

Hauptseminar

Rechnerarchitektur und Programmierung

Adapteva Parallella:

Crowd-funded Low-budget Open-source HPC

Tobias Frust (tobias.frust@mailbox.tu-dresden.de)

Tutor: Ronny Brendel (ronny.brendel@tu-dresden.de)

17th July, 2015

- The Adapteva Parallella Platform
 - Company History of Adapteva
 - Motivation
 - The Parallella Board
 - The Epiphany Coprocessor
- Implementation of the 1D-FFT with Filtering
 - Motivation
 - Implementation Details
 - Performance Results
- Comparison with GPUs
- Assessing the Parallella's Potential for HPC

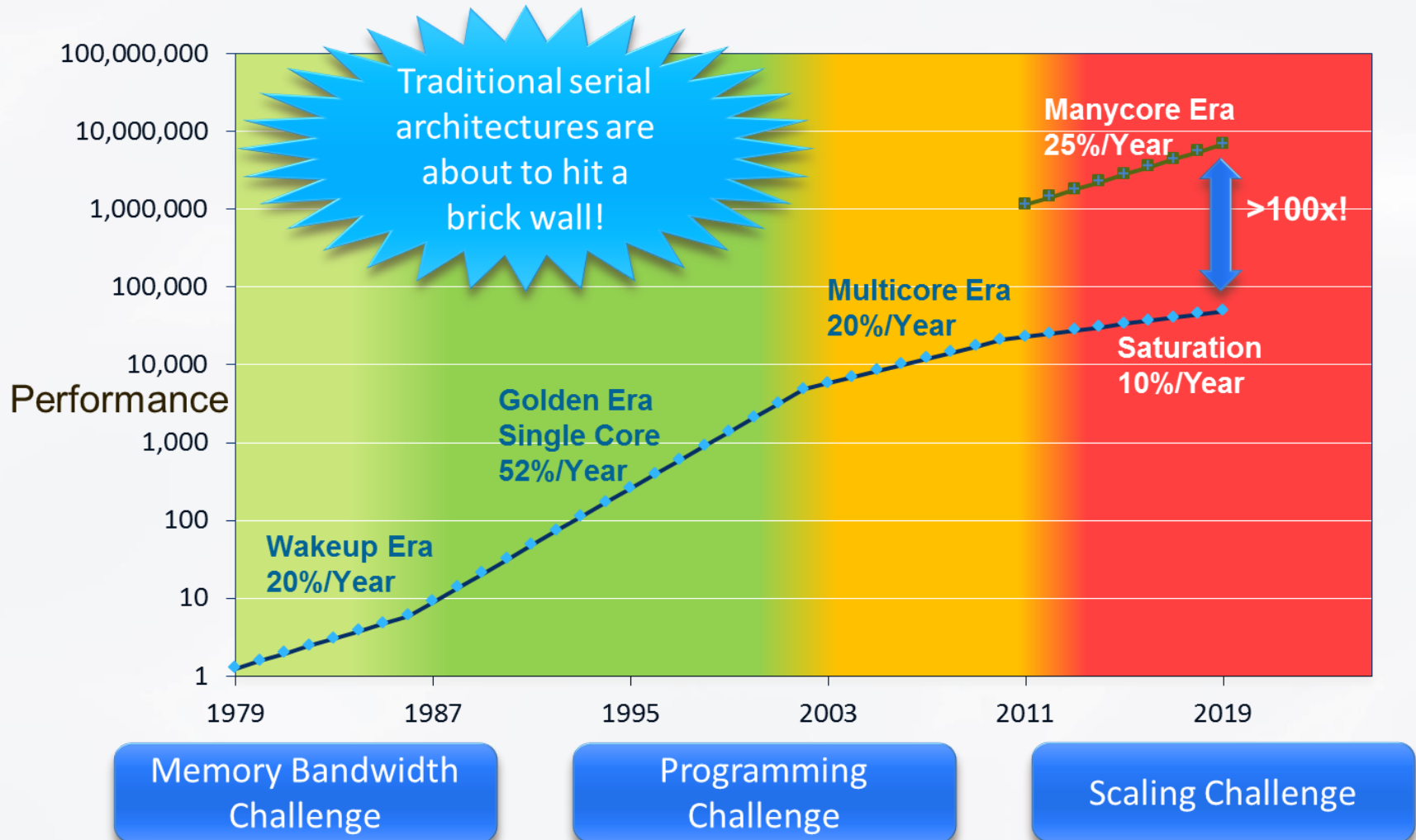
Company History of Adapteva

- Feb 2008: Founded by Andreas Oloffson
 - Goal: 10 times advancement in floating point processing energy efficiency
- Jun 2009: Tapeout of Epiphany-I prototype (65nm)
 - Secured 1.5M US\$ in Series-A funding from Bittware
- May 2011: Sampled Epiphany-III (65nm) product
- Aug 2012: Demonstration of 50 GFLOPS/Watt efficiency at 28nm
- Oct 2012: Launch of Parallella kickstarter project
- Jul 2013: first Parallella boards shipped to Kickstarter backers
- Apr 2014: Completed shipping of all Parallella boards to Kickstarter backers

