# CERN openlab ICE-DIP project: Data transfer for manycore processors

**Aram Santogidis**

› **23/09/2016**

*Background image: Shutterstock*

# Looking at the sky…

*Background image: Shutterstock*

# Big questions about the Universe



Heavy Elements: 0.03%

Neutrinos: 0.3%

Stars: 0.5%

Free Hydrogen and Helium: 4%

Dark Energy

Dark Matter

Dark Matter: 25%

Dark Energy: 70%

*Background image: Shutterstock*

# Higgs bosson



Particles of the Standard Model

Background image: Shutterstock

# Physics Experiments at CERN



Background image: Shutterstock

Background image: Shutterstock

# Engineers maintaining LHC

# The ATLAS experiment detector
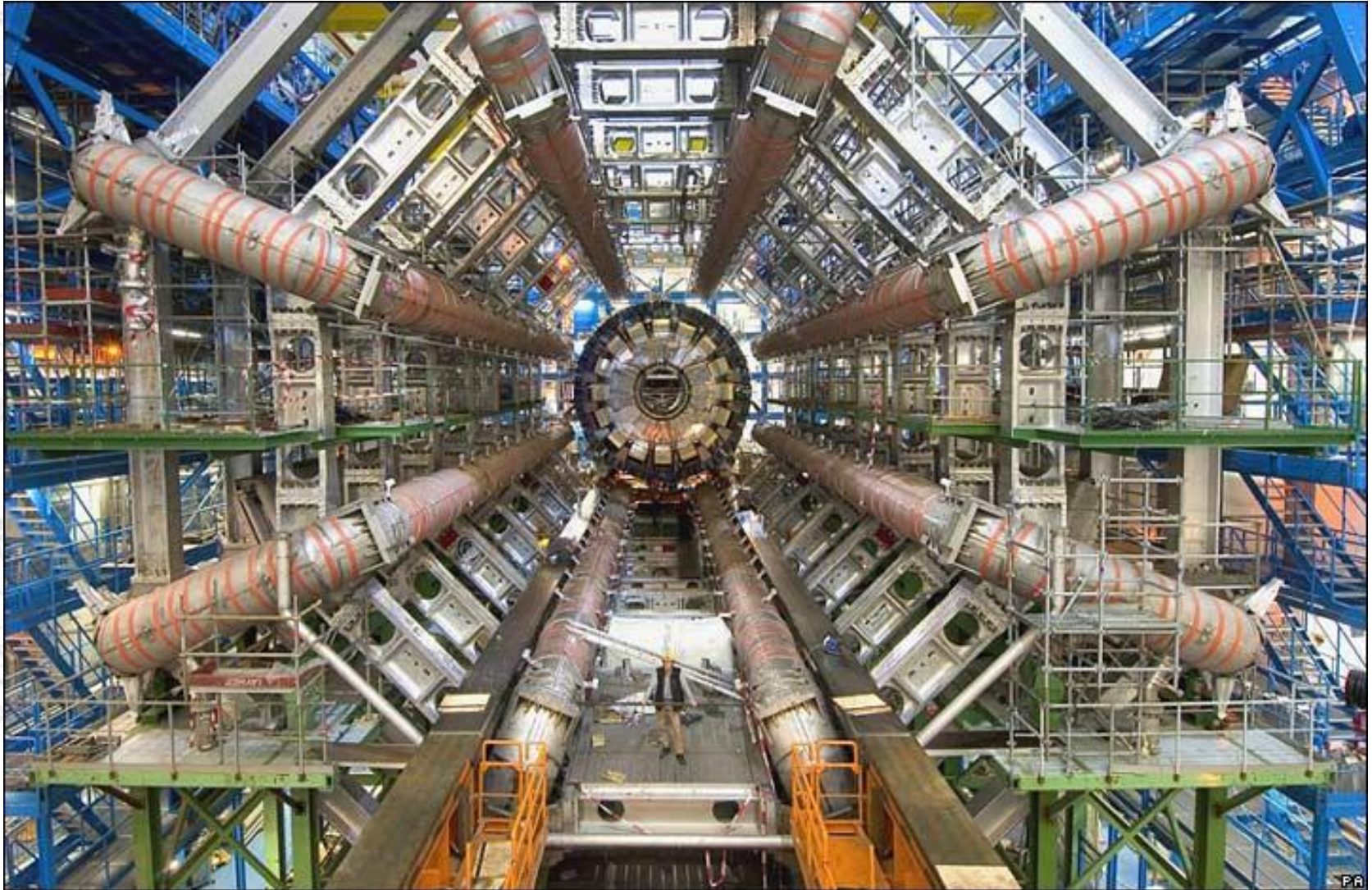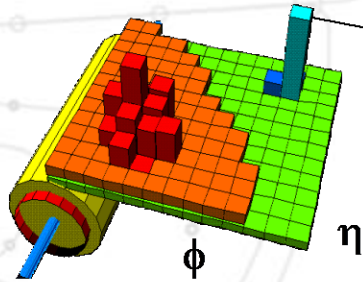
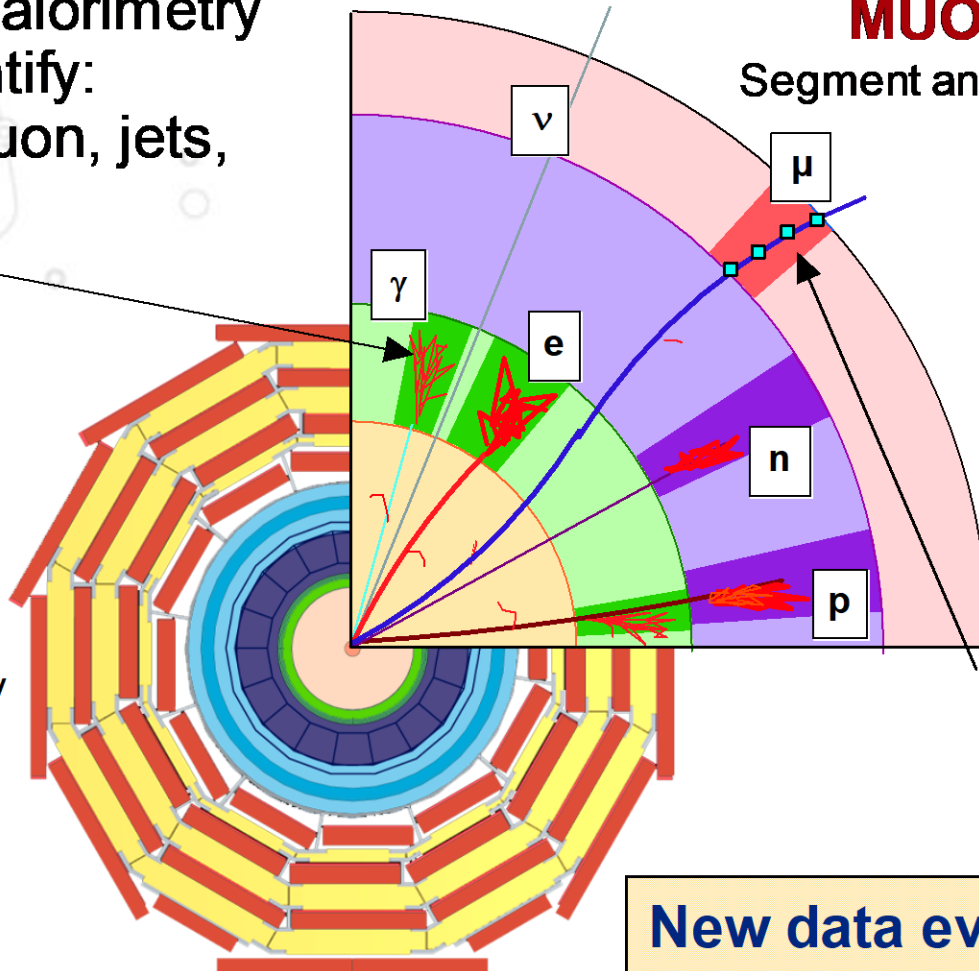# Particle trajectories in the detector

