TUП

Near Memory Accelerators for Efficient Inter-Tile Communication
in Distributed-Shared-Memory Architectures

Andreas Herkersdorf, Sven Rheindt,
Akshay Srivatsa, Thomas Wild
Technical University of Munich

MPSoC Forum 2019

---

TUП

## Big Picture: Hitting yet another wall!
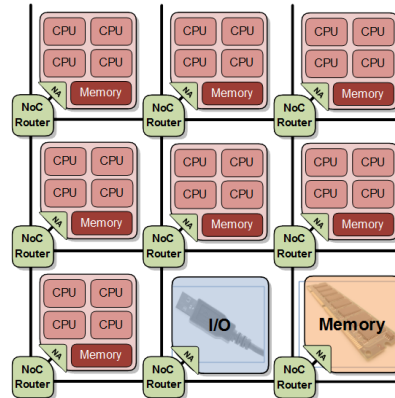
TUT

## The Walls of Computer Architecture

- 1st Wall – Mid 90s: the Memory Wall
- 2nd Wall – 2004: the Power Wall

Steps taken:
- Multi-/Many Core, Cache Hierarchies
- NoCs and Distributed Computing with application-specific accelerators

**3rd Wall – now: the Locality Wall [1]**
- Memory intensive, but cache-unfriendly characteristics
- Dominated by data access & movement



3

TUT

## How to break the locality wall?

**Increase Data-to-Task Locality**

- Data Migration
- Task Migration
- Near Memory Computing

**Research Questions:**

- Architecture of the future
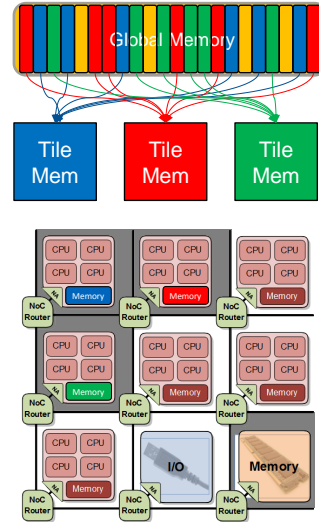- Best approach/combination
- Programming paradigm



4

TIM

# MPSoC 2018: Region-based Coherence [2,3]

Non-uniform memory accesses (NUMA)

Data placement influences performance
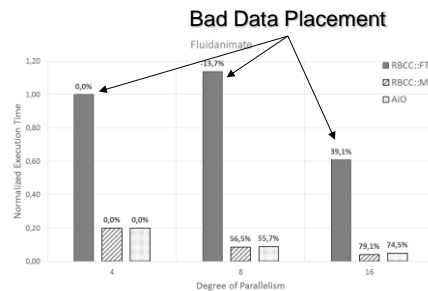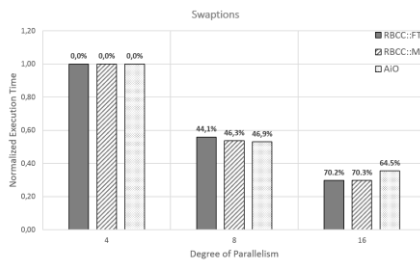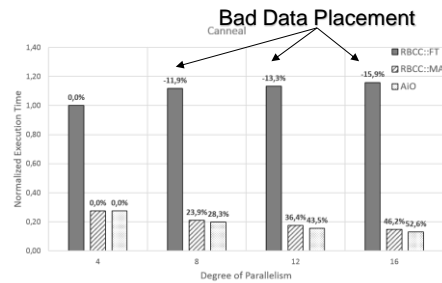→ Analyze the impact of data placement
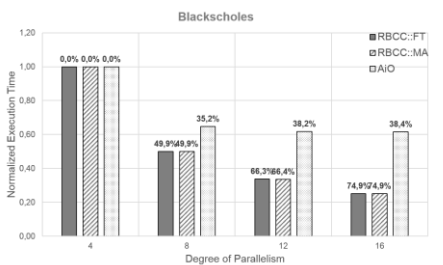
Placement Algorithms
- First Touch Policy
  - Place data into tile of first access