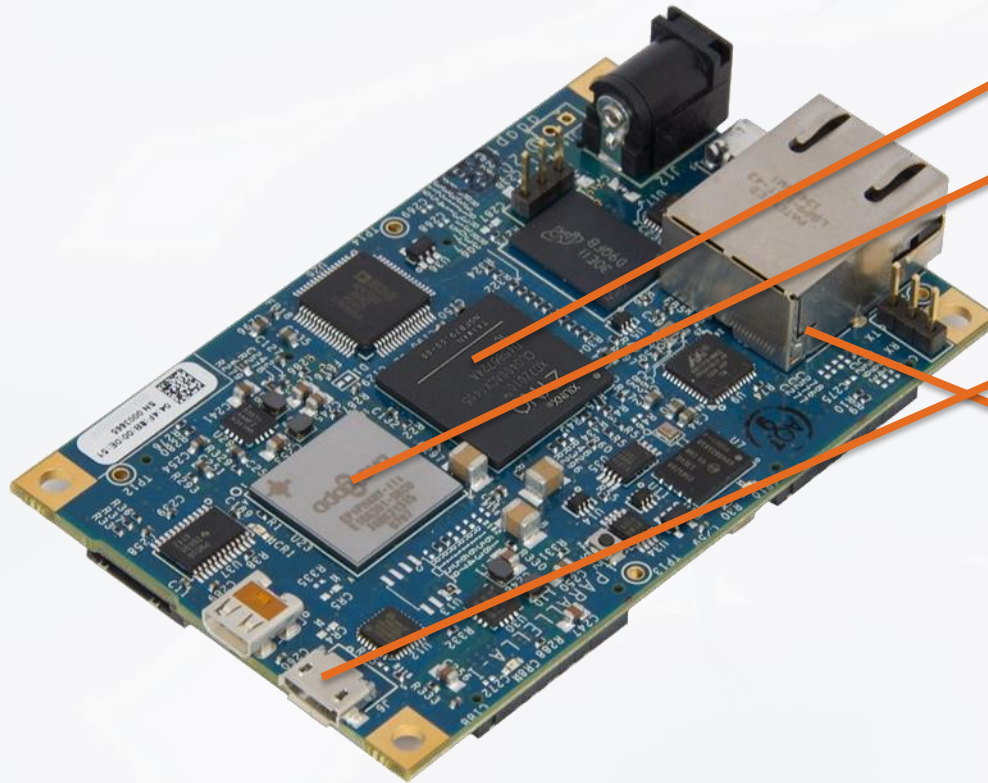
Parallella Kickstarter Program

- Slogan: *“The parallella project will make parallel computing accessible to everyone”*
 - Kickstarter project gained 898,921\$ from 4,965 backers
 - Attributes:
 - **Open Access:** no NDAs or special access needed
 - **Open Source:** platform based on free open source development tools and libraries
 - **Affordable:** Parallella high performance computer at costs below 100\$
- Close the knowledge gap in parallel programming
- Democratize access to parallel computing