Use prompt data (calorimetry and muons) to identify:
High $p_t$ electron, muon, jets, missing $E_T$

**MUON System**
Segment and track finding

**CALORIMETERs**
Cluster finding and energy deposition evaluation

$\phi$    $\eta$

$\nu$
$\gamma$
$e$
$\mu$
$n$
$p$

**New data every 25 ns
Decision latency ~ µs**

Aram Santogidis – ICE-DIP Project

# A large number of particle trajectories

*Background image: Shutterstock*

# Online processing system

**Detector**    **Front-end electronics**    **Read-out buffers**    **Event builder network**    **High Level Trigger Computing Farm**    **Data Storage**

# 1 PB/s     to     600 MB/s

# ICE-DIP 2013-2017:
# The Intel-CERN European Doctorate Industrial Program
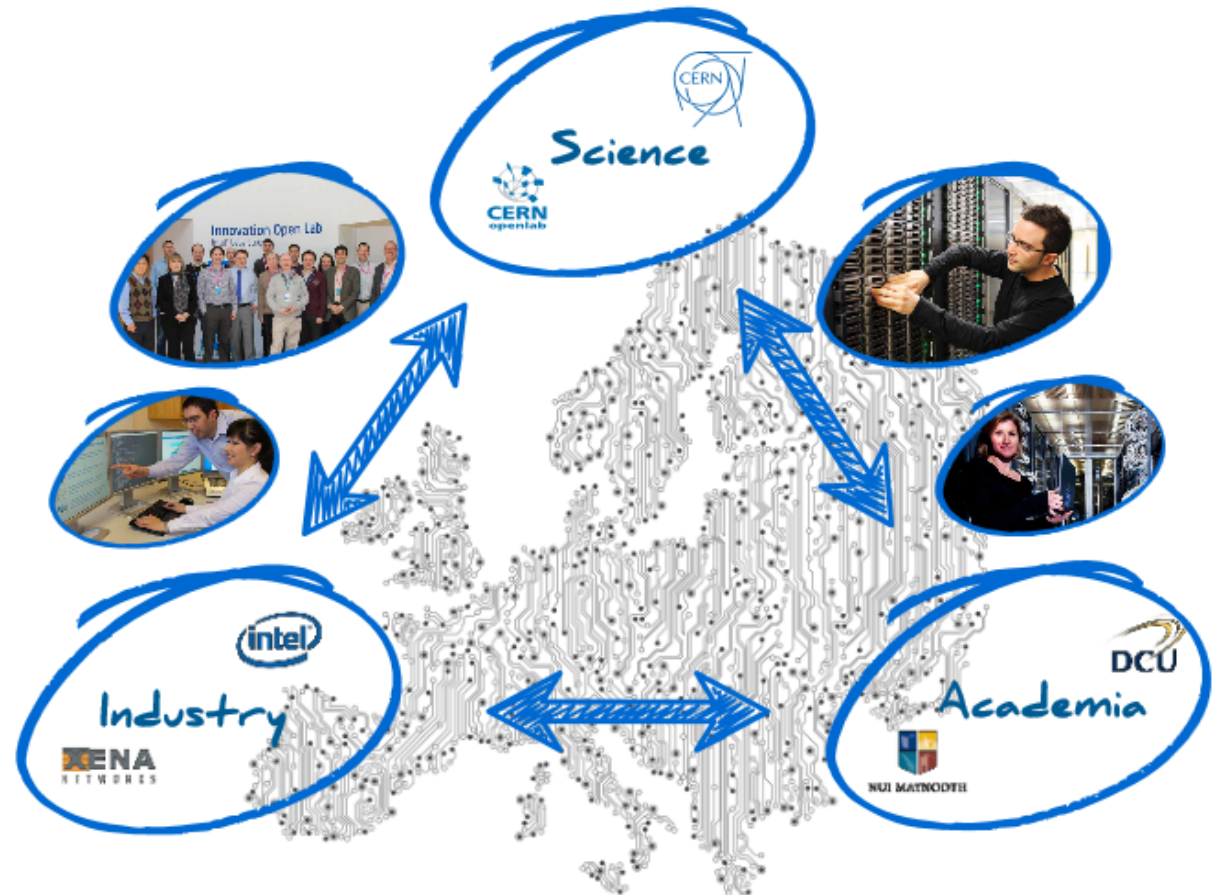
>> A public-private partnership to research solutions for next generation data acquisition networks, offering research training to five Early Stage Researchers in ICT
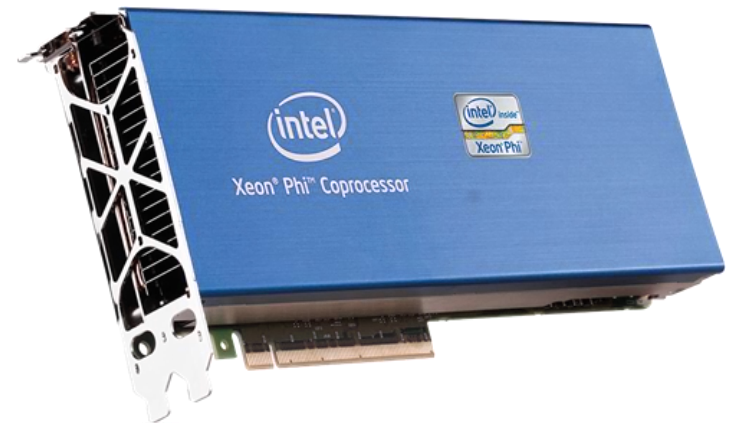
Science

Industry

Academia

Research topics:

- Silicon photonics systems
- Next generation data acquisition networks
- High speed configurable logic
- Computing solutions for high performance data filtering

