

# Towards 1000x Speedup for HEP Workloads with Heterogeneous Programmable Datacenters

Anton Burtsev, Alex Veidenbaum

[aburtsev@uci.edu](mailto:aburtsev@uci.edu), [alexv@ics.uci.edu](mailto:alexv@ics.uci.edu)

University of California, Irvine

March, 2018

# A Roadmap for HEP Software and Computing R&D for the 2020s

---

## HEP Software Foundation<sup>1</sup>

ABSTRACT: Particle physics has an ambitious and broad experimental programme for the coming decades. This programme requires large investments in detector hardware, either to build new facilities and experiments, or to upgrade existing ones. Similarly, it requires commensurate investment in the R&D of software to acquire, manage, process, and analyse the shear amounts of data to be recorded. In planning for the HL-LHC in particular, it is critical that all of the collaborating stakeholders agree on the software goals and priorities, and that the efforts complement each other. In this spirit, this white paper describes the R&D activities required to prepare for this software upgrade.

# Compute Ex #1: Exploratory Data Analysis

## Fast Access to Columnar, Hierarchically Nested Data via Code Transformation

Jim Pivarski, Peter Elmer

Physics Department

Princeton University

Princeton, NJ, 08544

pivarski@princeton.edu, peter.elmer@cern.ch

Brian Bockelman, Zhe Zhang

Computer Science and Engineering

University of Nebraska-Lincoln

Lincoln, NE 68588

bbockelm@cse.unl.edu, zhan0915@huskers.unl.edu

***Abstract***—Big Data query systems represent data in a columnar format for fast, selective access, and in some cases (e.g. Apache Drill), perform calculations directly on the columnar data without row materialization, avoiding runtime costs.

However, many analysis procedures cannot be easily or efficiently expressed as SQL. In High Energy Physics, the majority of

HEP data cross-linking is limited to small, disconnected graphs called “events” no larger than hundreds of kilobytes.

Furthermore, most of the well-known database systems use SQL or an SQL variant as the query language, but even the simplest HEP analysis functions are awkward and possibly

# Compute Ex #1: Exploratory Data Analysis

- Dataset:
  - 5.4 million events ( simulated Drell-Yan collisions)
  - Typical analysis will involve 10 such datasets
  - Float:  $5.4 \times 4 = 21.6\text{MB} \times 10 = 216\text{MB}$
  - Double: 432MB

### max p<sub>T</sub> in C++

```
float maximum = 0.0;
for (i=0; i < muons.size(); i++)
    if (muons[i]->pt > maximum)
        maximum = muons[i]->pt;
fill_histogram(maximum);
```

### eta of best by p<sub>T</sub> in C++

```
float maximum = 0.0;
Muon* best = nullptr;
for (i=0; i < muons.size(); i++)
    if (muons[i]->pt > maximum) {
        maximum = muons[i]->pt;
        best = muon; }
if (best != nullptr)
    fill_histogram(best->eta);
```

