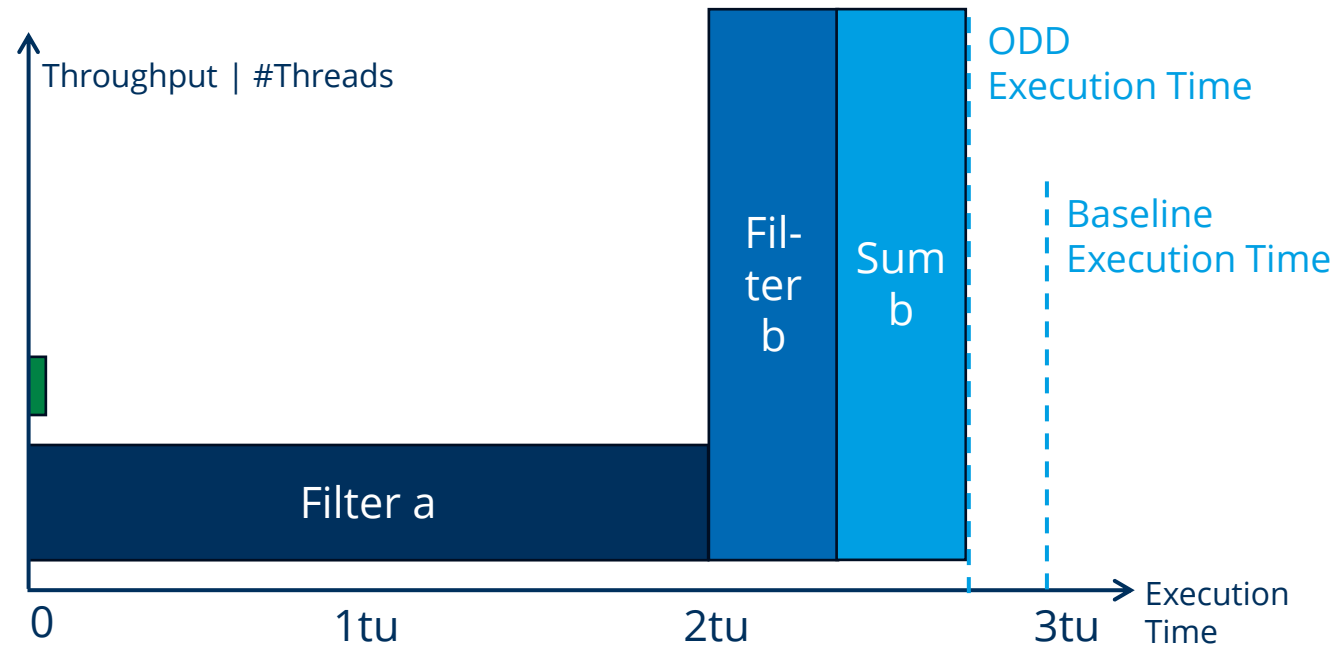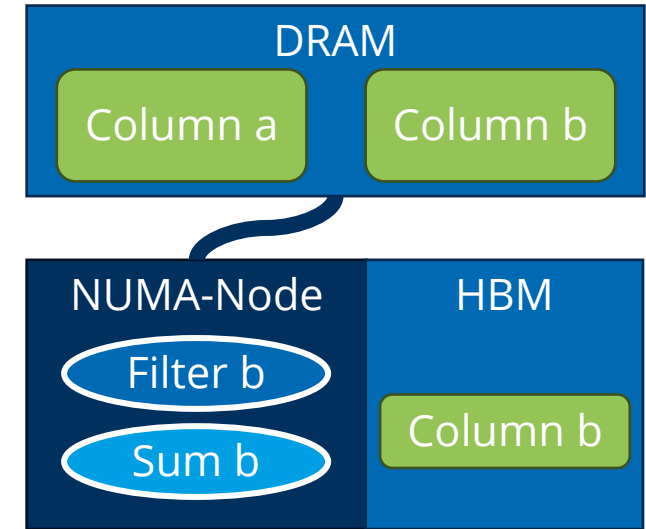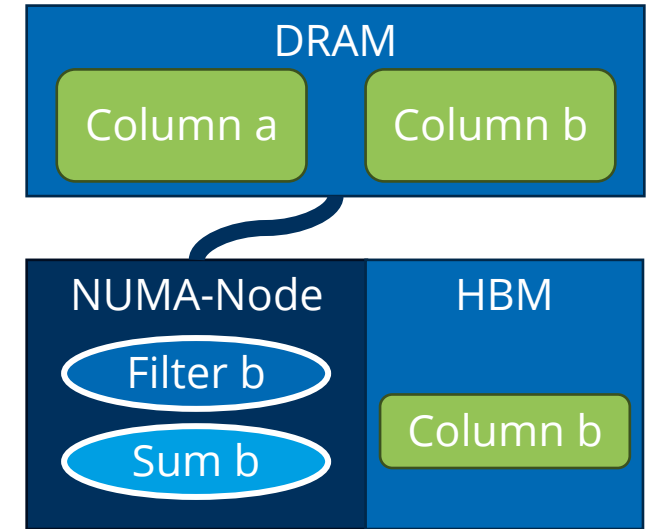