# High-performance computing hardware for high data rates

## Agenda

Parallel Computing:
    Possibilities & Challenges
Handling Data I/O at High Rates
Accelerating Synchrotron Tomography
Scaling to Cluster

## Authors

Suren A. Chilingaryan, KIT
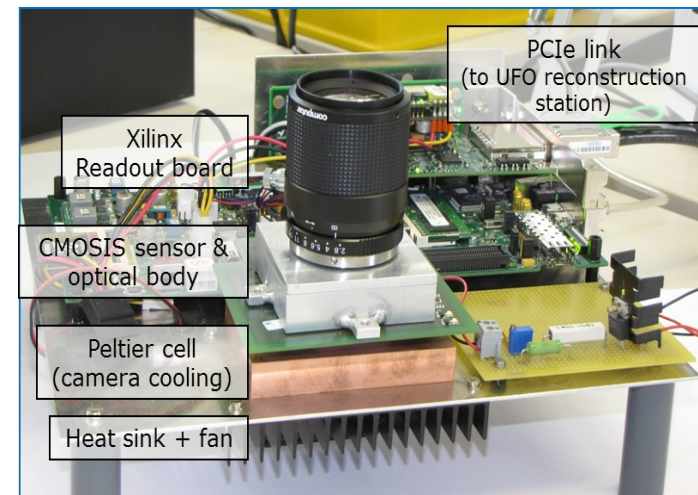Michele Caselle, KIT
Thomas van de Kamp, KIT
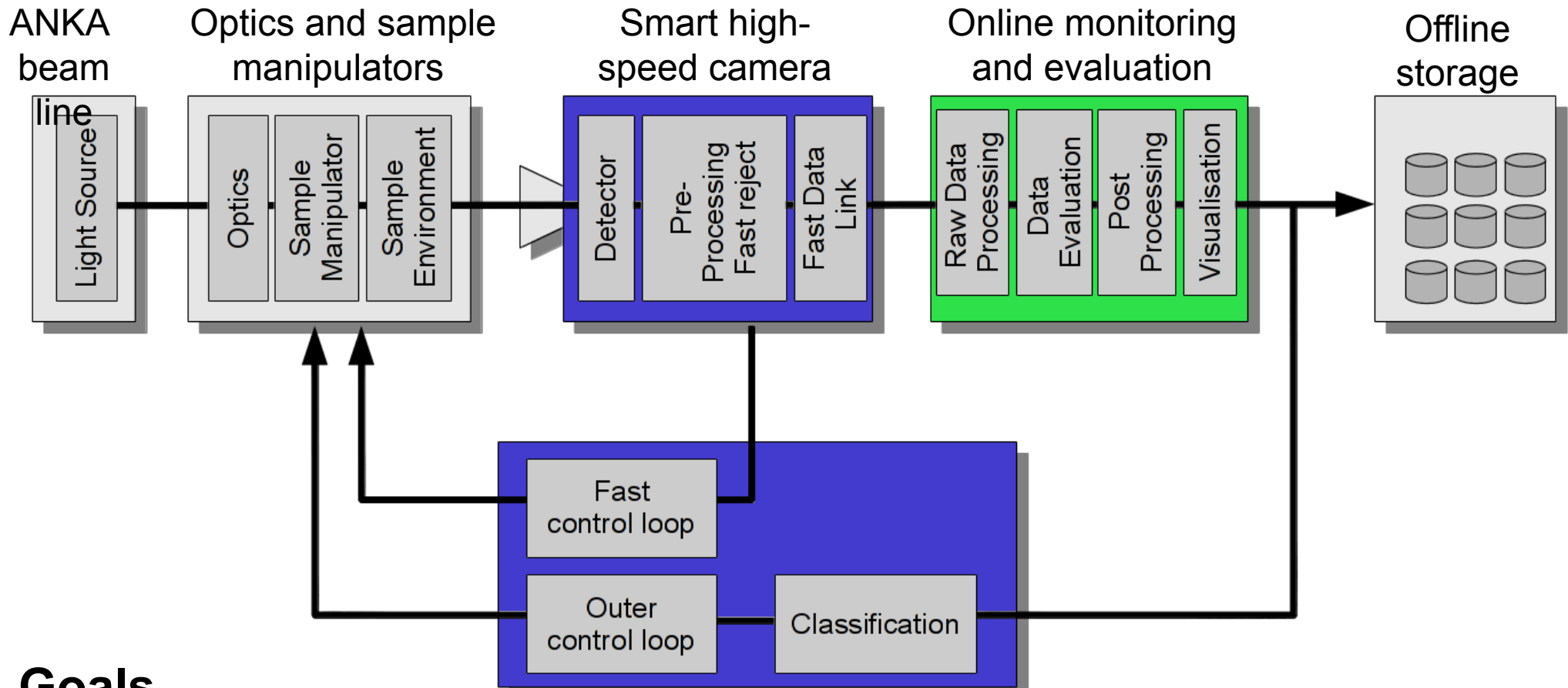Andreas Kopmann, KIT
Alessandro Mirone, ESRF
Uros Stevanovic, KIT
Tomy dos Santos Rolo, KIT
Matthias Vogelgesang, KIT

PCIe link (to UFO reconstruction station)
Xilinx Readout board
CMOSIS sensor & optical body
Peltier cell (camera cooling)
Heat sink + fan

# UFO

## Ultra Fast X-ray Imaging of Scientific Processes with On-Line Assessment and Data-Driven Process Control

KIT — Karlsruhe Institute of Technology



## Goals

- High speed tomography
- Increase sample throughput
- Tomography of temporal processes
- Allow interactive quality assessment

- Enable data driven control
  - Auto-tunning optical system
  - Tracking dynamic processes
  - Finding area of interest

S. Chilingaryan et. all

# Reconstruction Problem

## PCO.edge



Resolution: 2560 x 2160
Dynamic Range: 16 bit
Frame Rate: 100 fps

## Tomographic Reconstruction

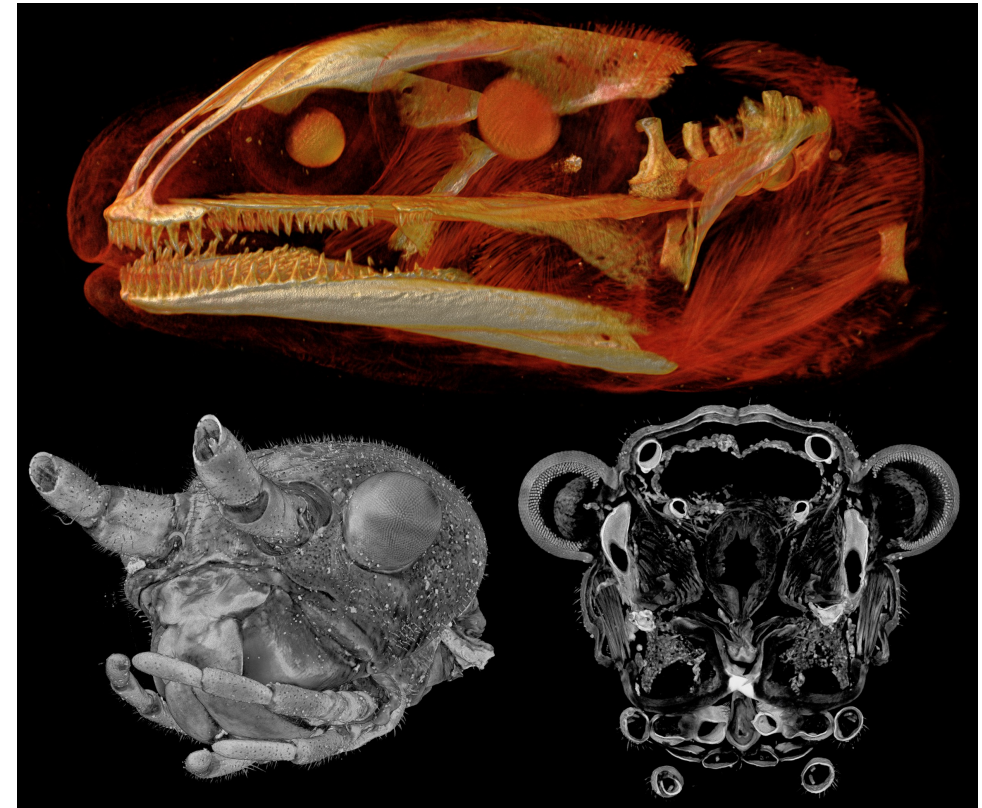3D image: $2000^3$
Projections: 2000
Acquisition time: 20 seconds