# Intel Xeon Phi Coprocessor

› **Up to 61 Cores**

› **PCIe card**

› **Different computation modes**

   ▪ Offloading

   ▪ Symmetric

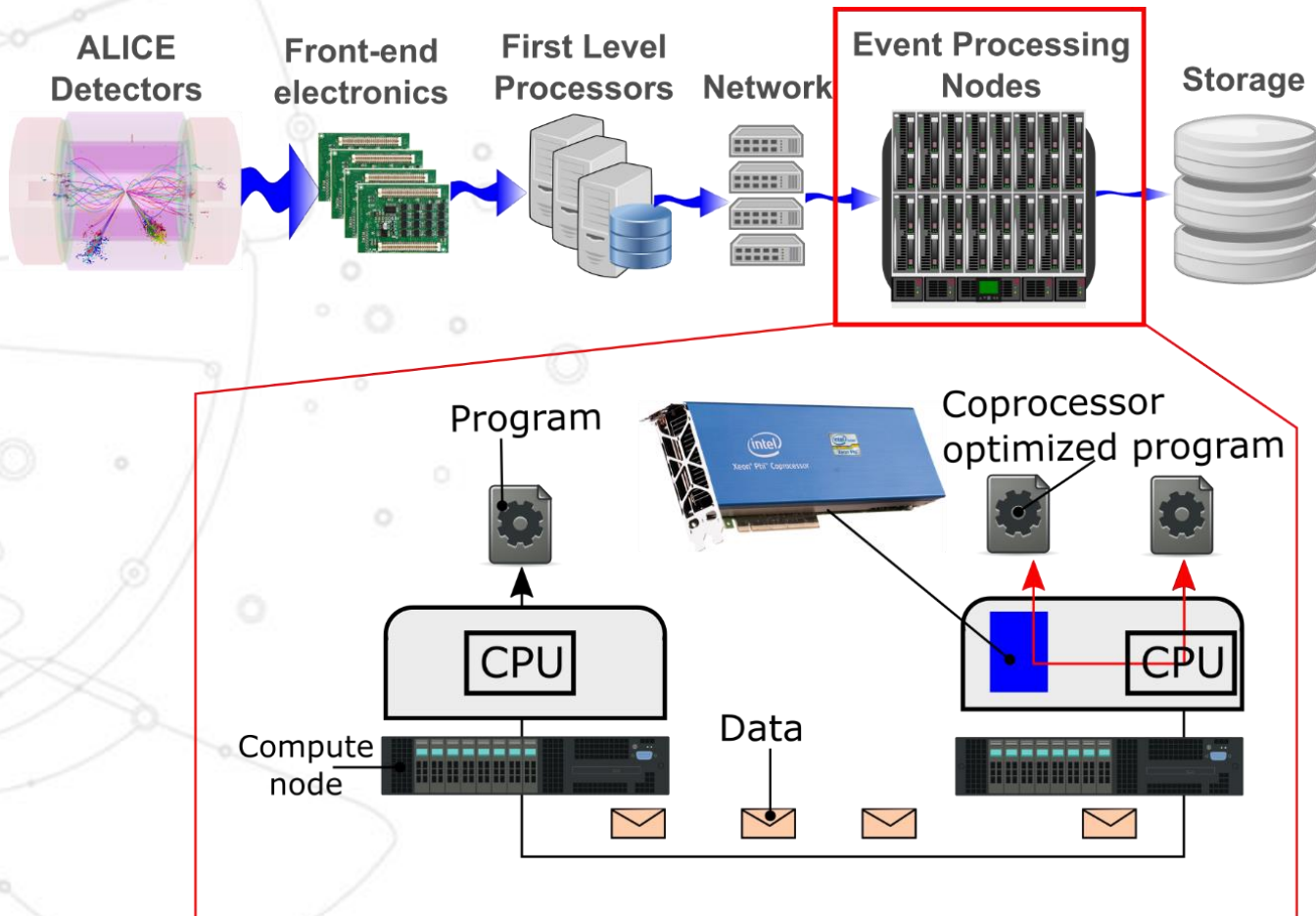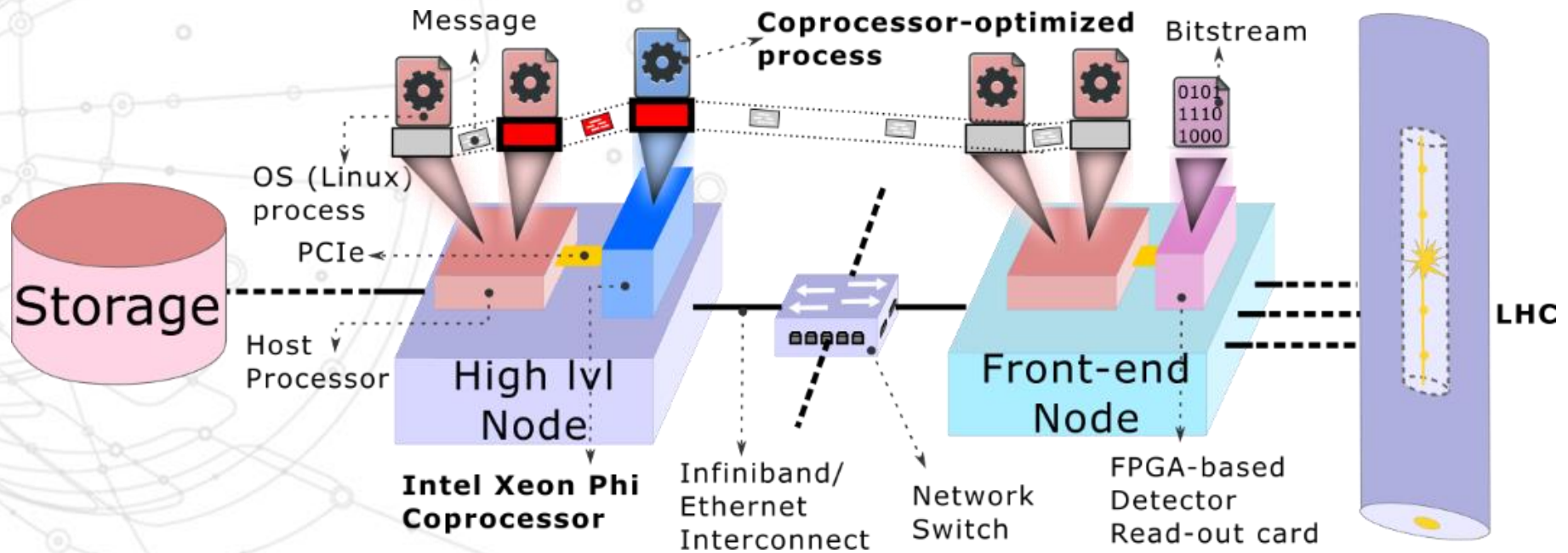# My focus in the online processing pipeline



*Figure 1. An overview of the ALICE data acquisition process with emphasis on the Event Processing Nodes.*