Motivation for Epiphany Multicore Architecture



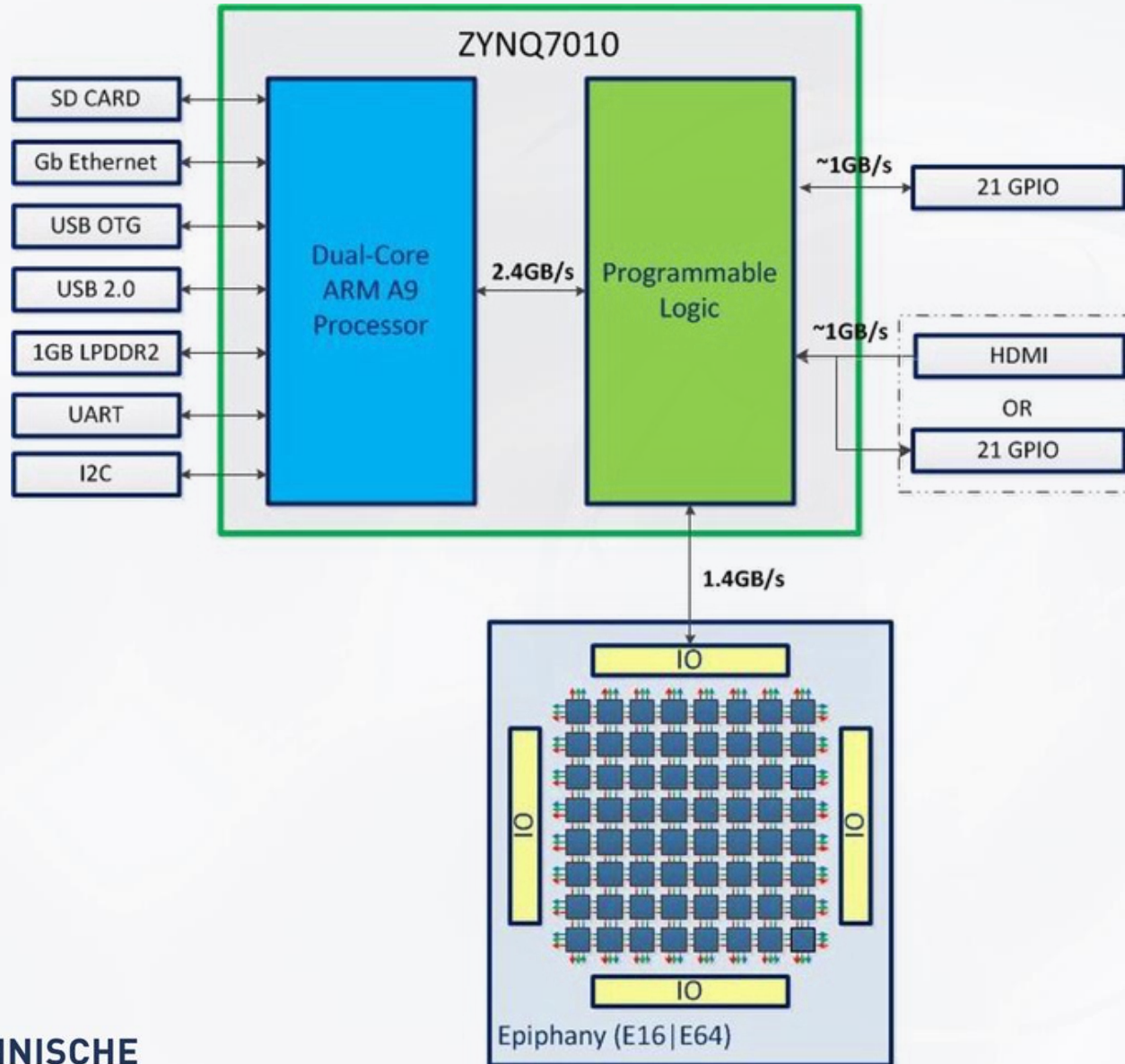
Parallella Board



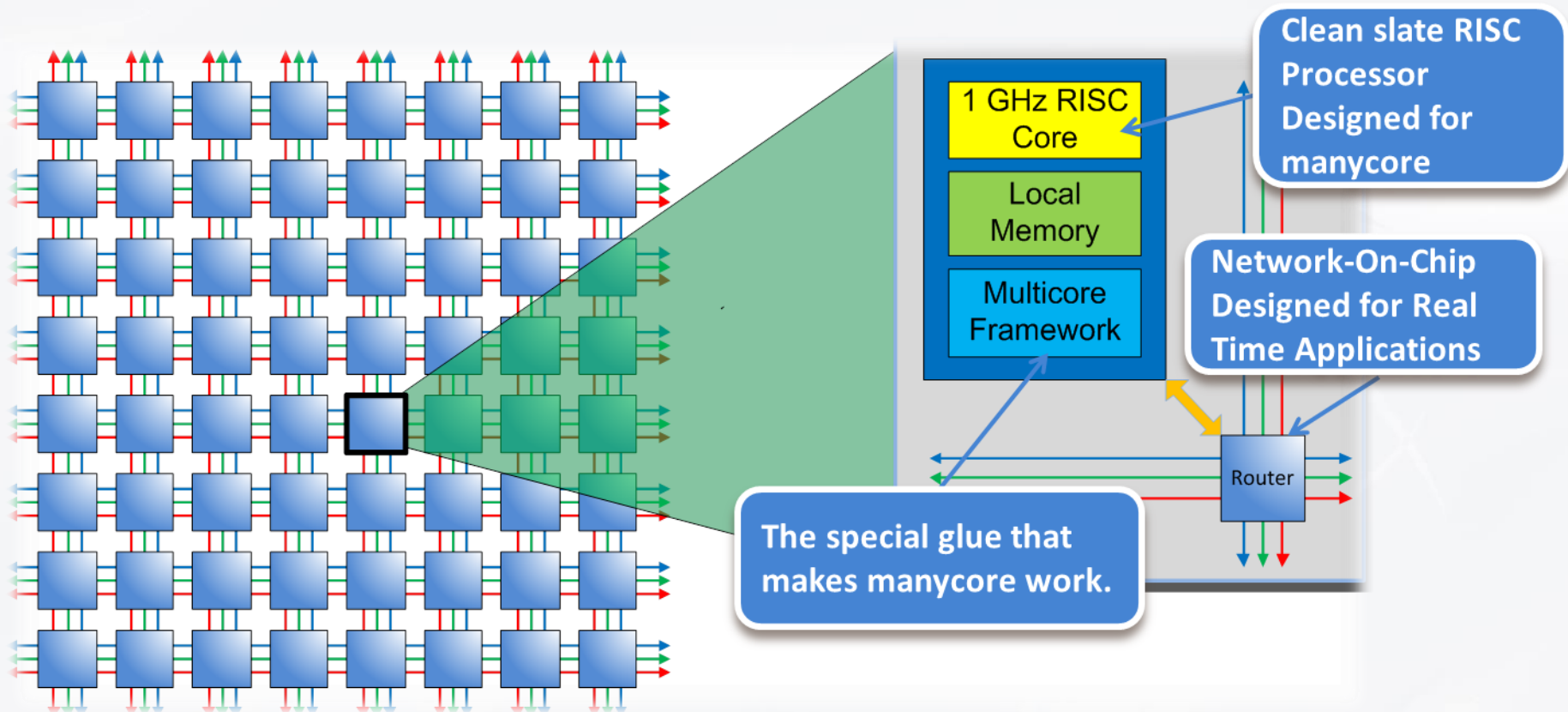
Tech specs:

- Dual-core ARM A9
- 16-core Epiphany Coprocessor
- 1 GB RAM
- MicroSD Card
- USB 2.0
- Gigabit Ethernet
- Linux
- Peak Performance: 26 GFLOP/s (single precision)
- 2 Watt Energy Consumption

System Overview of Parallella Board



High Level Overview of the Epiphany Architecture



Coprocessor to
ARM/Intel CPU

25mW per core

Ease To Use

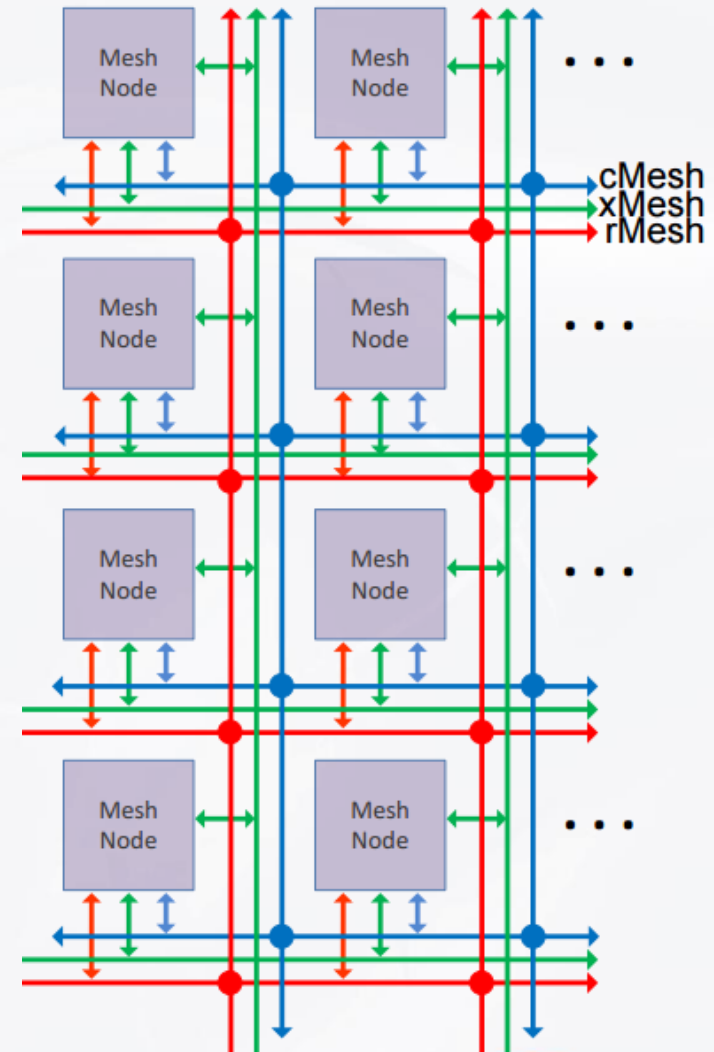


TECHNISCHE
UNIVERSITÄT
DRESDEN

eMesh Network on Chip

Attributes:

- Only nearest neighbor connections
- Connections to North, East, South and West
- 1.5 clock cycles latency per hop (write)
- Three separate orthogonal networks:
 - cMesh: used for on-chip write transactions
 - xMesh: used for off-chip write transactions
 - rMesh: used for all read requests