FBP Complexity: 144 Tflops
Xeon Performance: ~ 100 Gflops
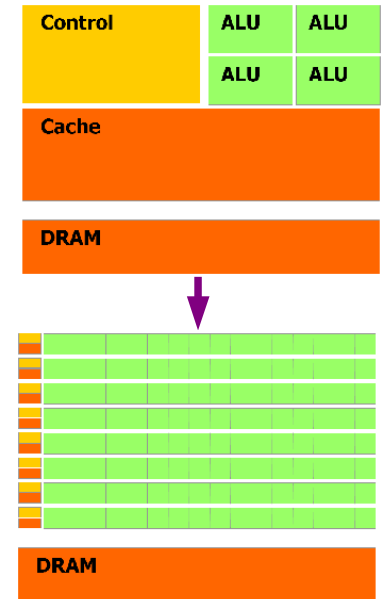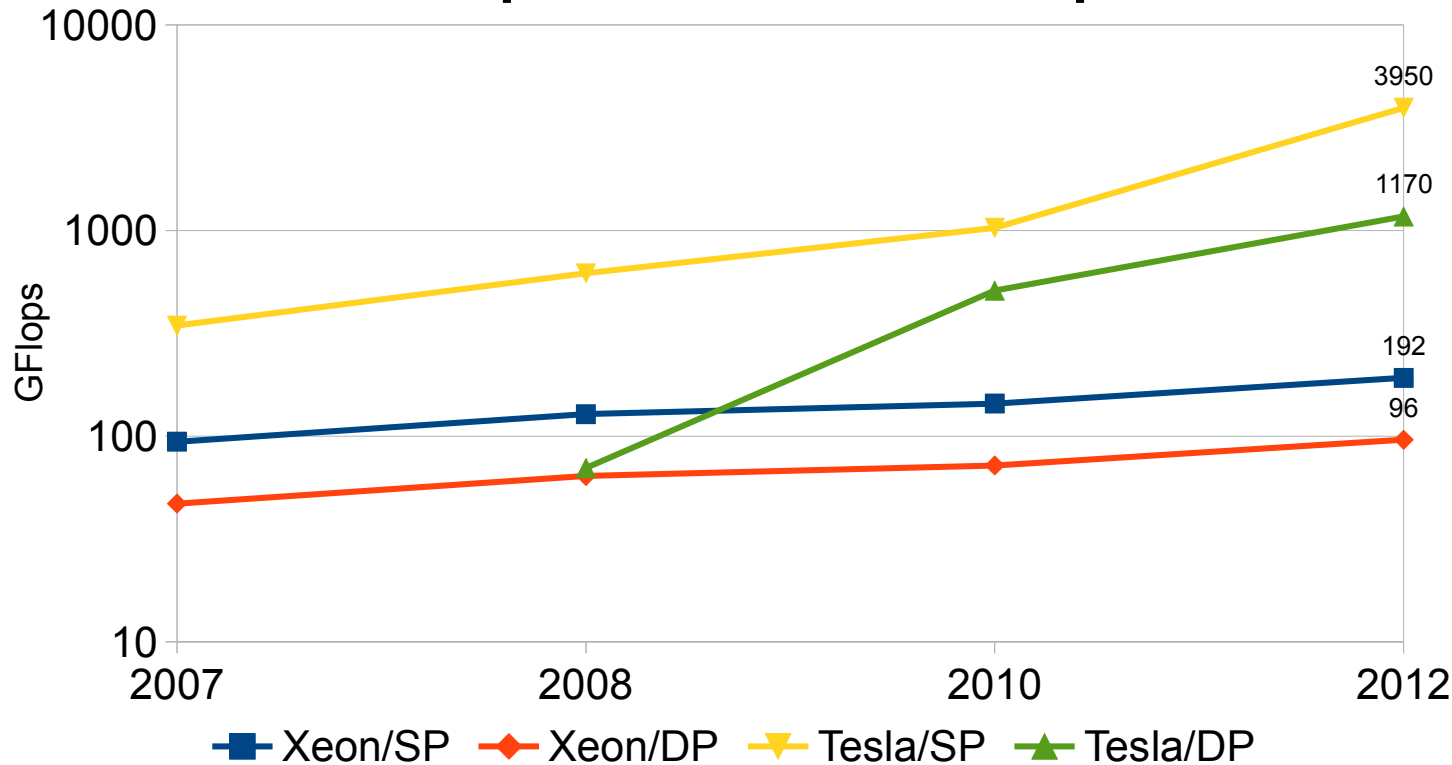Minimum time: ~ 15 minute on DP
Actually: ~ 1 hour

## 20 seconds acquisition
## 1 hour reconstruction



*Heads of a newt larva showing bone formation and muscle insertions (top) and a stick insect (bottom), acquisition time 2s.*

S. Chilingaryan et. all

# Parallel Architectures

## Rise of GPU performance as compared to Xeons



*Parallel Architecture*

| | E7-8870 | Xeon/Phi | Tesla K20 | GeForce Titan | AMD HD7970 | Power7+ |
|---|---|---|---|---|---|---|
| **SP** | 192 | 2020 | 3950 | **4500** | 3790 | 265 |
| **DP** | 96 | 1010 | 1170 | **1300** | 950 | 132 |
| **Mem** | 34.11 | **320** | 250 | 288 | 264 | 68 |
| **Max dev.** | 8 | ? | 8 | 8 | >=4 | 32 |
| **Price** | $4,800.00 | $2,800.00 | $3,200.00 | $1,000.00 | **$400.00** | $$$$$$$ |

S. Chilingaryan et. all

# Efficiency

## Matrix Multiplication



Core i7-980X — 94%
HD7970 — 63%
GTX580 — 63%
GTX680 — 43%
GTX Titan — 72%

Legend: ■ Peak ■ Measured
x-axis: 0, 1000, 2000, 3000, 4000, 5000 GFlops

## 1D Fast Fourier Transform



Core i7-980X — 52%
HD7970 — 3%
GTX 580 — 18%
GTX 680 — 9%
GTX Titan — 10%

x-axis: 0, 50, 100, 150, 200, 250, 300, 350, 400, 450, 500 GFlops

## Memory Bandwidth



y-axis: GB/s — 0, 50, 100, 150, 200, 250, 300

Core i7-980X — 98%
HD7970 — 69%
GTX580 — 86%
GTX680 — 78%

Legend: ■ Peak ■ Measured

GTX Titan performance is taken from Anandtech

S. Chilingaryan et. all

# GPU-programing considerations

- **Special programming tools and techniques are required**
  - Multiple different and ever-changing architectures
  - Branching, non-fp operations are very expensive
  - Optimized mathematical libraries are some times missing
- **Limited amount of memory and expensive data transfers**
  - x16 PCIe gen2 (8 GB/s), gen3 (16 GB/s)
  - Specially allocated (pinned) memory required for a full performance and to overlap computations and data transfers
- **Reduced caches, low memory to computation ratio, strict access patterns**
  - 177 GB/s per Teraflop for Xeon, 60 - 70 GB/s  per Teraflop for GPUs
  - Varying cache hierarchies on different architectures
  - Special access patterns are required for better performance. For instance, bandwidth of matrix transpose (GTX280 with 142 GB/s memory bandwidth)
    - **2 GB/s –** for naive approach
    - **17 GB/s –** if shared memory is used
    - **80 GB/s –** if care taken for shared memory banks and global memory partitions
- **I/O problem**
  - 100 MB/s sequential write while camera produces ~ 1 GB/s
  - Handling big data sets not fitting in the memory (up to 500 GB)
- **Various problems with growing number of GPUs connected to a system**