# Coprocessor-Host communication in the online processing system

Background image: Shutterstock

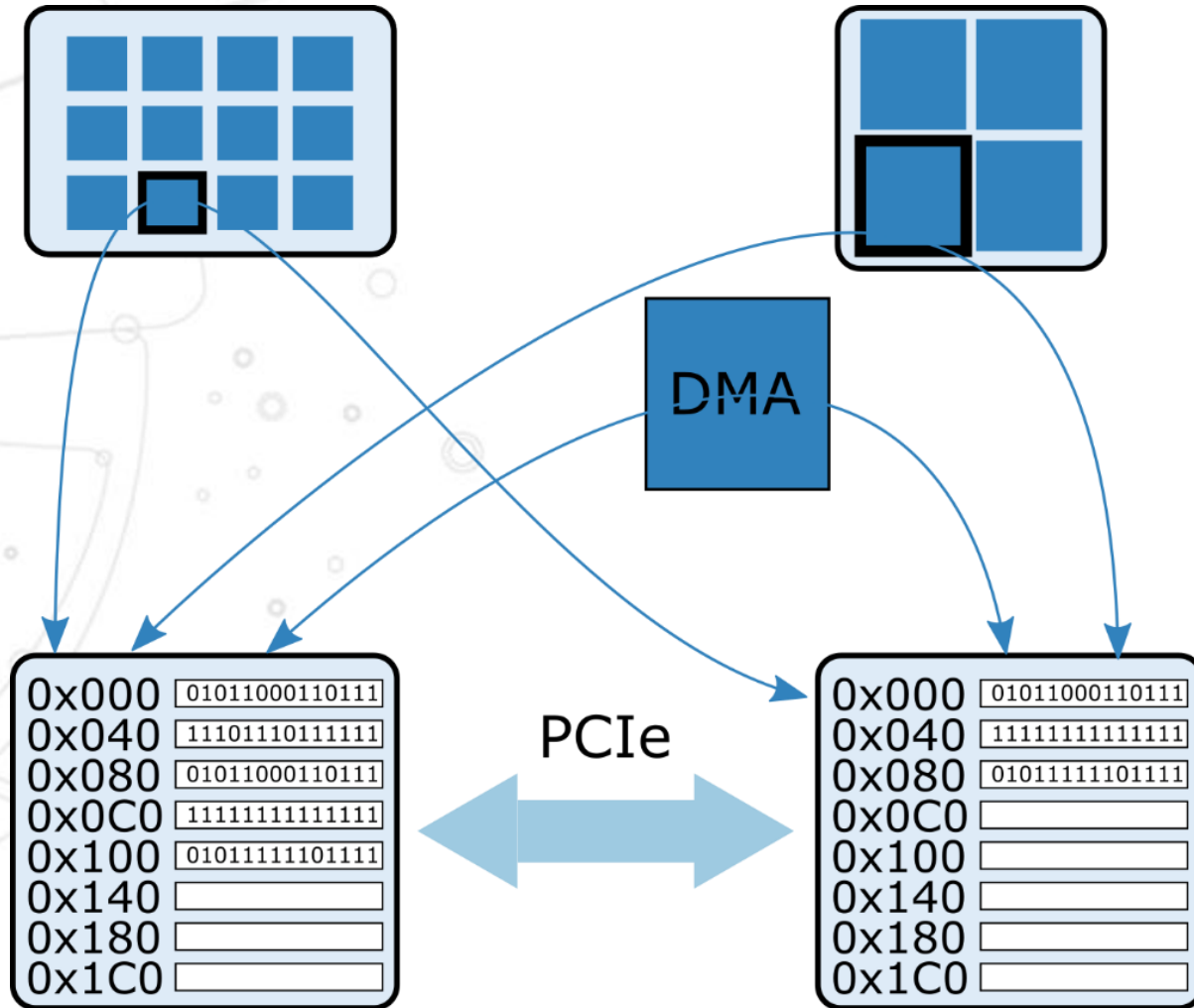# Data transfer libraries for the data acquisition system of the ALICE experiment