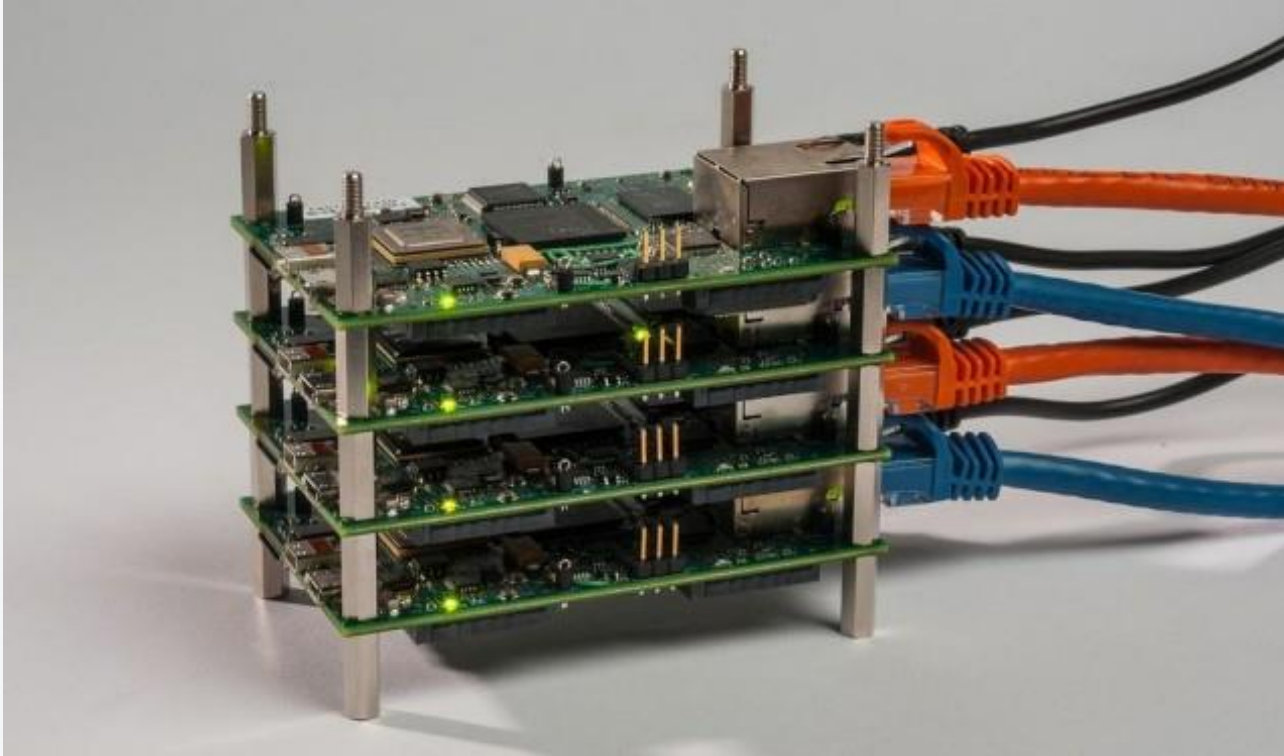


Memory Scheme of the Epiphany Architecture

Attributes:

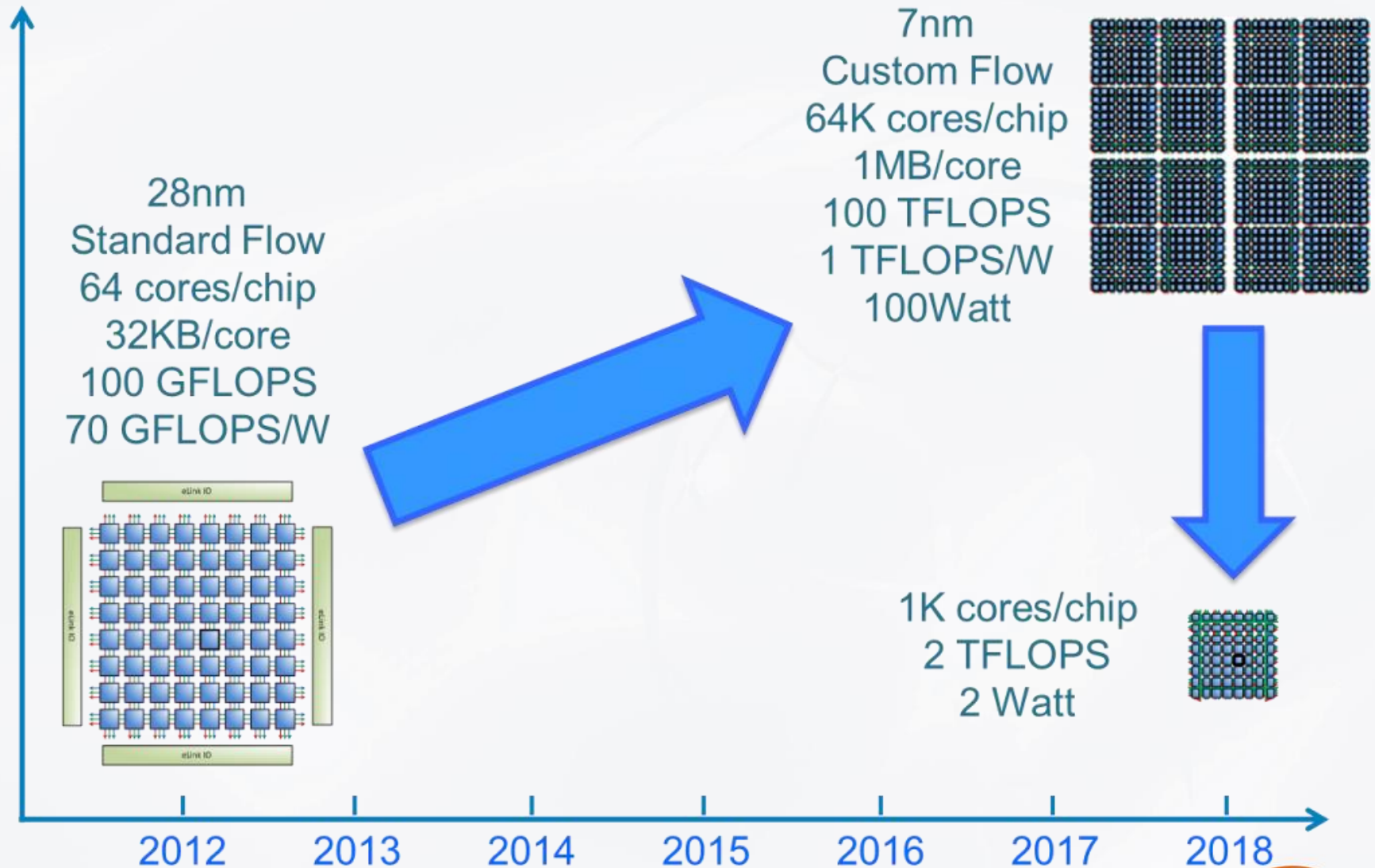
- Flat 32 bit address space split into 4096 1-MiB chunks
- Each core is assigned his own 1-MiB chunk
 - But has transparent access to memory of any other core
- Optimal performance only if data is placed in local memory banks
- Magic address translation – each core can access all memory locations