S. Chilingaryan et. all

# Parallel Programming Environments

- **CUDA –** The oldest GPU programming technology from NVIDIA
- **OpenCL –** Open standard technology close to CUDA, but working with wide range of hardware platforms including CPUs and GPUs
- **OpenAC –** Declarative technology similar to OpenMP
- MATLAB and other mathematical packages with integrated GPU support. Only some operations are parallelized and necessity to transfer over slow PCIe bus to execute non-parallelized operations kill the performance.
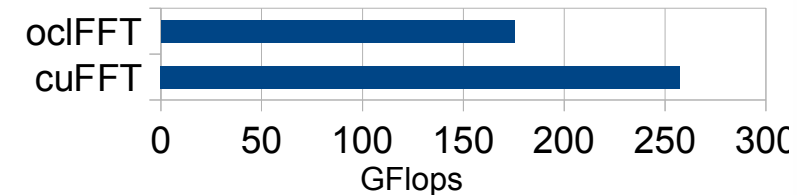
## CUDA

- **Supports latest NVIDIA technologies**
  - GPUDirect – direct transfers between GPU and IB, etc. Integration with MPI frameworks
  - Dynamic parallelism – GPUs are ablee to spawn new jobs
- **NVIDIA provides a set of highly optimized libraries (BLAS, FFT, Lapack, Reduction, etc.)**
- **Only NVIDIA GPUs are supported**

S. Chilingaryan et. all

# Programming Environments

## OpenCL

- Syntax is very similar to CUDA (easy porting)
- Well written code is as fast as CUDA
- Works with CPUs and GPUs from multiple vendors (Intel, AMD, IBM, NVIDIA)
- Still no way to run code simultaneously on NVIDIA GPU and CPU, but possible with AMD cards
- Many libraries existing, but generally slightly slower than CUDA counterparts. Some libraries are only available commercially
- No GPUDirect, significantly limited options to use pinned memory (i.e. slower data transfers)

## OpenACC

- Existing applications may be easily parallelized. Also developing new code is easy compared to OpenCL/CUDA
- No free compilers are existing at the moment. Though there is a similar technology OmpSS developed at Barcelona Supercomuter Center.
- At current level, technology does not support shared memory and some other technologies available with direct programing (i.e. it is slower)

S. Chilingaryan et. all

# GPUDirect and Frame Grabbing

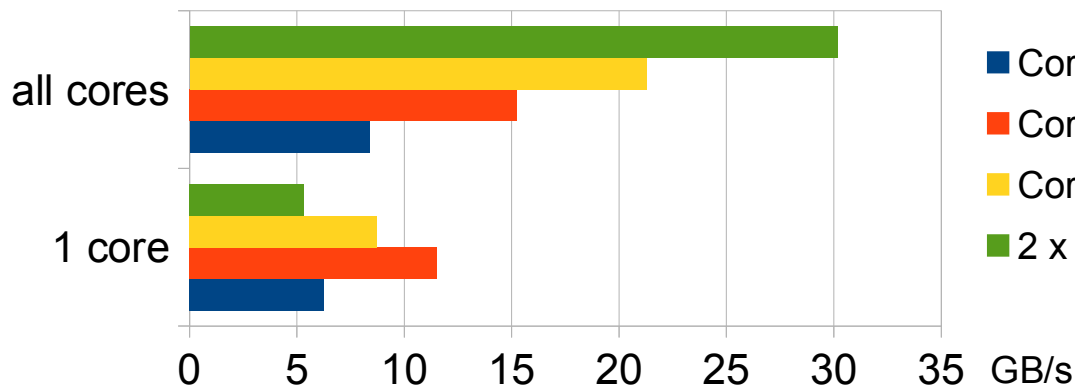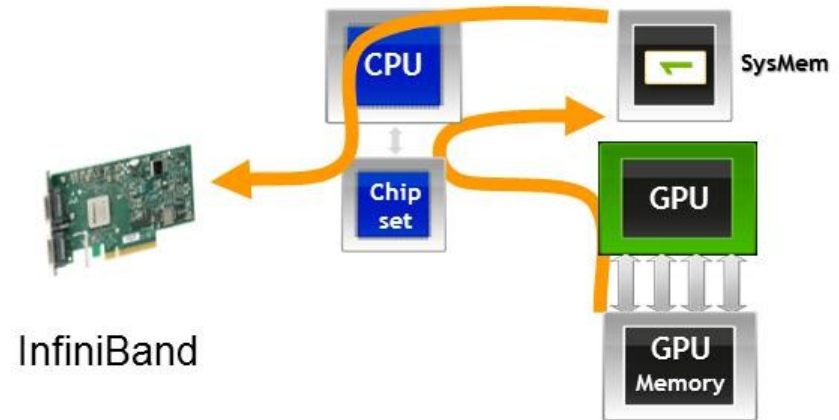## Without GPUDirect

Same data copied three times:
1. GPU writes to pinned sysmem1
2. CPU copies from sysmem1 to sysmem2
3. InfiniBand driver copies from sysmem2