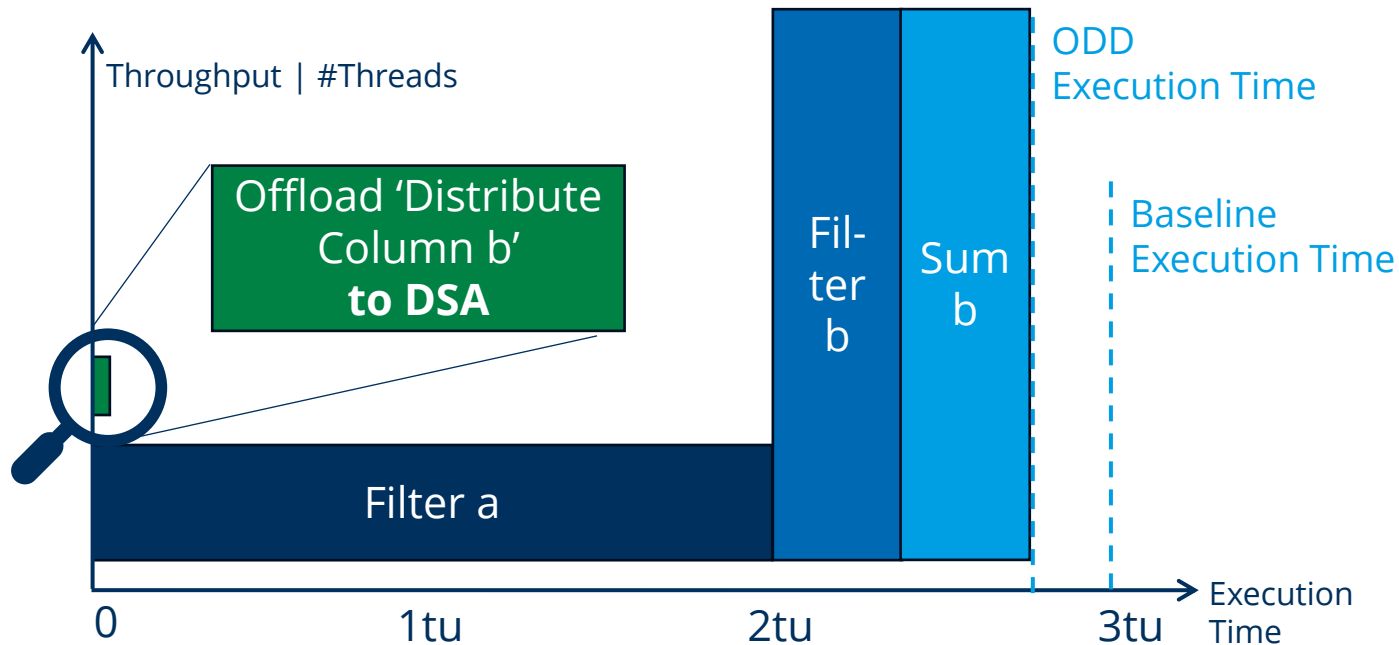Dresden University of Technology / André Berthold