Example: Parallella as a Cluster



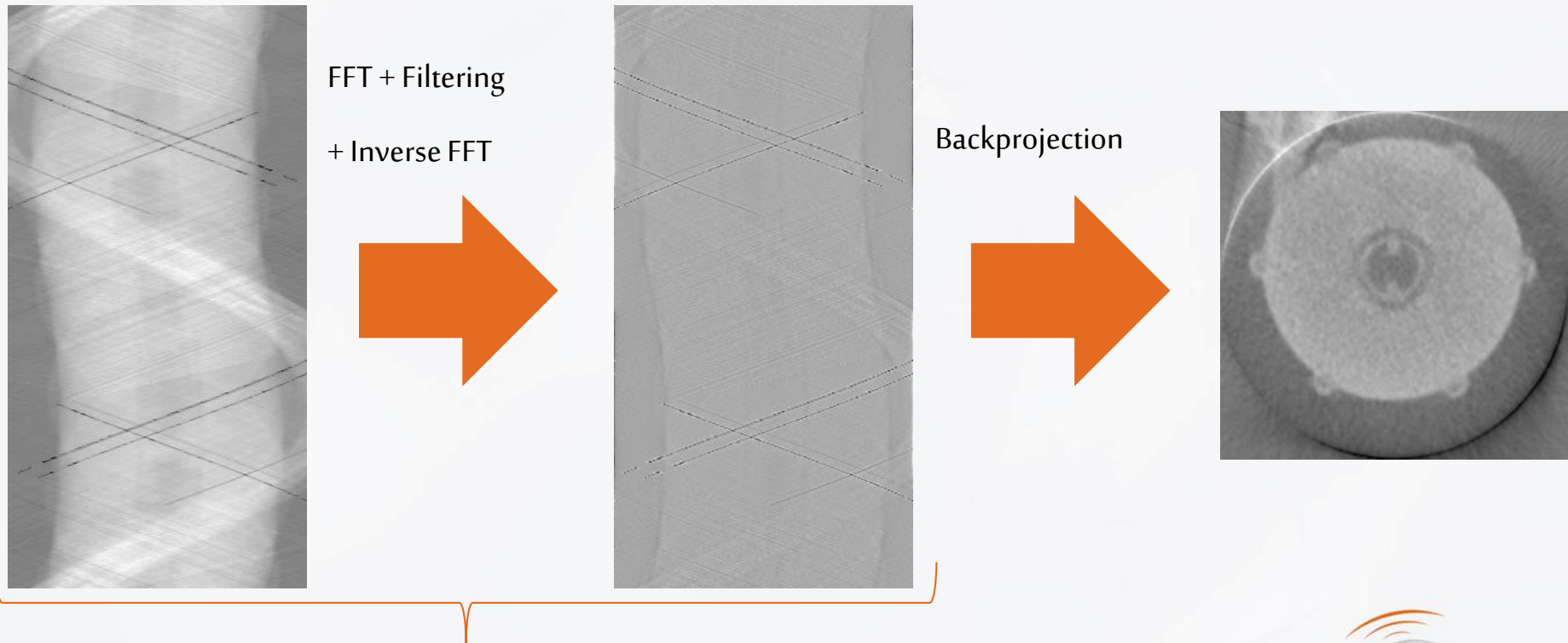
- Boards can be combined to clusters
 - Gigabit-Ethernet interconnect
- With thousands of boards put together you can build a supercomputer at low cost

Future Ideas for Epiphany Architecture



Motivation: Implementation of the 1D-FFT on Epiphany

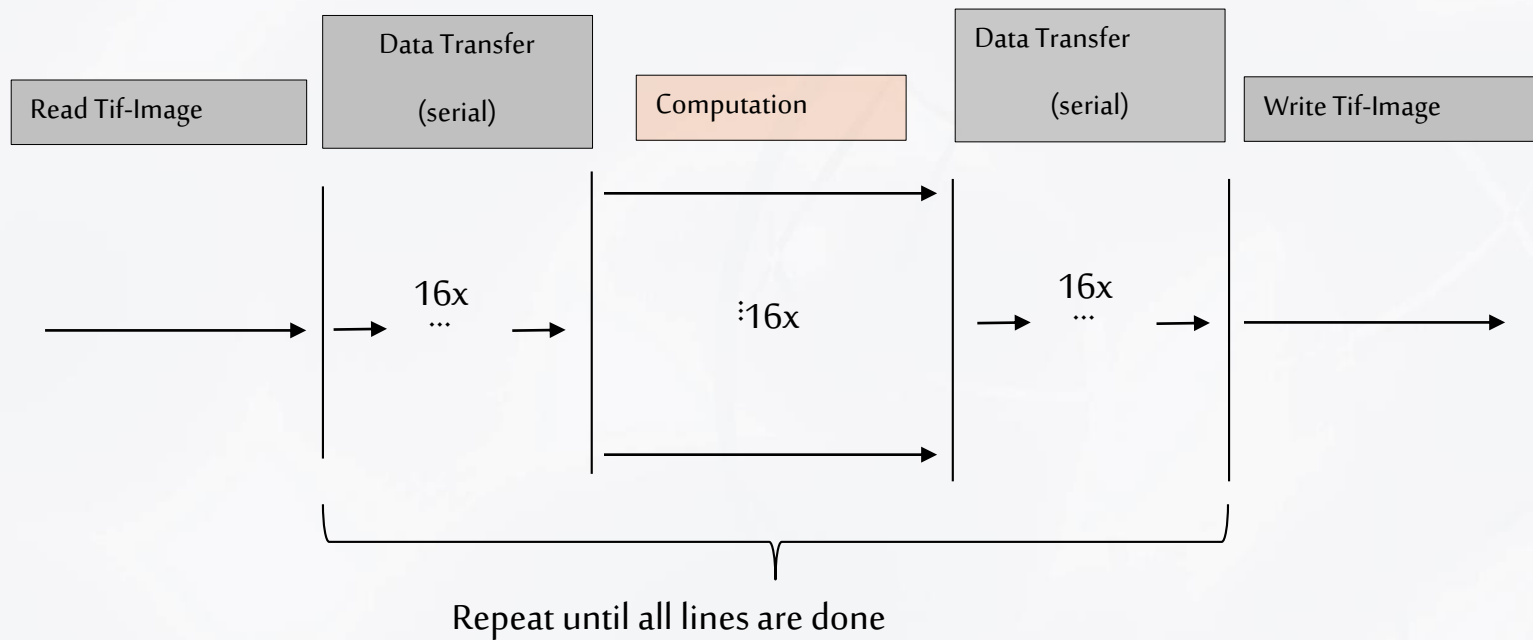
- Fast Fourier Transform is basic operation for many applications
- Algorithm to compute the discrete fourier transform
- Concrete example: Filtered backprojection as CT-Reconstruction algorithm