InfiniBand

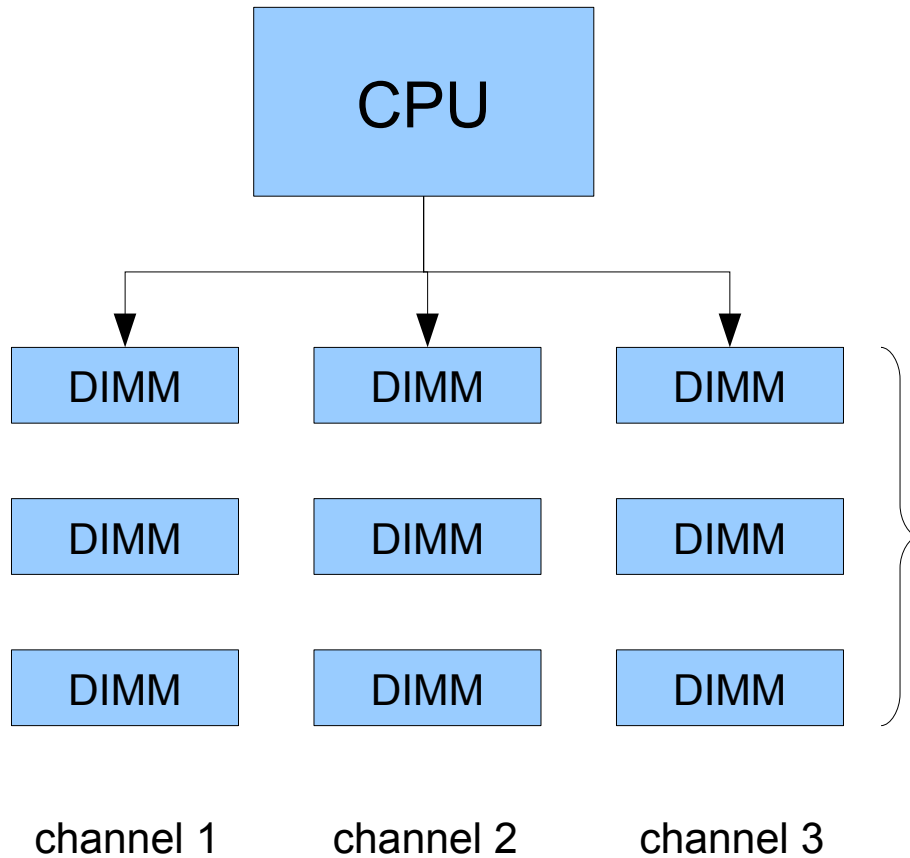## With GPUDirect

Data only copied twice
Sharing pinned system memory makes
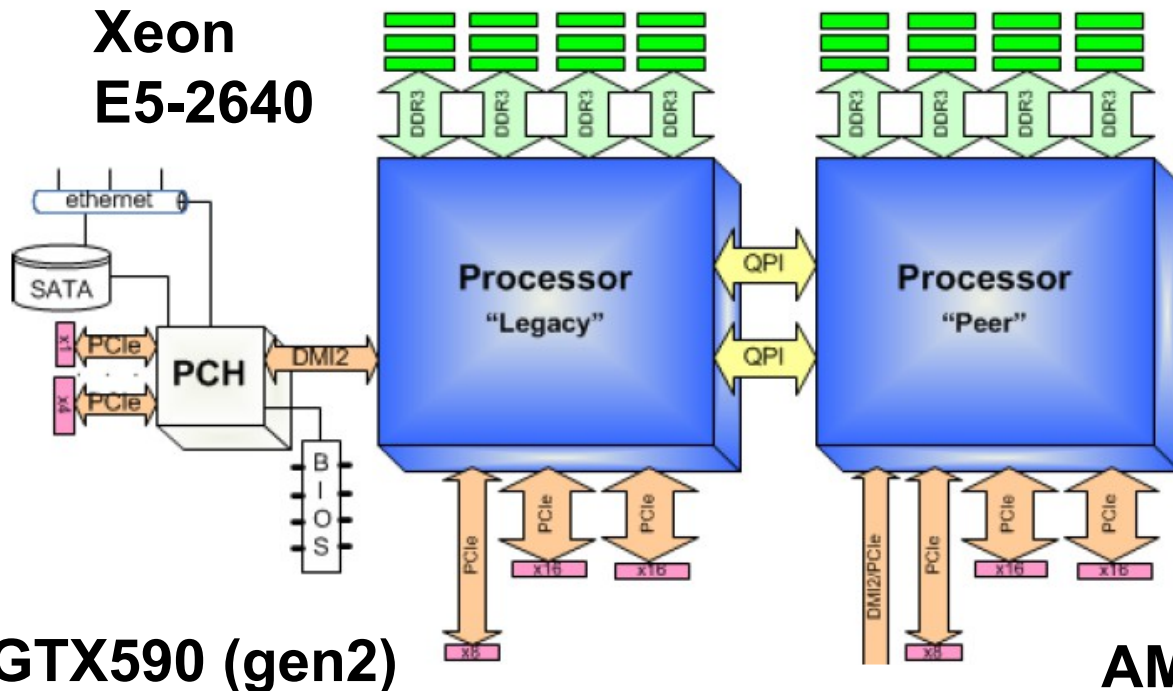sysmem-to-sysmem copy unnecessary



InfiniBand



- Core i7 950
- Core i7-980X
- Core i7 3820
- 2 x E5-2640

And we get ~ 1 GB/s from camera. With 3 memcpy it is already on the border.

S. Chilingaryan et. all

# Memory: Space vs. Speed

CPU

| DIMM | DIMM | DIMM |
| DIMM | DIMM | DIMM |
| DIMM | DIMM | DIMM |

channel 1  channel 2  channel 3

| DPC DIMM per Channel | Xeon X5500 | Xeon E5-2600 |
| --- | --- | --- |
| **1 DIMM** | 10.6 GB/s | 12.8 GB/s |
| **2 DIMMs** | 8.5 GB/s | 12.8 GB/s |
| **3 DIMMs** | 6.4 GB/s | 8.5 GB/s |

S. Chilingaryan et. all

# NUMA Architecture and Data Transfers



**Xeon E5-2640**

**Bandwidth**
**QPI**: 14.4 GB/s
**DDR3**: 10.6 GB/s (PC1333)
**PCIe**: 16 GB/s (gen3 x16)

QPI bus is even not enough to feed both GPU cards
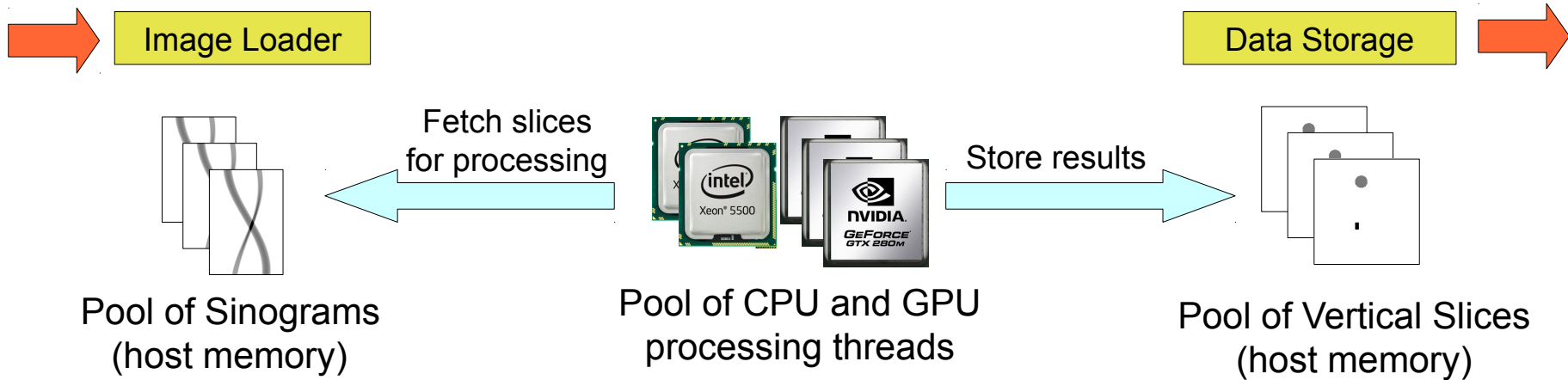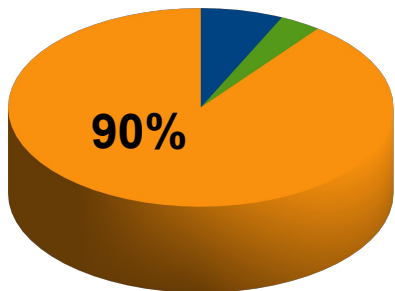
**GTX590 (gen2)**

**AMD HD7970 (gen3)**

Pinned memory support is limited with OpenCL

NVIDIA does not support gen3 mode on X79 boards (workaround exits for Win, but not Linux)

S. Chilingaryan et. all

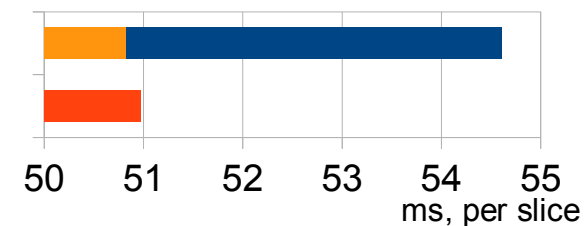# Filtered Back Projection on GPU

KIT
Karlsruhe Institute of Technology

Image Loader

Data Storage

Fetch slices for processing

Store results

Pool of Sinograms
(host memory)

Pool of CPU and GPU
processing threads

Pool of Vertical Slices
(host memory)

PCIe Data Transfer

PCIe Data Transfer

**GPU thread**

**1ˢᵗ Stage**

**2ⁿᵈ Stage**

Double buffering

H

W

Double buffering

Filtering

Texture

## Ratio of operations

**90%**

- Back Projection
- Filtering
- PCIe Transfer

- Transfer
- Compute

## Overlapping Efficiency

50  51  52  53  54  55
ms, per slice

S. Chilingaryan et. all

# Tuning for hardware architectures

**GT200**
Base version
Uses texture
engine

**Fermi** +100%
High computation power, but low speed of texture unit
Reduce load on texture engine: use shared memory to cache the fetched data and, then, perform linear interpolation using computation units.