> **High level messaging patterns (e.g. pub-sub)**

> **Used to create distributed systems**
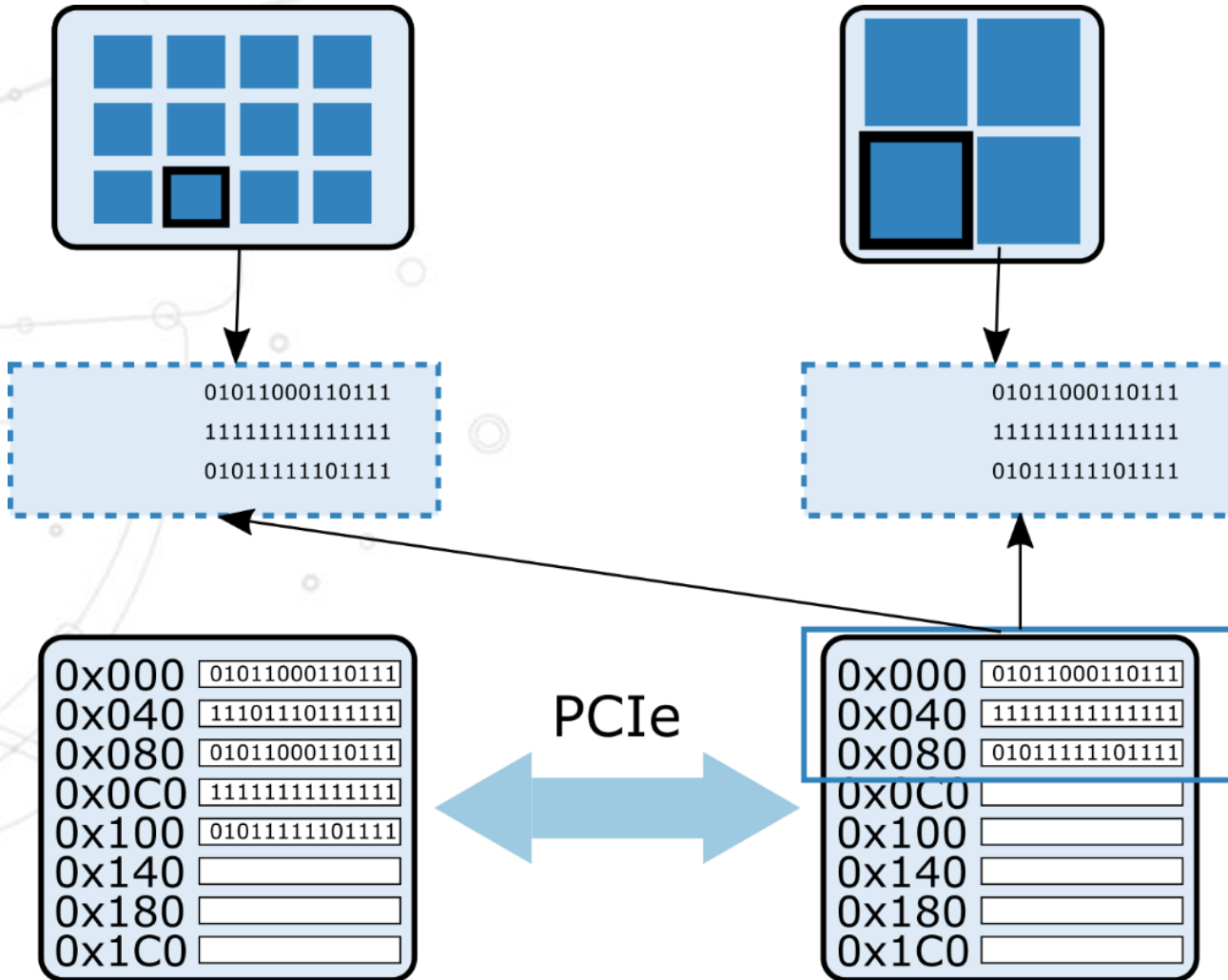
> **Provides good performance**

# Symmetric Communications Interface (SCIF) For Intel® Xeon Phi

› **Communication over PCIe with minimum overhead**