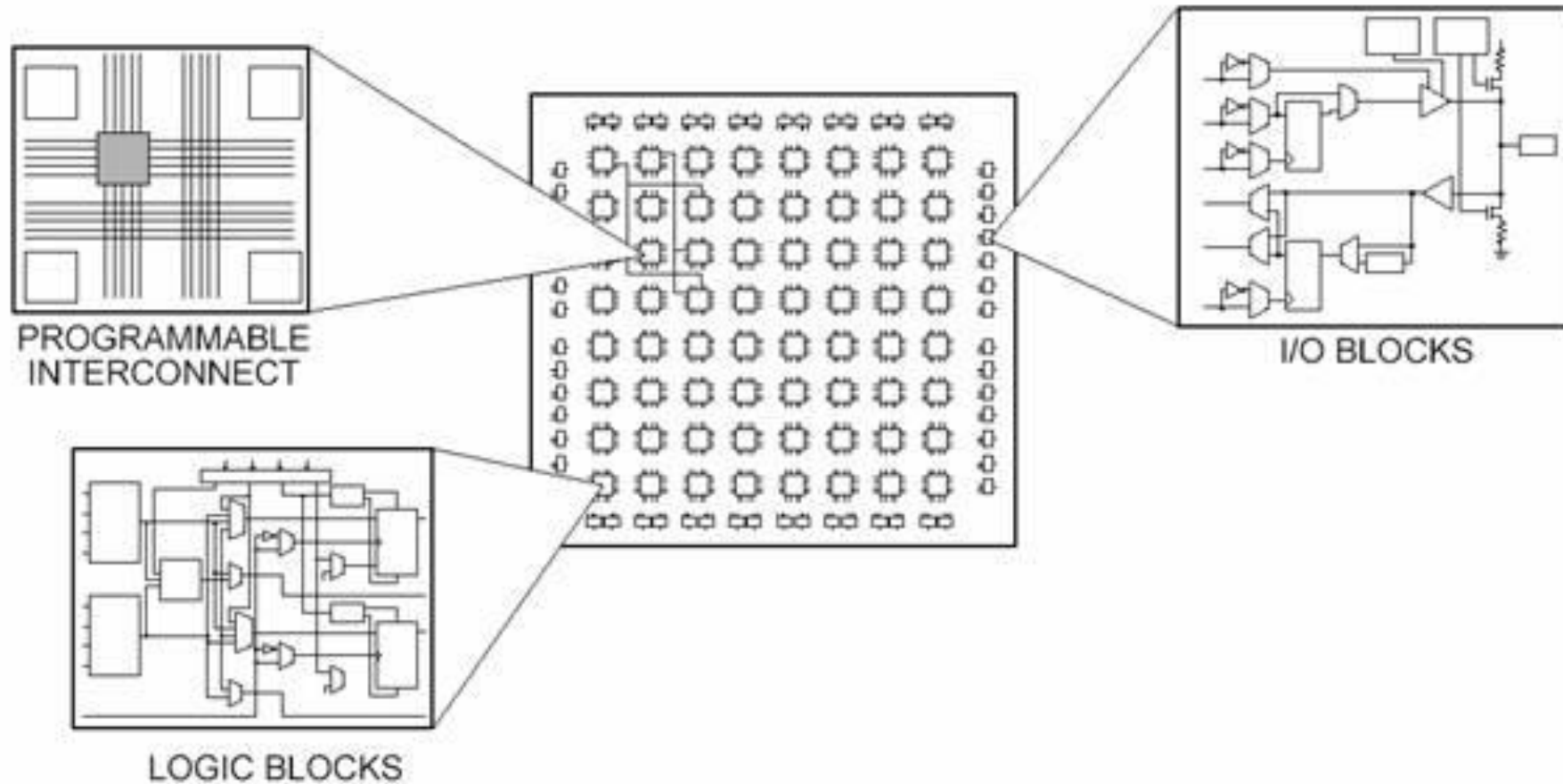
### mass of pairs in C++

```
int n = muons.size();
for (i=0; i < n; i++)
    for (j=i+1; j < n; j++) {
        Muon* m1 = muons[i];
        Muon* m2 = muons[j];
        double mass = sqrt(
            2*m1->pt*m2->pt*(
                - cosh(m1->eta - m2->eta) -
                cos(m1->phi - m2->phi)));
        fill_histogram(mass); }
```

### p<sub>T</sub> sum of pairs in C++

```
int n = muons.size();
for (i=0; i < n; i++)
    for (j=i+1; j < n; j++) {
        Muon* m1 = muons[i];
        Muon* m2 = muons[j];
        double s = m1->pt + m2->pt;
        fill_histogram(s); }
```