Implementation details

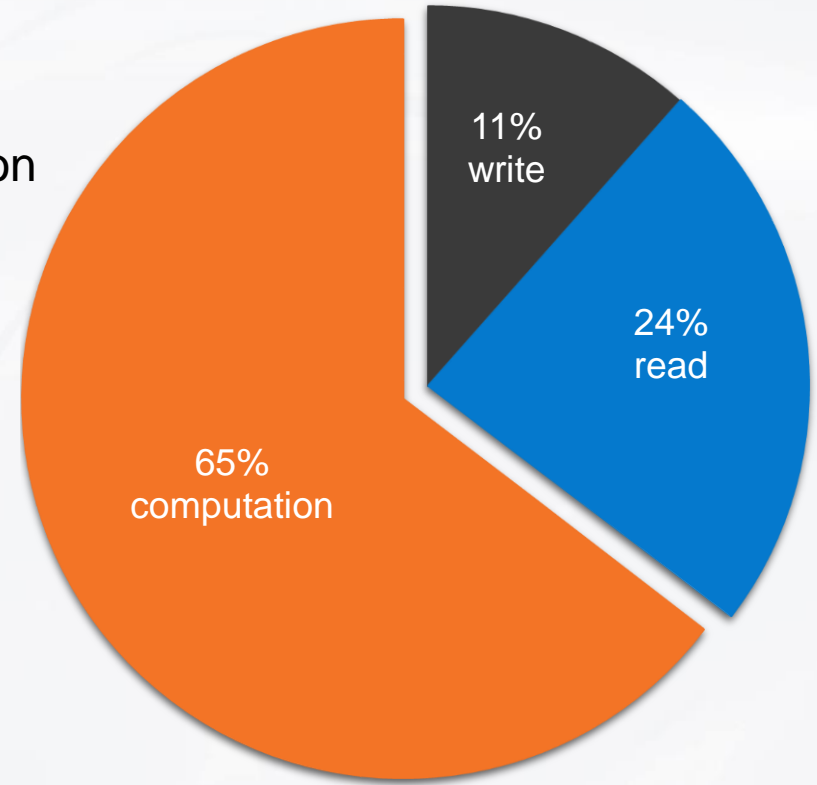
Attributes:

- Sinogram is divided into lines
 - 16 cores calculate 64 lines in parallel to transfer as much data as possible



Data Transfer vs Calculation Time (1 core)

- Data transfer is the main bottleneck
- To increase performance, more calculation per memory transfer is necessary
- Due to Amdahls law speedup cannot exceed 2.5 with 16 cores
- Maximum speedup of 2.9 possible