**Kepler** +75%
Low bandwidth of integer instructions, but high register count
Uses texture engine, but processes 16 projections at once and 16 points per thread to enhance cache hit rate

**VLIW** +530%
Executes 5 independent operations per thread
Computes 16 points per thread in order to provide sufficient flow of independent instructions to VLIW engine

**GCN** +95%
High performance of texture engine and computation nodes
Balance usage of texture engine and computation nodes to get highest performance

S. Chilingaryan et. all

Back Projection: Evolution of GPU architectures

S. Chilingaryan et. all

# Adding more GPU devices

## Initialization Time



*Y-axis: time (seconds), values 5, 7, 9, 11, 13, 15, 17, 19, 21, 23, 25*
*X-axis: Number of GTX590 cores used, values 1 through 8*

## NVIDIA

▶Maximum 8-9 GPU cores (not cards) per system. System will not turn on otherwise

▶Lan Option ROM have to be turned off in the BIOS

▶The PCIe slots, where storage adapters inserted, have to be disabled

▶ASTRA Lab reported to run 13 GPU cores with modified BIOS

▶To run more than 5 GPUs, NVIDIA driver have to be force to use MSI interrupts. Crashes will occur otherwise

## AMD

▶4 GPU cards (single core) working fine, no configuration modifications required

▶Dual-core card are working in a single-core mode only

S. Chilingaryan et. all

# Handling large data sets



Using SSD drives may significantly increase random access performance to the data sets which are not fitting in memory completely. The big arrays of magnetic hard drives will not help unless multiple readers involved.
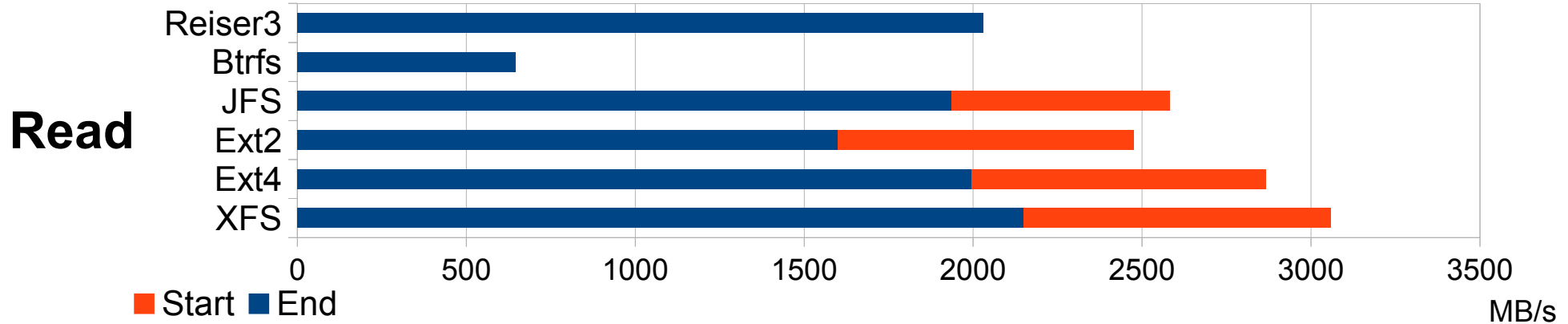
S. Chilingaryan et. all

# Streaming data: file system caches

Data

↓

Buffer Cache

↓

Storage

Default data flow in Linux



**Write** / **Read** — Buffered, Direct, AIO (MB/s)

■ Buffered  ■ Direct  ■ AIO

▶ Buffer cache significantly limits maximal write performance

▶ Kernel AIO may be used to program IO scheduler to issue read requests without delays

Optimizing I/O for maximum streaming performance using a single data source/receiver

# Data Streaming

- ▶ Used file system matter. And it should be adapted to raid configuration (strip, read-ahead)
- ▶ Unless really big number of disks used, the start of partition will be faster than the end)
- ▶ fallocate may significantly improve performance (allocation unit may be increased during FS creation/mount, XFS supports allocation sizes up to 1GB)
- ▶ Ext4 does not support partitions more than 16TB yet
- ▶ Real-time feature of XFS is unstable, data is loss is likely

19

S. Chilingaryan et. all

# Processing Pipeline



## 4 stage pipeline
I/O + Computations

1. Reading data from fast  SSD Raid-0 (random reads are effective)
2. Scheduling and preprocessing using SIMD instructions of x86 CPUs
3. Reconstructing on GPUs
4. Storing to Raid on magnetic disks (sequential writes are effective)

S. Chilingaryan et. all

# Building a server

- **Too much external hardware is required**
  - High speed network
  - Storage system (and SSD cache separately preferably)
  - High speed Frame Grabber for Camera
  - Normally 4-6 high speed PCIe slots per server
  - Space for 1-2 GPUs only
- **System cooling is complicated**
  - both GPUs, HDDs, and SSDs produce a lot of heat
- **Extensibility**
  - There is no space to add more storage / computing power

S. Chilingaryan et. all

# UFO Computing Infrastructure



**Camera**

PCO.edge
PCO.dimax
PCO.4000

CameraLink
850MB/s

External PCIe x16 (8 GB/s)

SFF8088 (2.4 GB/s)

Ethernet
10 Gb/s

**Storage**

LSDF
Large Scale Data Facility

**SuperMicro 7046GT-TRF** (Dual Intel 5520 Chipset)
CPU: 2 x Xeon X5650 ( total 12 cores at 2.66 Ghz)
GPUs: 4 x GTX590 External
Memory: 96 GB / 12 DDR3 slots (192GB max)
Network: Intel 82598EB (10 Gb/s)
Camera Link Frame Grabber (850 MB/s)
Storage: Areca ARC-1880-ix-12 SAS Raid
   16 x Hitachi A7K200 (Raid6)
    8 x Samsung 840 Pro 510 (Raid0)