# FPGA Filed-programmable gate array





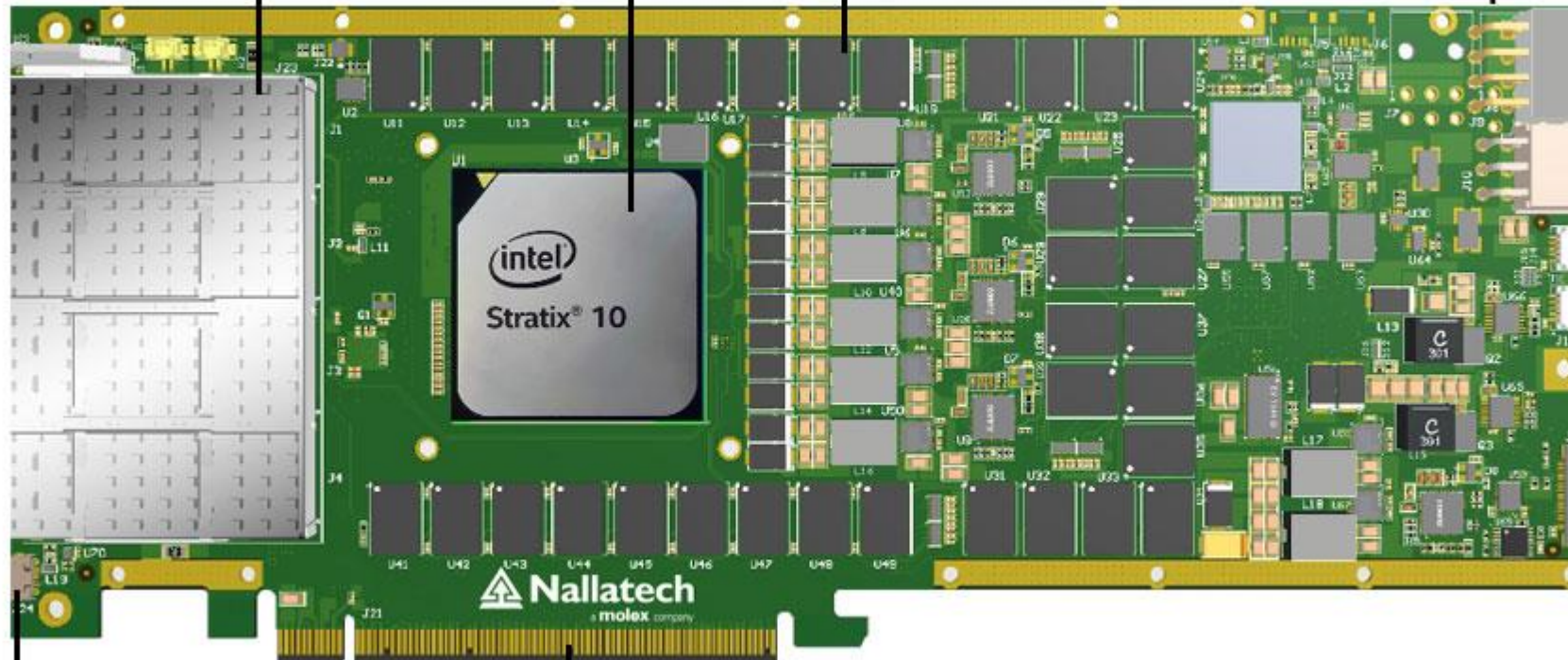
# Intel Stratix 10 FPGA

Network I/O  
4x QSFP28

FPGA  
Intel  
Stratix 10

Memory  
4x 8GB DDR4

12V AUX  
1x 8-pin  
1x 6-pin

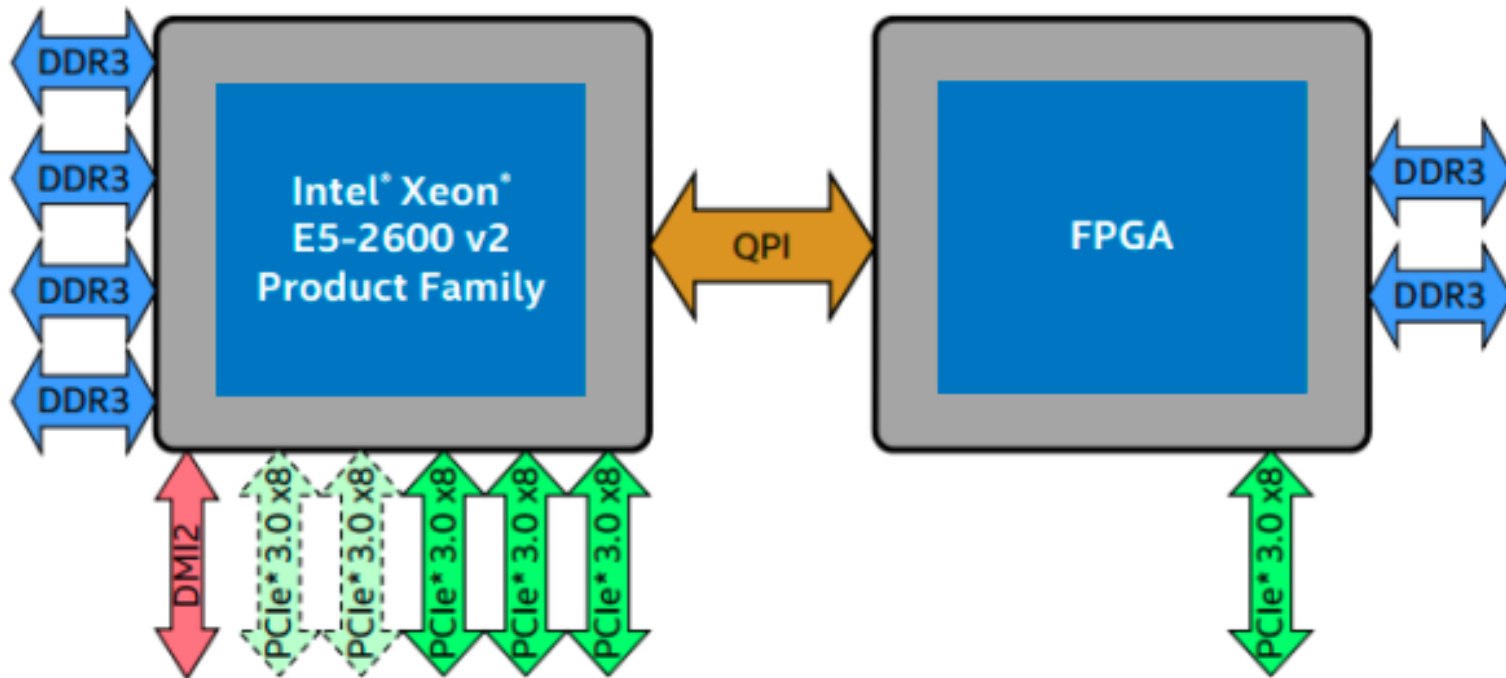


USB #1

PCIe  
16-lane Gen 3

USB #2

# Intel HARP: Cache-coherent FPGA



Processor	Intel® Xeon® E5-26xx v2 Processor
FPGA Module	Altera Stratix V
QPI Speed	6.4 GT/s full width (target 8.0 GT/s at full width)
Memory to FPGA Module	2 channels of DDR3 (up to 64 GB)
Expansion connector to FPGA Module	PCIe 3.0 x8 lanes - maybe used for direct I/O e.g. Ethernet
Features	Configuration Agent, Caching Agent,, (optional) Memory Controller
Software	Accelerator Abstraction Layer (AAL) runtime, drivers, sample applications



### max $p_T$ in C++

```