Funded by

# On-the-fly Data Distribution
## Optimized using DSA

**Query**

Filter a ⟶ Filter b ⟶ Sum b

Throughput | #Threads

Offload 'Distribute Column b' **to DSA**

Filter a

Fil-ter b

Sum b

ODD Execution Time

Baseline Execution Time

0    1tu    2tu    3tu    Execution Time

DRAM

Column a    Column b

NUMA-Node    HBM

Filter b

Sum b    Column b

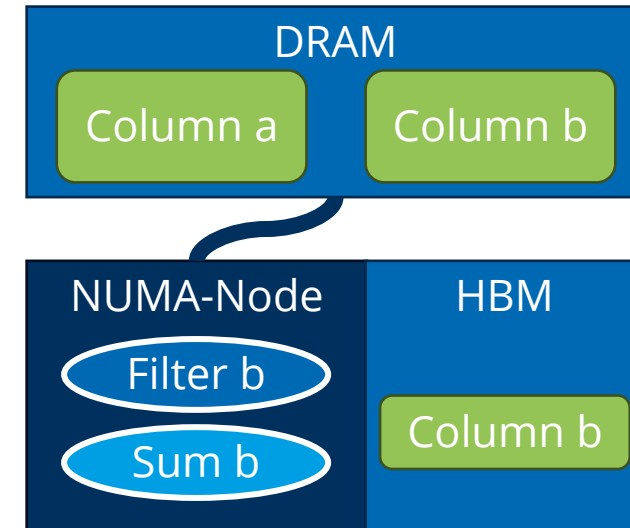# On-the-fly Data Distribution
## Optimized using DSA