## External GPU Box

Internal ■  External ■

GT/s

## SSD Raid

Read ■  Write ■

MB/s

## SAS Attached Storage

32 disks ■  16 disks ■

sequential write, MB/s

S. Chilingaryan et. all

# Filtered Back Projection Performance



11 GB data set

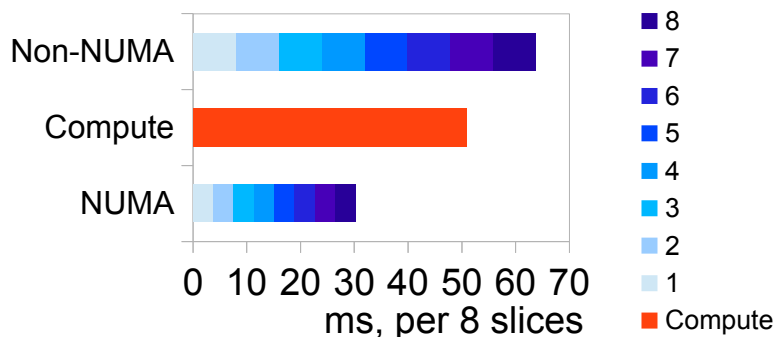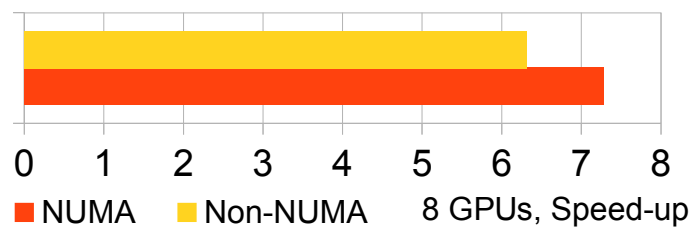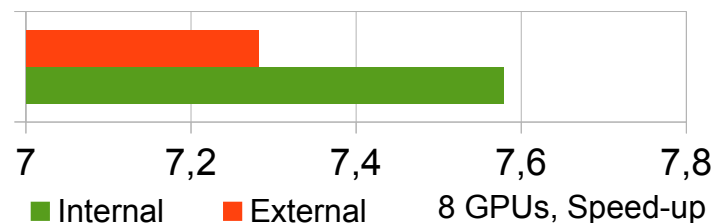S. Chilingaryan et. all
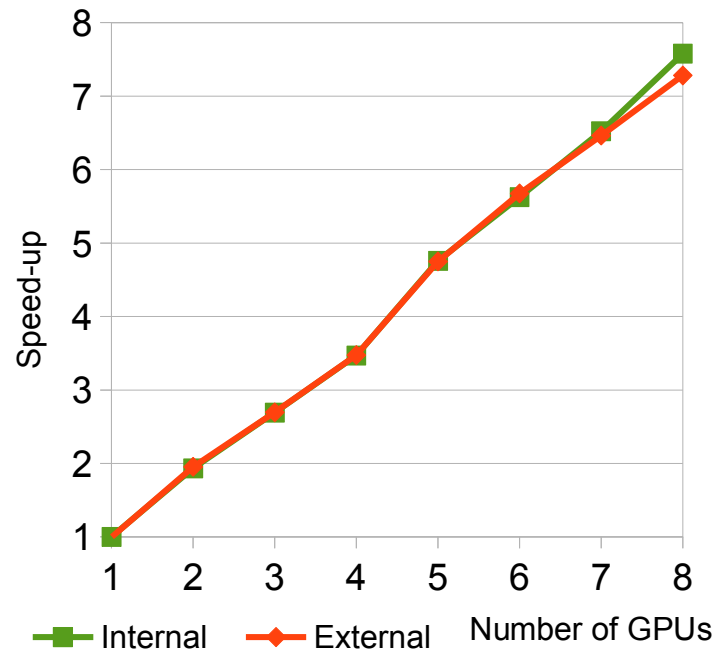
# PCIe Extension Box



1 x PCIe x16 2.0
4 x GTX590
8 GPU cores

External GPU Enclosure
by One Stop Systems

## Scalability



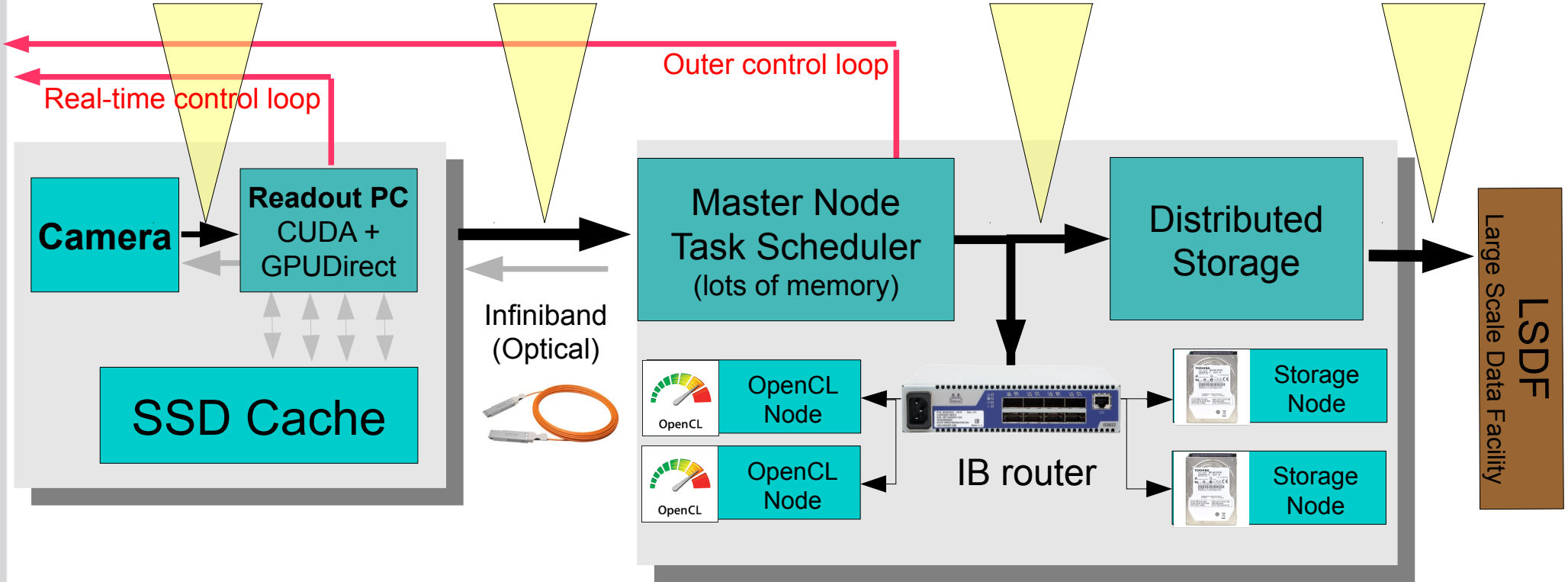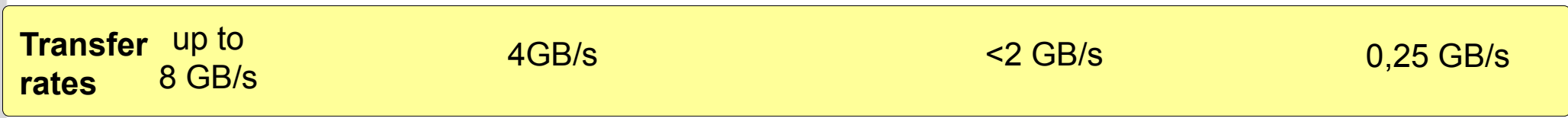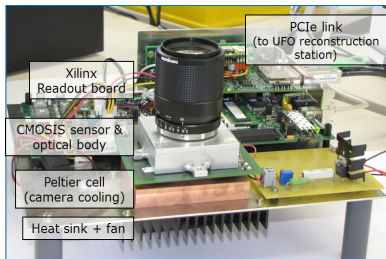**Non-NUMA / NUMA — 1 slice, transfer time, ms**



1 slice, transfer time, ms

### With external box configuration



ms, per 8 slices

Legend: 8, 7, 6, 5, 4, 3, 2, 1, Compute



Internal / External — 8 GPUs, Speed-up



NUMA / Non-NUMA — 8 GPUs, Speed-up

S. Chilingaryan et. all

# Scaling up to Cluster

S. Chilingaryan et. all

# Infiniband: Connection and Protocols

**Mellanox ConnectX 3 VPI**



SDP is obsolete by OpenFabric alliance,
but we have patches for latest kernels.