› **POSIX-like interface (listen, connect, send, receive)**

› **RDMA capability (registering memory and read/write to remote address space)**

› **memory mapped IO along the lines of POSIX mmap**

# SCIF mmap

Background image: Shutterstock
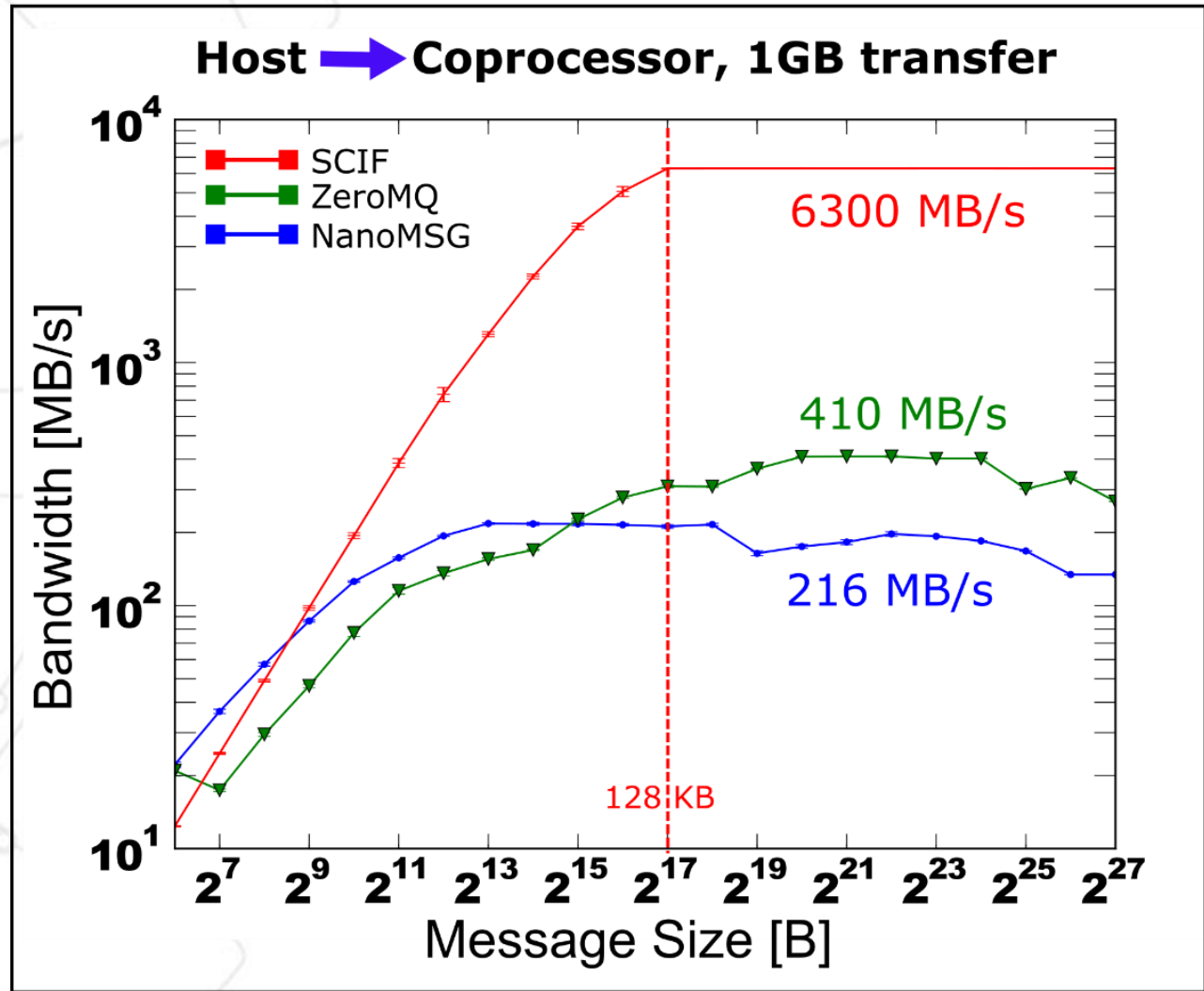
# An opportunity for improvement