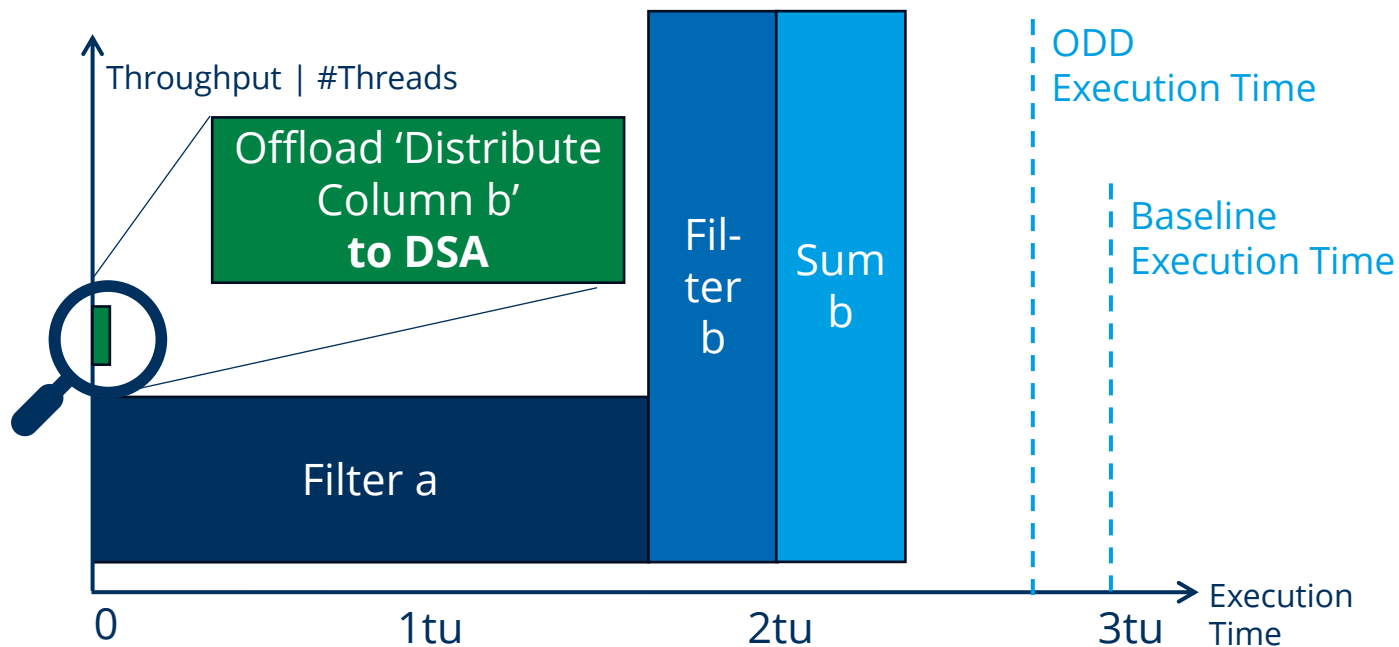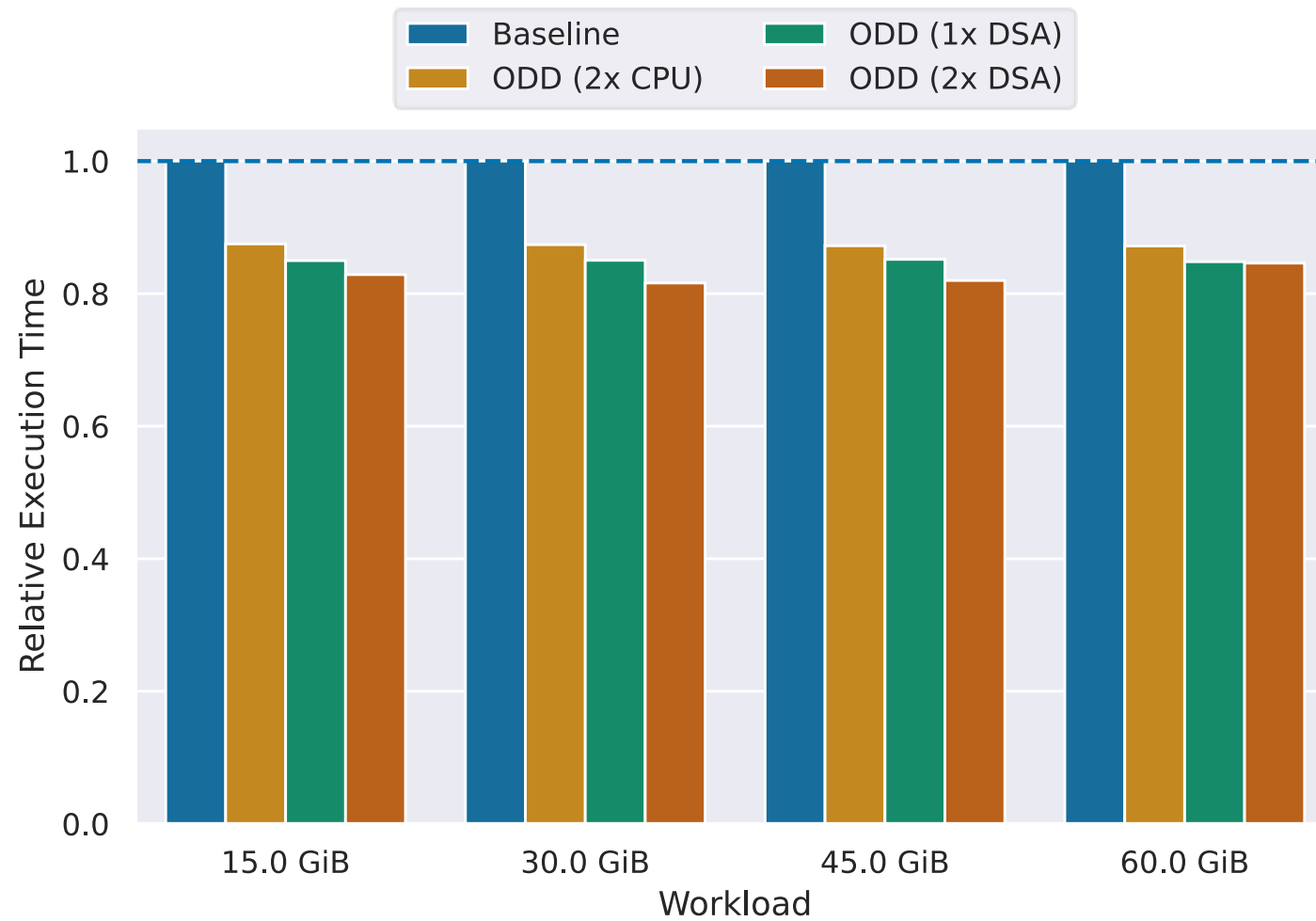