S. Chilingaryan et. all

# Storage Protocols

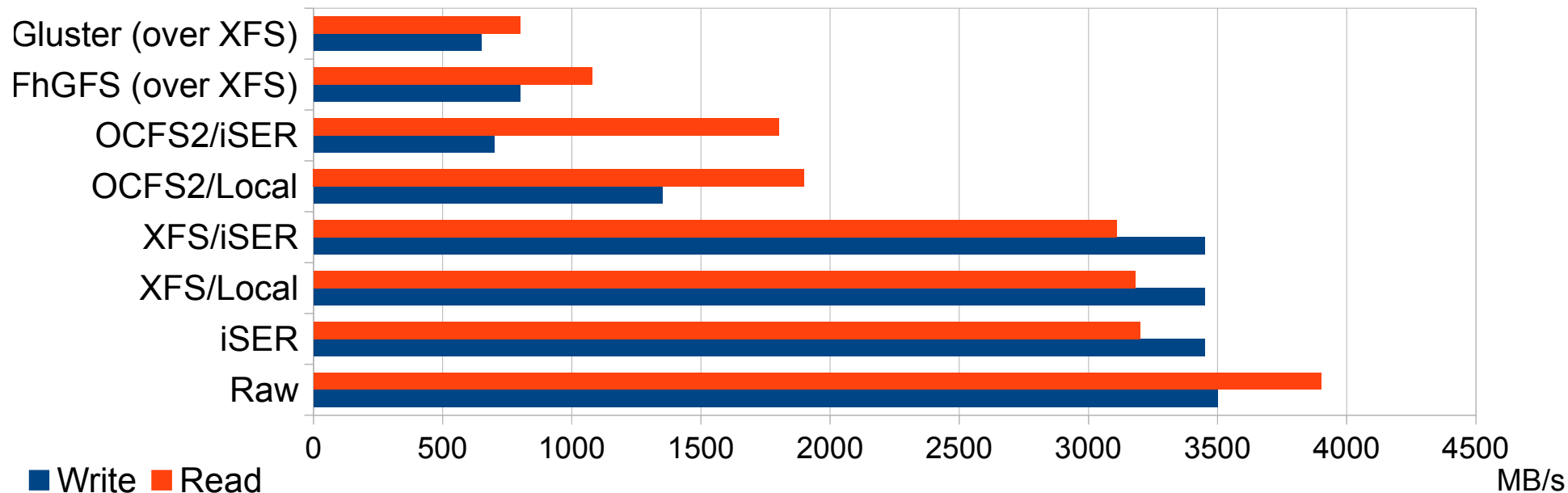Network FS
NFS
Samba
SSHFS
Slow

Cluster FS
Lustre (patched kernel)
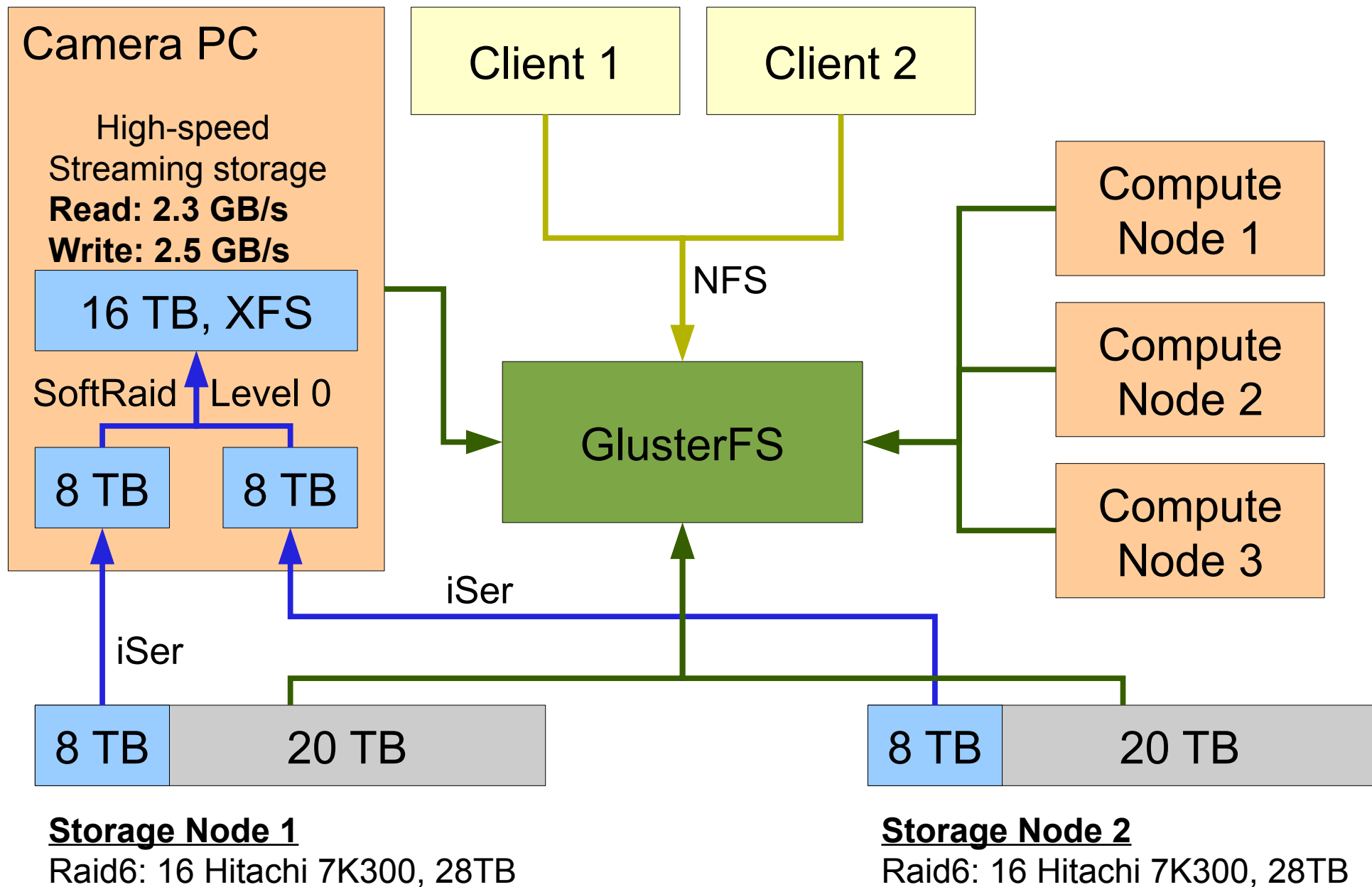Gluster
FhGFS (close-sourced)
Slow if few nodes
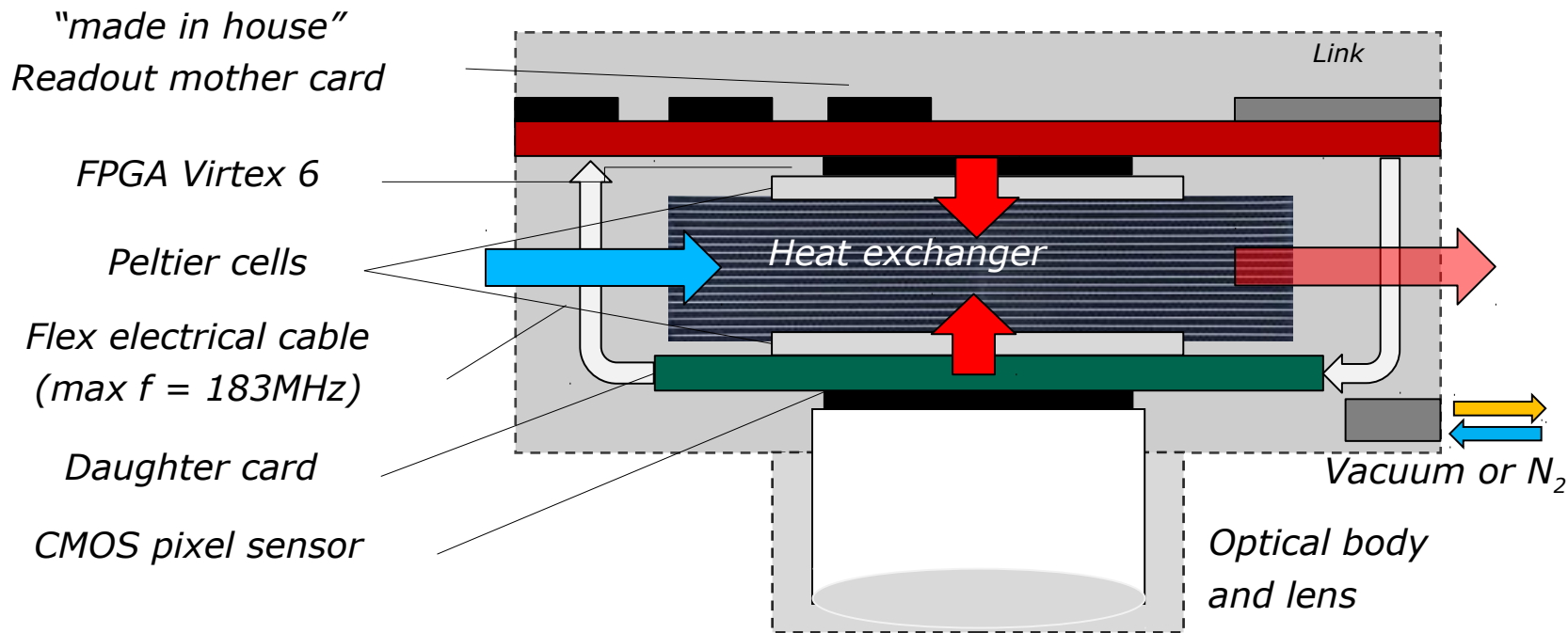
Network Devices
ISCSI (slow)
iSER