**Figure 2** . *This plot corresponds toa performance test of transfer-ring1GB payload in chunks from 4KB to 128MB[1].*

# The features of the new transport mechanism over SCIF

› **Streaming semantics (along the lines of TCP)**

› **Cacheline-aligned RDMA transfers only**

› **Lock-free one-sided communication**

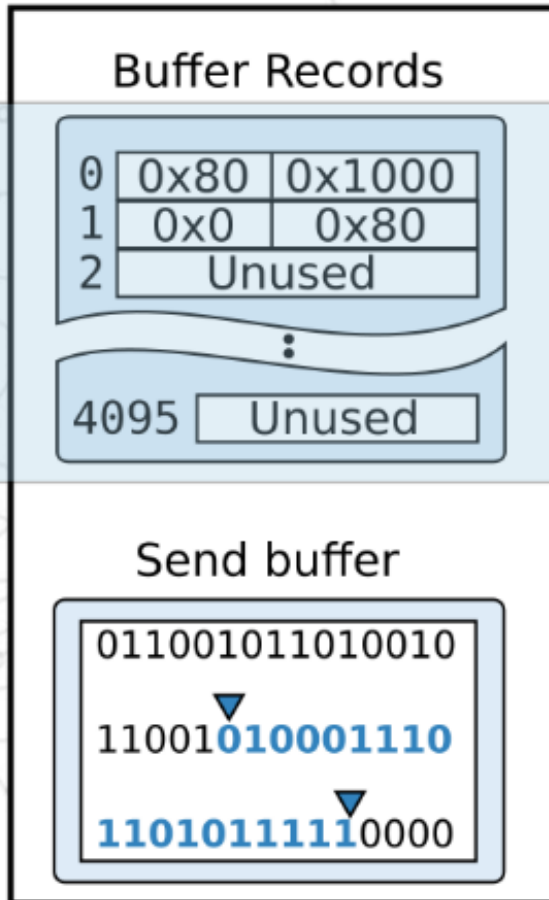› **Maximize data transfer throughput (bandwidth)**

Background image: Shutterstock

# Trans4SCIF synchronization algorithm

**Sender**'s registered address space

**Receiver**'s registered address space

Buffer Records

| 0 | 0x80 | 0x1000 |
|---|------|--------|
| 1 | 0x0 | 0x80 |
| 2 | Unused | |

| 4095 | Unused |
|------|--------|

Writer Records

| 0 | INF | INF |
|---|-----|-----|
| 1 | 0x0 | 0x56 |
| 2 | 0x80 | 0xb1 |

| 4095 | 0x800 | 0xbf0 |
|------|-------|-------|

shared with scif_mmap

Send buffer

011001011010010

11001010001110

110101111110000

Recv buffer

011001011010010

11001010001110

110101111110000

Background image: Shutterstock

# Trans4SCIF synchronization algorithm



0x56 bytes of user data

# Trans4SCIF synchronization algorithm



Sender

Buffer Records

sender_idx ▶

| | | |
|---|---|---|
| 0 | **0x80** | 0x1000 |
| 1 | 0x0 | 0x0 |
| 2 | 0x0 | 0x0 |

⋮

| | |
|---|---|
| 4095 | Unused |

Send buffer

**01100101101001**

**110010**100011101

1101011111000011

0x56 bytes of user data

Receiver

Writer Records

| | | |
|---|---|---|
| 0 | INF | INF |
| 1 | INF | INF |
| 2 | INF | INF |

⋮

| | | |
|---|---|---|
| 4095 | INF | INF |

Recv buffer

**01100101101001**

**110010**10001110

11010111110000

# Trans4SCIF synchronization algorithm



scif_fence_signal(**0x56**)

**Sender**

Buffer Records

sender_idx

| | | |
|---|---|---|
| 0 | **0x80** | 0x1000 |
| 1 | 0x0 | 0x0 |
| 2 | 0x0 | 0x0 |