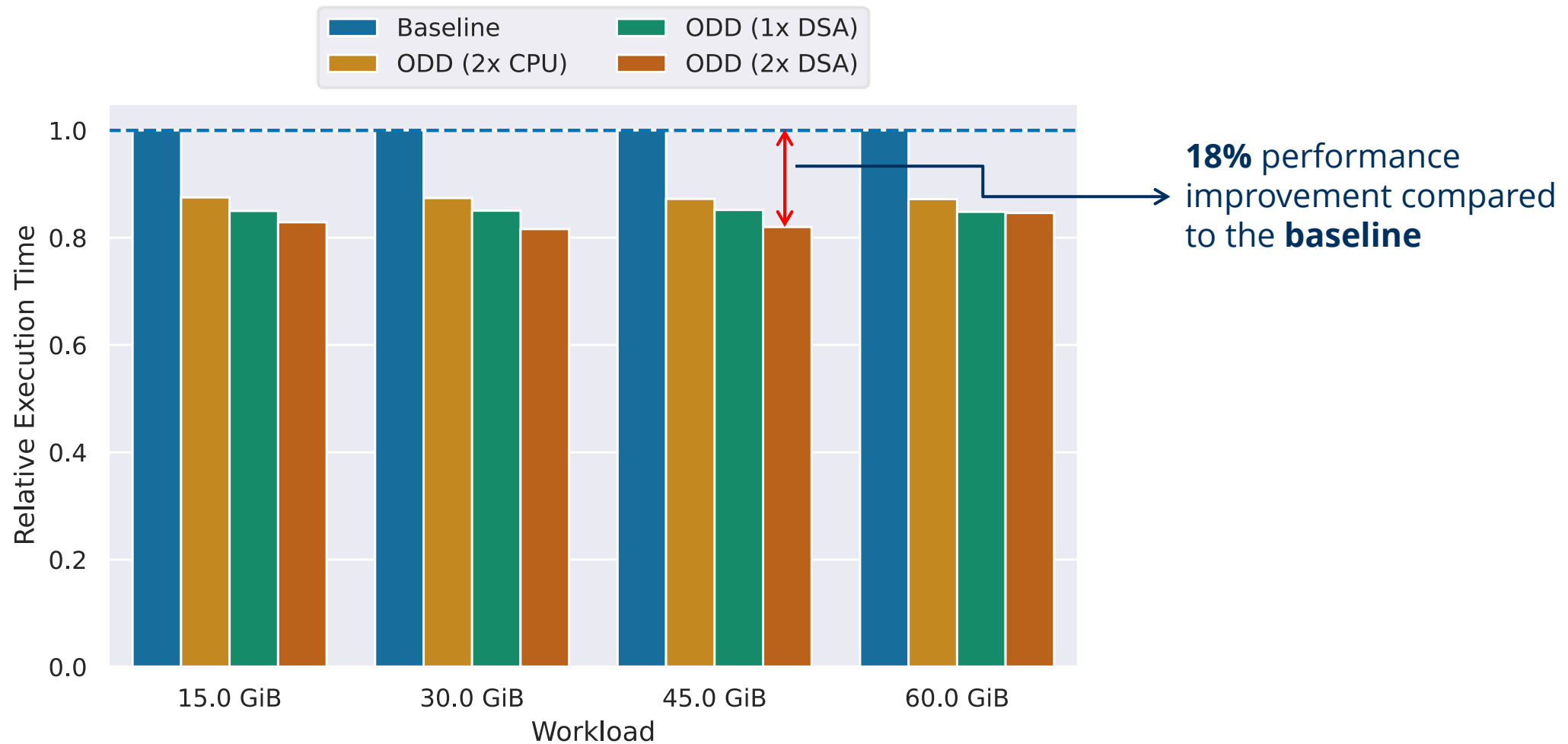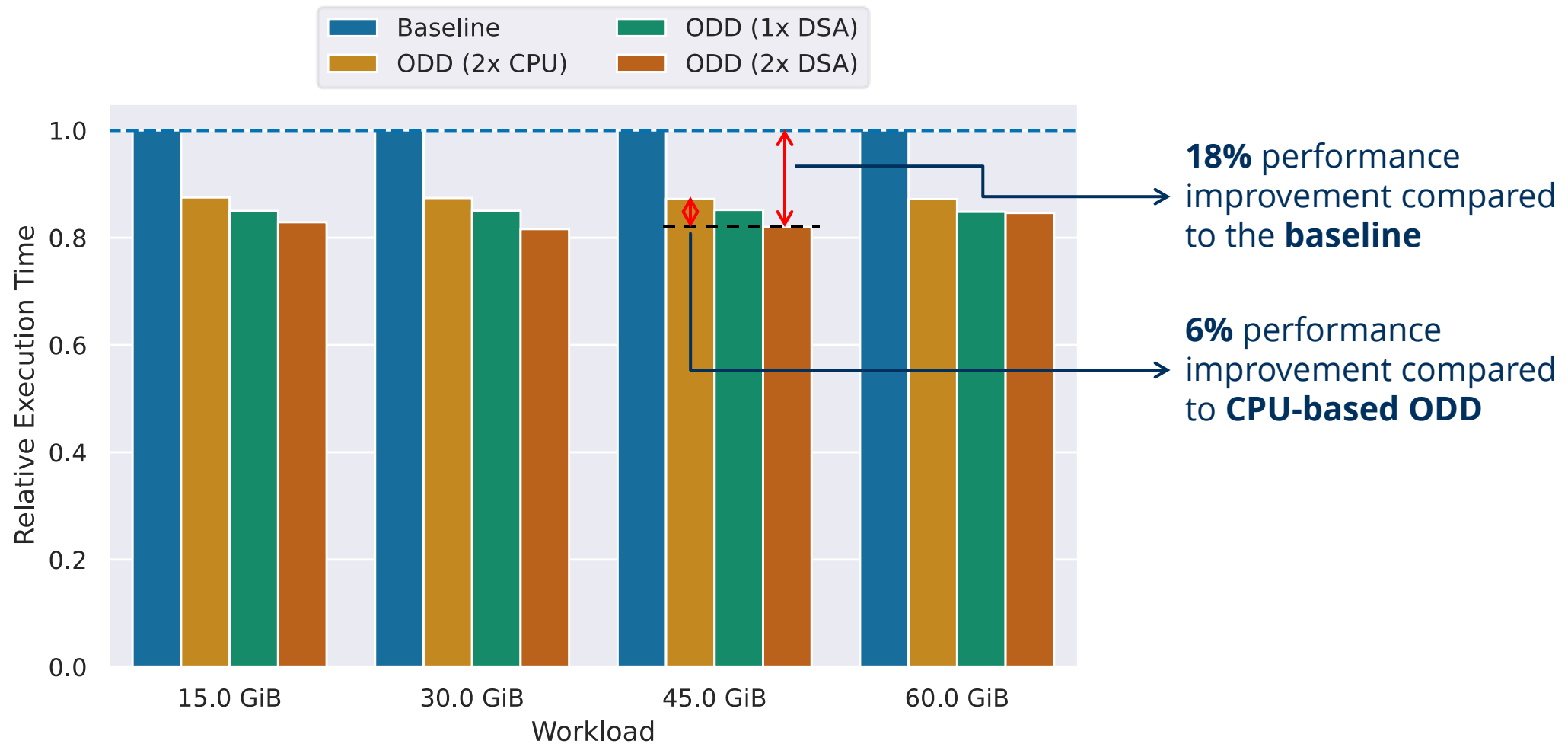Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