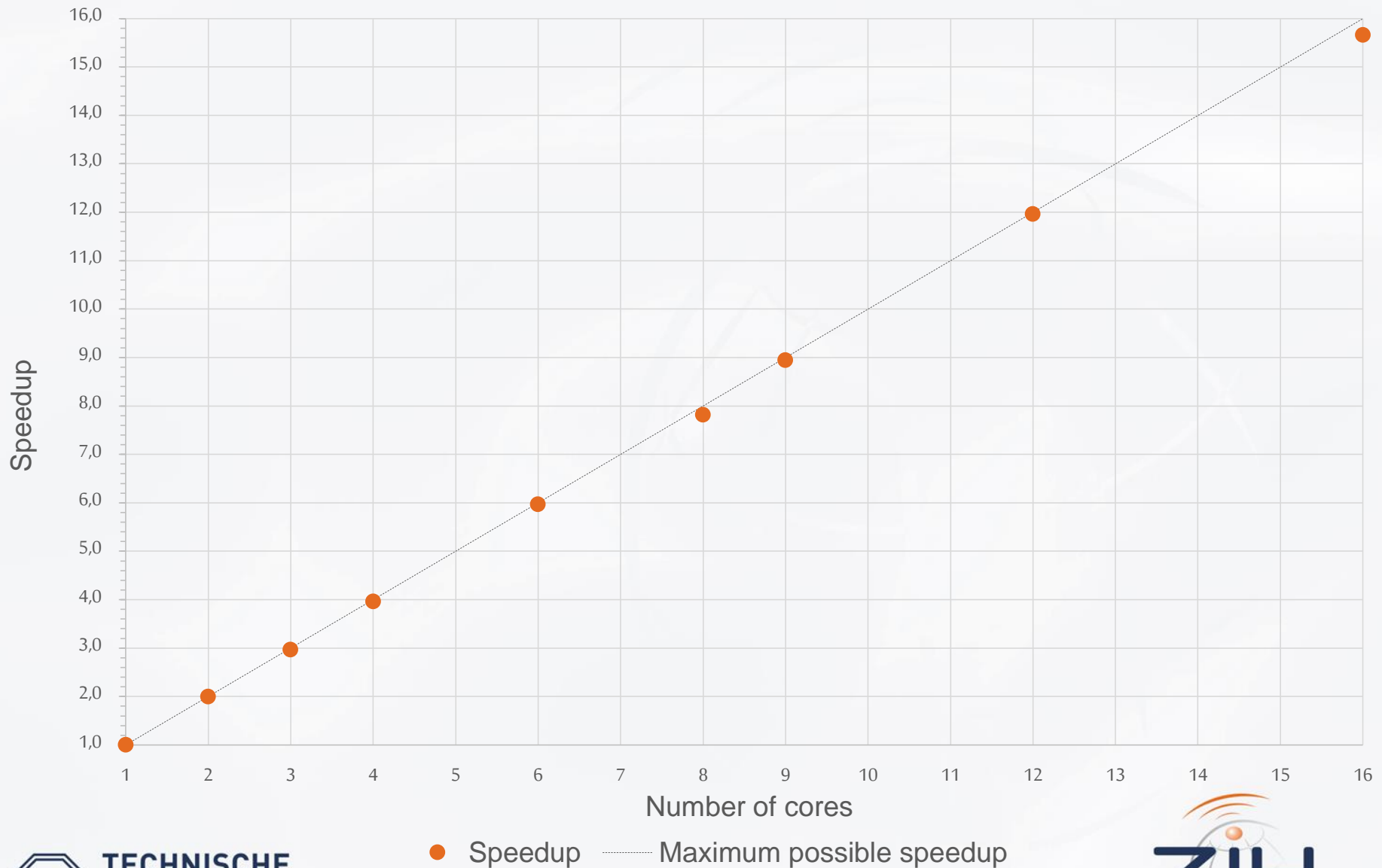
Performance Results

- Results measured for 512 FFTs of size 256

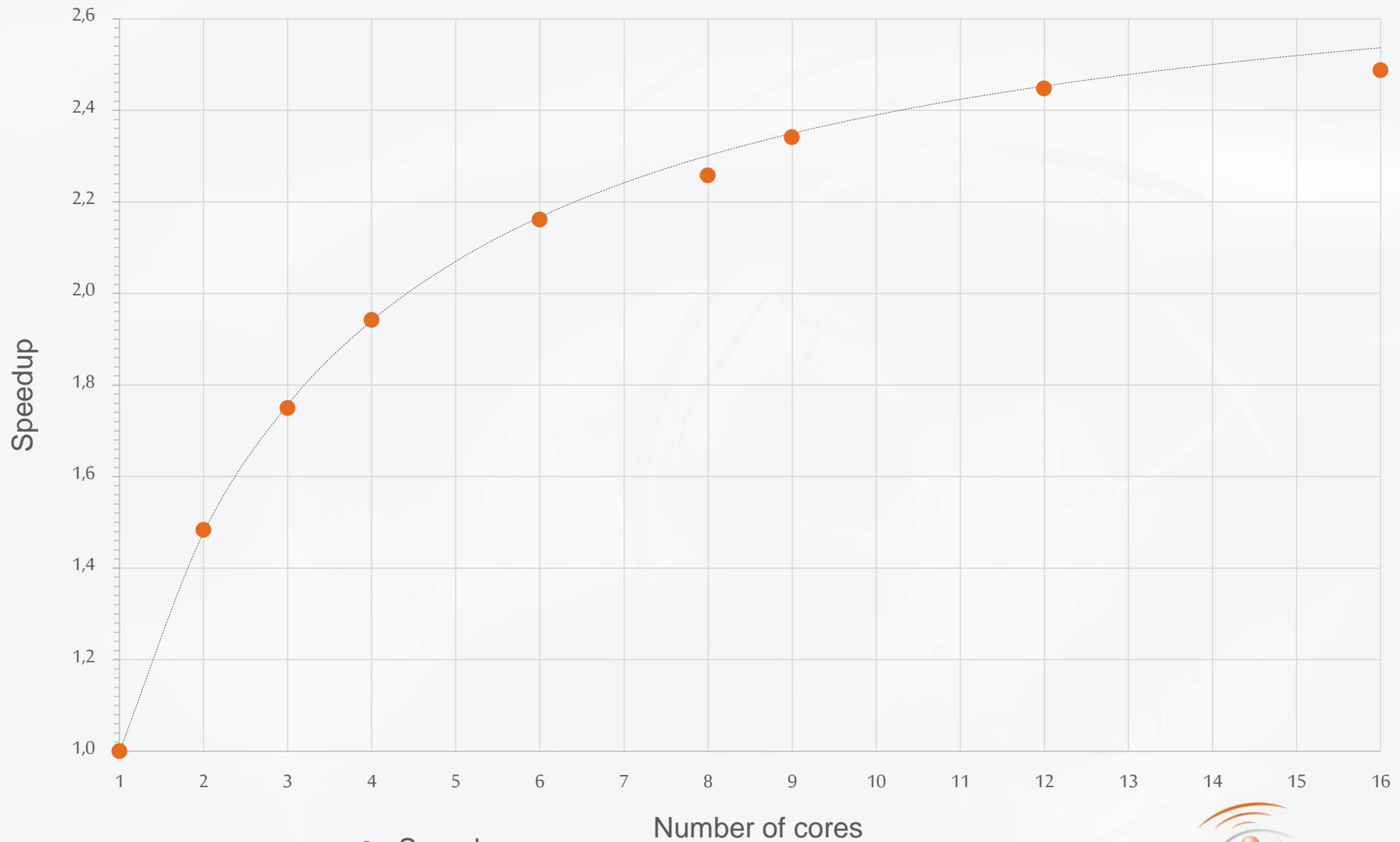
	1 core	16 cores
Write Time	70 ms	72 ms
Computation Time	394 ms	25 ms
Read Time	146 ms	148 ms
Number of Floating Point Operations	$13,9 \cdot 10^6$	$13,9 \cdot 10^6$
MFLOP/s in Computation Part	35 MFLOP/s	556 MFLOP/s
MFLOP/s in whole Program	23 MFLOP/s	57 MFLOP/s
Transfer Rate: Write	14,6 MByte/s	14,2 MByte/s
Transfer Rate: Read	7,01 MByte/s	6,9 MByte/s

- Small memory transfer rate to local memory banks
- Read = read + write → half performance compared to write