- Most Accessed Policy
  - Place data to preferred tile
  - Known after task execution, but exploitable for periodic task invocation
  - Maximize local accesses
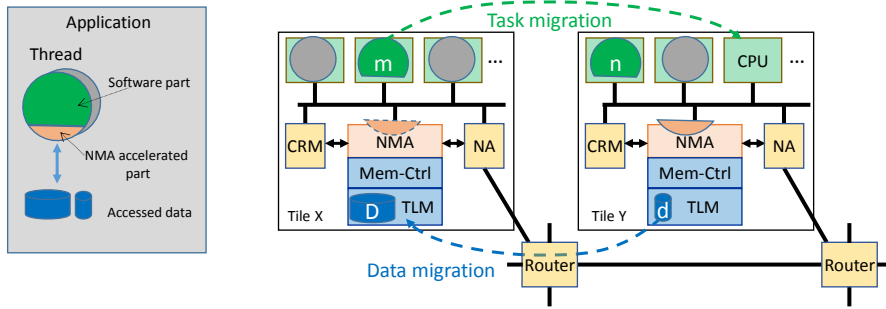


5

TIM

# Impact of Data Placement



**Bad Data Placement**

**Bad Data Placement**

6

3

Пॖ

# Near Memory Acceleration (NMA)

**Topical research subject**
- Eliminating / reducing CPU-to-memory accesses via local interconnect for memory-intense sub-functions
- Mostly addressing accelerator-centric 3D-stacked memory systems [4] [5]

**Our scope**
- Role of NMA in distributed-shared memory architectures
- **Implies necessity for considering task and data placement / migration**
- **SHARQ & Graphcopy**
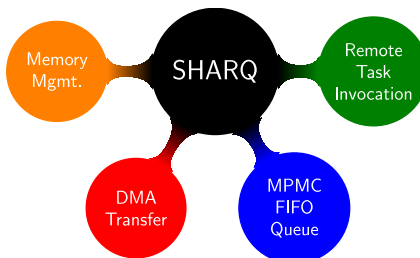


7

Пॖ

# SHARQ – An Overview [6]

Software-Defined Hardware-Managed Queues

Combine **flexibility of software** queues with **performance of hardware** acceleration
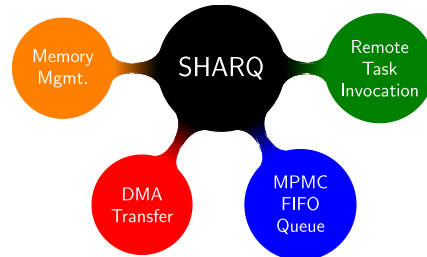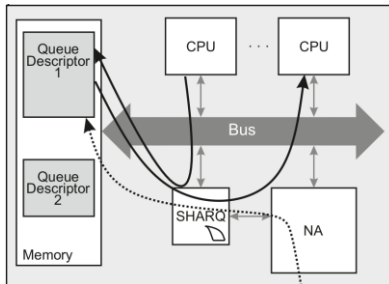
**Software-Defined**

Dynamic allocation and definition of arbitrary sized queues

**Hardware-Managed**



8

ТЛП

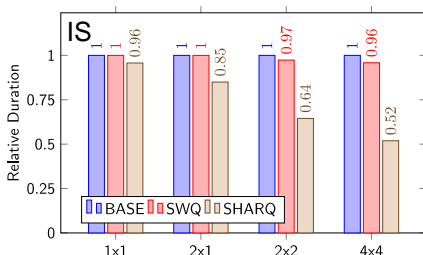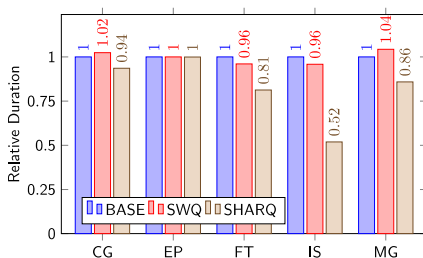## Selected SHARQ Features





**Queue & Memory Management**
- SHARQ access to queue descriptor
- Separation into allocator stack & bounded buffer
- No software involvement / up-calls

**Remote Task Invocation**
- Ensure processing of elements
- Scheduled by SHARQ on demand
- Specify max. number of handler tasks
- Contract between hardware & software