⋮

| | |
|---|---|
| 4095 | Unused |

Send buffer

**01100101101001**

**110010**100011101

11010111111000011

**Receiver**

Writer Records

sender_idx

| | | |
|---|---|---|
| 0 | **0x0** | INF |
| 1 | INF | INF |
| 2 | INF | INF |

⋮

| | | |
|---|---|---|
| 4095 | INF | INF |

Recv buffer

**01100101101001**

**110010**10001110

11010111110000

0x56 bytes of user data

# Trans4SCIF synchronization algorithm



Copy 0x56 bytes of user data

Background image: Shutterstock

# Trans4SCIF synchronization algorithm



**Sender**

Buffer Records

| 0 | 0x80 | 0x1000 |
| 1 | 0x0 | **0x80** |
| 2 | 0x0 | 0x0 |

sender_idx → 0

receiver_idx → 1

4095 | Unused

Send buffer

```
01100101101001

110010100011101

1101011111000011
```

**Receiver**

Writer Records

| 0 | 0x0 | 0x56 |
| 1 | INF | INF |
| 2 | INF | INF |

receiver_idx → 0

sender_idx → 1

4095 | INF | INF

Recv buffer
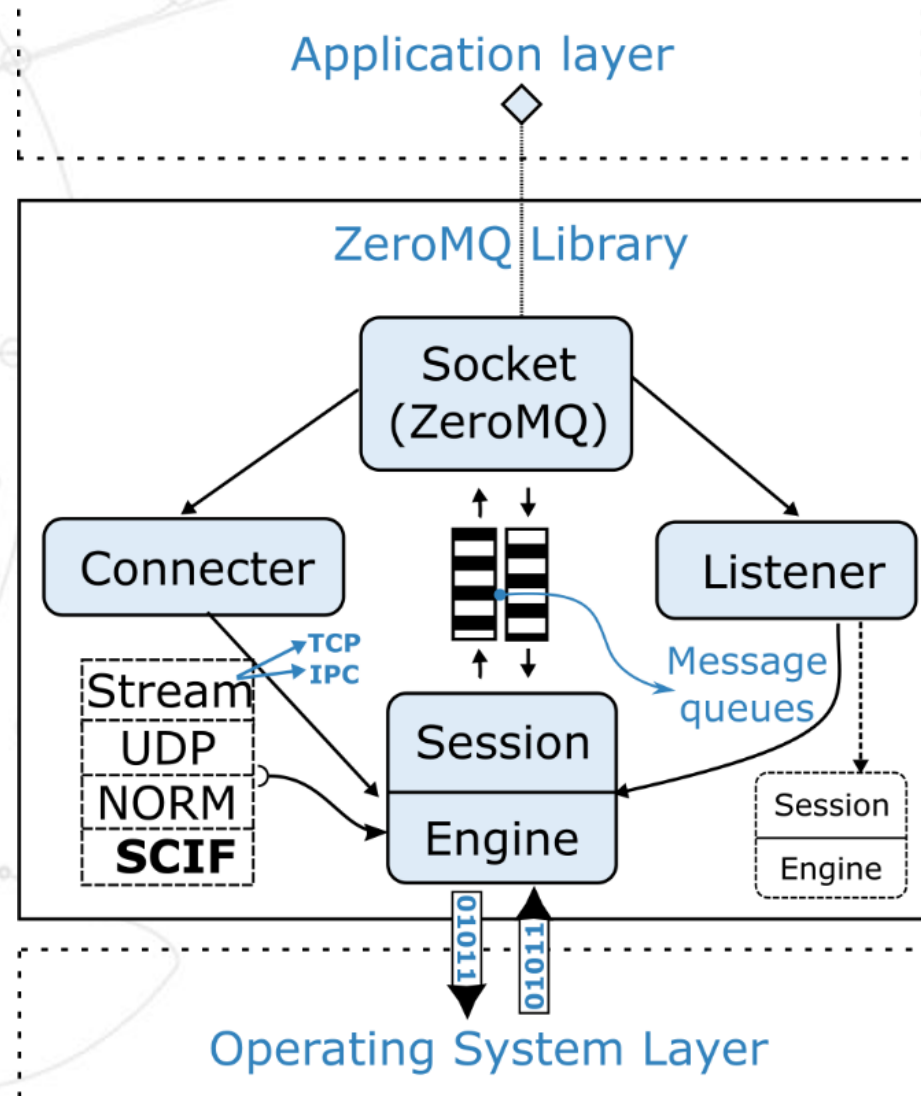
```
01100101101001

110010100011110

11010111110000
```

Copy 0x56 bytes of user data

# The performance of Trans4SCIF



Trans4SCIF (16MB buffer)

Background image: Shutterstock

# The performance of ZeroMQ extend with SCIF via Trans4SCIF

# The Trans4SCIF library the Intel Xeon Phi Coprocessor

› **Easy-to-use socket-like interface**

- Send/Recv

› **E.g. up to 3 GB/s data throughput (4x imporovementwith the cost of 32 MB of memory space**

› **ZeroMQ extension for SCIF**

› **In principle can be re-used by other RDMA based transports (e.g OmniPath)**

# Than you for your attention.

*Background image: Shutterstock*