OCFS2

S. Chilingaryan et. all

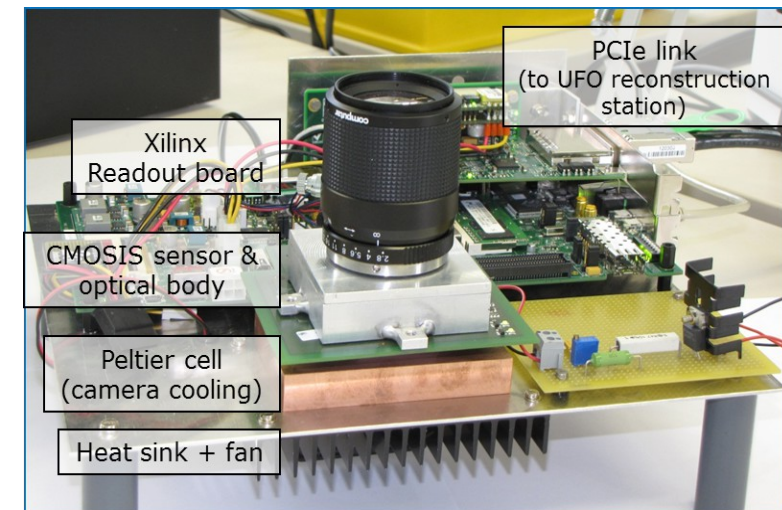# UFO Storage Subsystem



S. Chilingaryan et. all

# High-speed Programmable Camera

- "made in house" Readout mother card
- FPGA Virtex 6
- Peltier cells
- Flex electrical cable (max f = 183MHz)
- Daughter card
- CMOS pixel sensor

Heat exchanger

Link

Vacuum or $N_2$

Optical body and lens

- **High speed CMOS sensor**
- **1Mpix, 5000 fps, 10 bits**
- **Self-trigger & Data compression**
- **On-line elaborations and control**
- **Full Programmability**
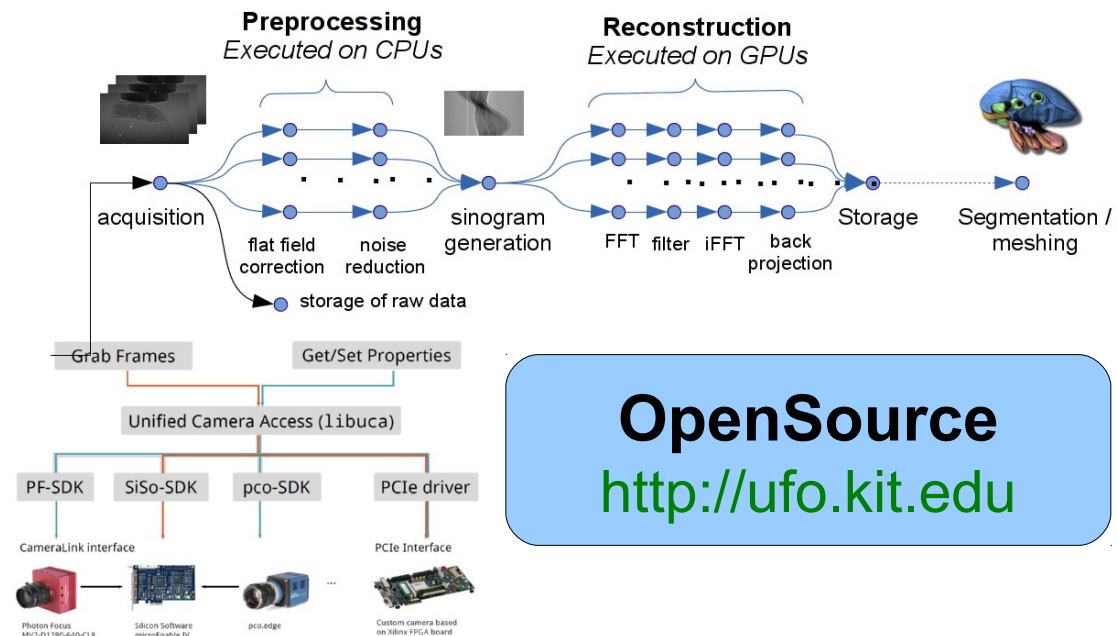- **Direct connection to Infiniband-cluster**



- PCIe link (to UFO reconstruction station)
- Xilinx Readout board
- CMOSIS sensor & optical body
- Peltier cell (camera cooling)
- Heat sink + fan

First Prototype

S. Chilingaryan et. all

# Summary

- We are in the age of parallel architectures
- Getting good performance is rather easy, getting ultimate-performance managing multi-gigabyte streams of data is complicated and needs care on multiple levels. The hardware should be carefully selected according to the planned tasks and data rates. The software should be tunned to the selected hardware.
- Streams about 500 MB/s may be processed with a single reconstruction station, cluster is required to handle more data in near real-time.
- Hybrid CUDA/OpenCL system is probably the best approach.
- UFO Parallel Processing Framework is provided to help you to come along some of these difficulties and will be presented in next talk.

### Features

➢Easy Algorithm Exchange
➢Camera Abstraction
➢Pipelined Processing
➢Glib/GObject, scripting language support with introspection
➢OpenCL + automated management of OpenCL buffers



**OpenSource**
http://ufo.kit.edu