9

ТЛП

## Evaluation: NAS-Benchmarks & Scalability





**Benchmark Results**

- Communication intensive kernels highly benefit from SHARQ

- Only EP (embarassingly parallel) does not profit

- SHARQ has good scaling behaviour

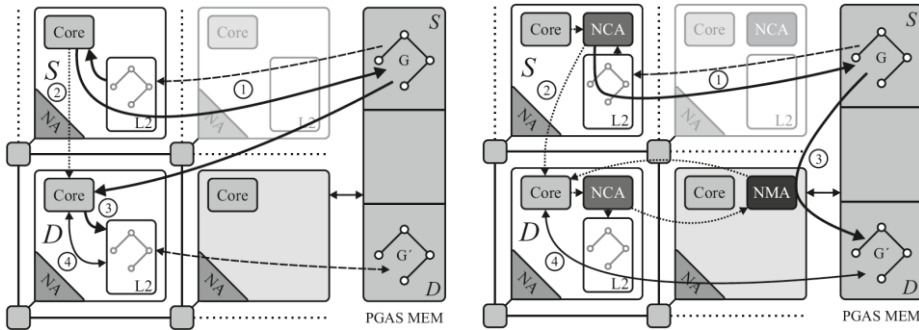- Especially multi-tile systems benefit

10

## Near-Memory Accelerated Graph Copy

**Pegasus**
- Avoids (de)serialization
- Graph copy algorithm in software

**Near-Memory Graph Copy**
- Reduces NoC traffic
- Graph copy algorithm in hardware



11

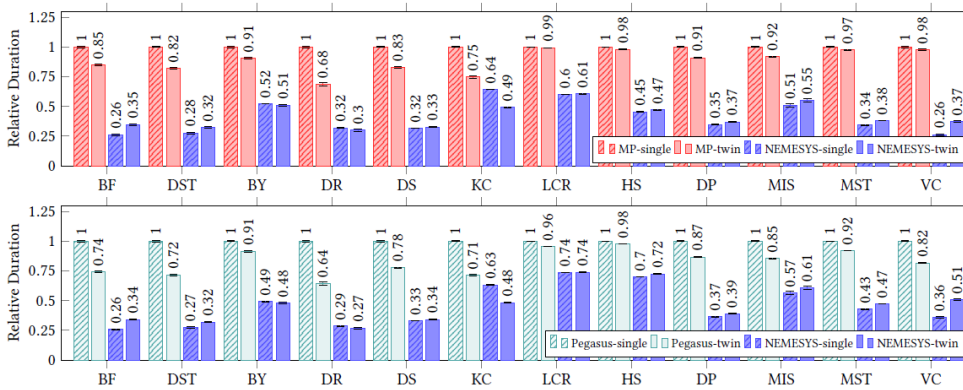## Near-Memory Graph Copy Results

**IMSuite Graph Algorithm Kernels**



Figure 9: Runtime measurements of the IMSuite benchmarks in the 4x4 configuration. Top: NEMESYS vs. Message-passin
(MP) normalized to MP-single. Bottom: NEMESYS vs. Pegasus normalized to Pegasus-single.
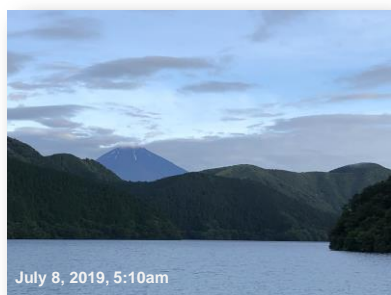
12

TIM

## Summary

Near Memory Acceleration in Distributed-Shared-Memory architectures is an effective means to tackle the locality wall
- NMA primarily means bringing the processing closer to the data
  - … with all its implications / dependencies on data / task placement
  - … or bring / keep the data closer to the processor
    - with Region-Based Cache Coherence

This presentation: Near Memory Acceleration applied to Multicore OS or runtime middleware support functions

13

TIM



July 8, 2019, 5:10am

# Thanks for your attention!

14

TUM

# References

[1] Peter Kogge. 2017. Memory Intensive Computing, the 3rdWall, and the Need for Innovation in Architecture. https://memsys.io/wp-content/uploads/2017/12/The_Wall.pdf

[2] Akshay Srivatsa, Sven Rheindt, Thomas Wild, Andreas Herkersdorf. "Region Based Cache Coherence for Tiled MPSoCs". 2017 30th IEEE International System-on-Chip Conference (SOCC), 2017

[3] Akshay Srivatsa, Sven Rheindt, Dirk Gabriel, Thomas Wild, Andreas Herkersdorf. CoD: Coherence-on-Demand – Runtime Adaptable Working Set Coherence for DSM-based Manycore Architectures. (SAMOS 2019)

[4] S. F. Yitbarek, T. Yang et al, "Exploring specialized near-memory processing for data intensive operations", in DATE, pp. 1449-1452, 2016.

[5] F. Schuiki, M. Schaffner et al, "A Scalable Near-Memory Architecture for Training Deep Neural Networks on Large In-Memory Datasets", 2018.

[6] Sven Rheindt, Sebastian Maier, Florian Schmaus, Thomas Wild, Wolfgang Schröder-Preikschat, Andreas Herkersdorf. SHARQ: Software-Defined Hardware-Managed Queues for Tile-Based Manycore Architectures. (SAMOS 2019)

[7] Manuel Mohr and Carsten Tradowsky. 2017. Pegasus: Efficient Data Transfers for PGAS Languages on non-cache-coherent many-cores. DATE.