Speedup of Parallelisable Calculation Part



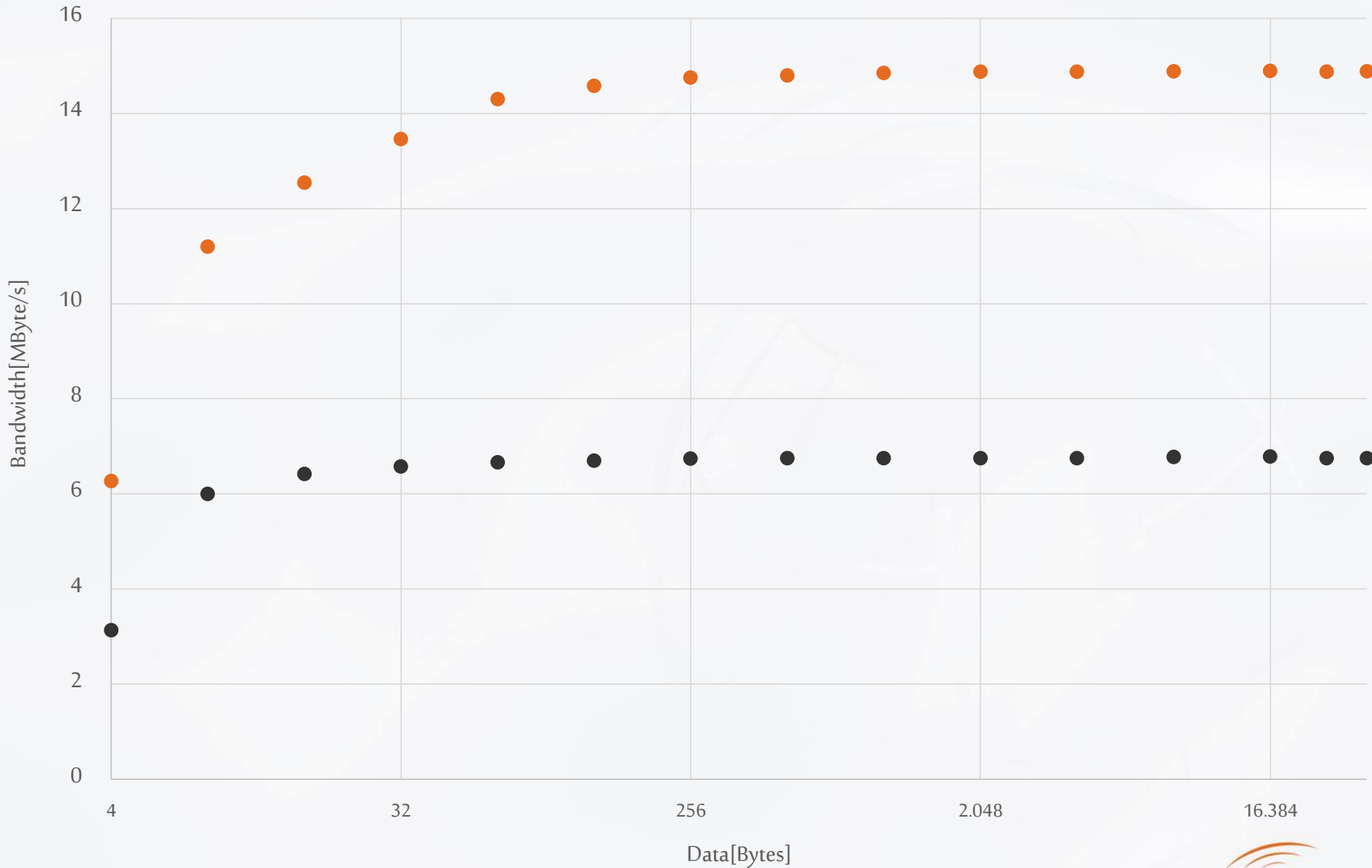
Overall Speedup



● Speedup

..... Maximum possible speedup according to Amdahl's Law

Excursion: Memory Bandwidth



Bandwidth Test of Adapteva

ARM Host	--> eCore(0,0) write speed	=	14.13 MB/s
ARM Host	--> eCore(0,0) read speed	=	6.55 MB/s
ARM Host	--> ERAM write speed	=	75.14 MB/s
ARM Host	<-- ERAM read speed	=	105.46 MB/s
ARM Host	<--> DRAM: Copy speed	=	194.72 MB/s
eCore (0,0)	--> eCore(1,0) write speed (DMA)	=	1280.04 MB/s
eCore (0,0)	<-- eCore(1,0) read speed (DMA)	=	390.01 MB/s
eCore (0,0)	--> ERAM write speed (DMA)	=	240.24 MB/s
eCore (0,0)	<-- ERAM read speed (DMA)	=	87.61 MB/s

Comparison with Current Architecture - GPUs

- Comparison with *Intel Core i7 2600* combined with *Nvidia Geforce 750Ti* in terms of energy efficiency
- Identical implementation with use of *cufft*-Library from Nvidia

Peak Performance Single Precision	1,472 TFLOP/s
Number of Streaming Multiprocessors (SMs)	5
Number of cuda cores	640
Memory Bandwidth	86,4 GByte/s
Architecture	Maxwell
Memory bus interface	128 Bit
Manufacturing Process	TSMC 28nm
Thermal Design Power (TDP)	60 Watt
Launch Date	02/18/14
Die size	148 mm ²

Comparison of Parallella with GeForce 750 Ti

- To stay fair comparison in terms of energy efficiency and chip area
- Energy consumption/Picture [mJ]:

Number of Pictures	Parallella	GPU
1	50	10200
100	50	114
1000	50	18,6

- Initialization phase of cuFFT is very time consuming
- Due to lower memory size, and smaller memory bandwidth of Parallella, GPU can outperform Parallella in terms of energy efficiency in this case

Comparison in Terms of Chip Area

- Example configuration of Epiphany with 1024 cores:

Number of cores	1024
Manufacturing Process	28 nm
Max. Operating Frequency	700 MHz
Peak Performance (Single Precision)	1,408 GFLOP/s
Die size	131 mm^2
Thermal Design Power	20 Watt

- Comparable to NVIDIA Geforce 750 TI in terms of die size
 - Better energy efficiency compared to 16 core Epiphany
 - Approximately 6 times less energy consumption/picture than Epiphany-16
 - better energy efficiency in this case compared to GPU

Assessing the Parallella's Potential for HPC

● Pros:

- Affordable
- Nice for gathering experience in parallel programming
- Direct access to local memory (no managed cache)
- No cache coherency latency
- Energy efficient

● Cons:

- Memory transfer bottleneck
- Small local memory size
- Expandable documentation and small user base



**TECHNISCHE
UNIVERSITÄT
DRESDEN**



Zentrum für Informationsdienste
und Hochleistungsrechnen

Comparison of Parallella with GeForce 750 Ti

- To stay fair comparison in terms of energy efficiency and chip area
- Energy consumption/Picture [mJ]:

Number of Pictures	Parallella	GPU
1	510	10200
100	510	114
1000	510	18,6

- Initialization phase of cuFFT is very time consuming
- Due to lower memory size, and smaller memory bandwidth of Parallella, GPU can outperform Parallella in terms of energy efficiency in this case