# On-the-fly Data Distribution with DSA - Benchmark

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 68

Funded by

# On-the-fly Data Distribution with DSA - Benchmark



18% performance improvement compared to the baseline

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Funded by

# On-the-fly Data Distribution with DSA - Benchmark



**18%** performance improvement compared to the **baseline**

**6%** performance improvement compared to **CPU-based ODD**

# Summary

DSA allows memory operation offloading …

Demystifying Intel Data Streaming Accelerator for In-Memory Data Processing
Dresden University of Technology / André Berthold

Slide 71

Funded by

# Summary

DSA allows memory operation offloading ...



... with **higher throughput** compared to one CPU.

# Summary

DSA allows memory operation offloading ...



... with **higher throughput** compared to one CPU.

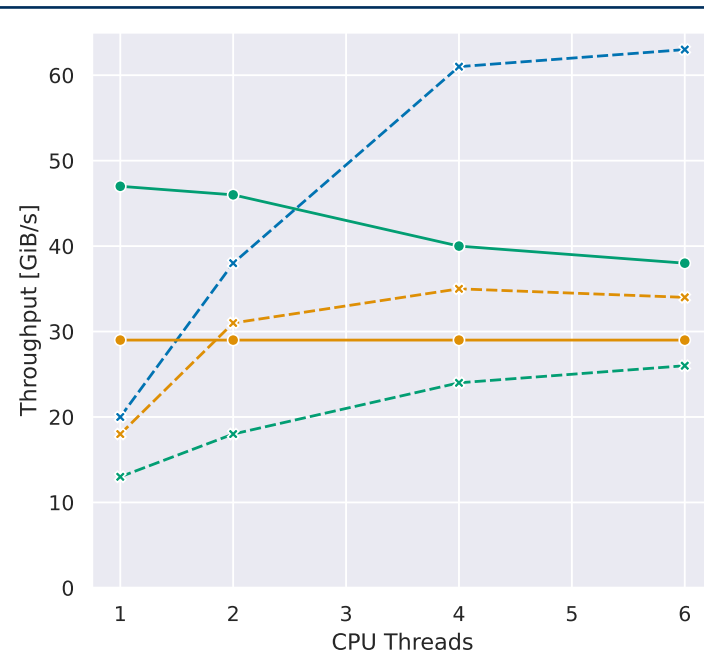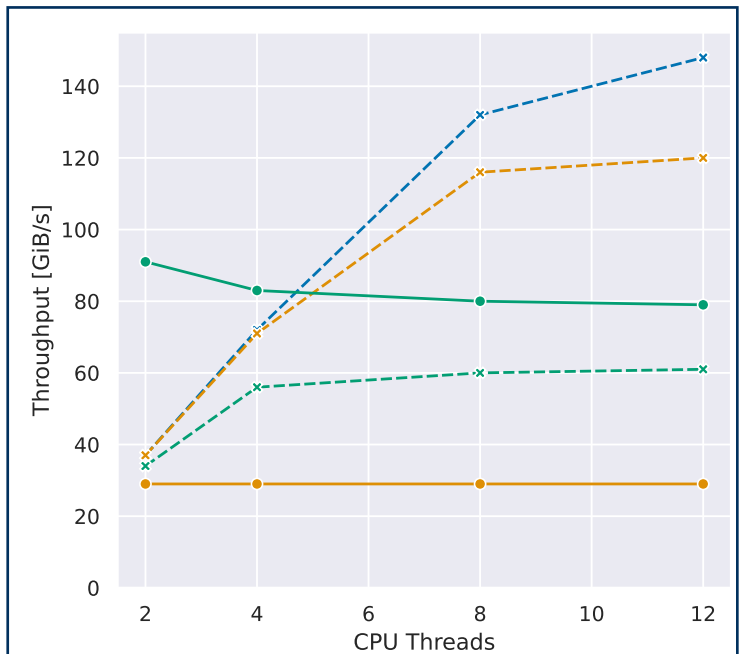... with **high-priority memory access**

# Summary

DSA allows memory operation offloading ...



... with **higher throughput** compared to one CPU.

... with **high-priority memory access**

... where DSA interferes with CPU,

but CPU has **minor influence on DSAs** throughput.