float maximum = 0.0;
for (i=0; i < muons.size(); i++)
    if (muons[i]->pt > maximum)
        maximum = muons[i]->pt;
fill_histogram(maximum);
```

### eta of best by $p_T$ in C++

```
float maximum = 0.0;
Muon* best = nullptr;
for (i=0; i < muons.size(); i++)
    if (muons[i]->pt > maximum) {
        maximum = muons[i]->pt;
        best = muon; }
if (best != nullptr)
    fill_histogram(best->eta);
```

### mass of pairs in C++

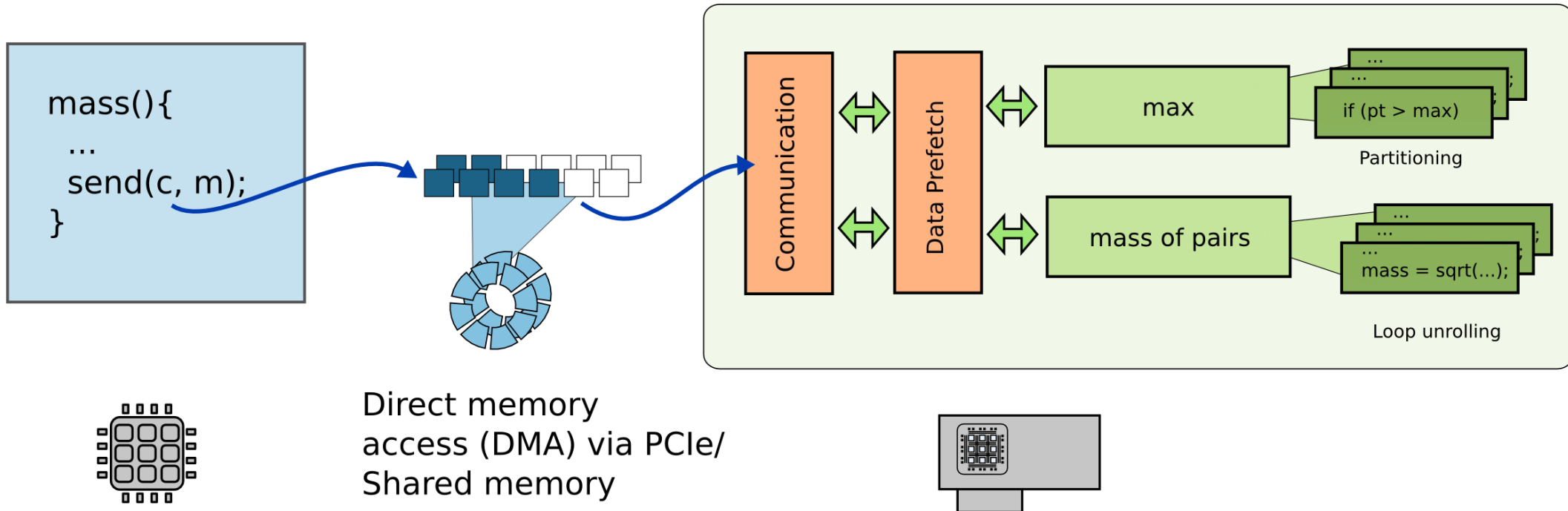
```
int n = muons.size();
for (i=0; i < n; i++)
    for (j=i+1; j < n; j++) {
        Muon* m1 = muons[i];
        Muon* m2 = muons[j];
        double mass = sqrt(
            2*m1->pt*m2->pt*(
                - cosh(m1->eta - m2->eta) -
                cos(m1->phi - m2->phi)));
        fill_histogram(mass); }
```

### $p_T$ sum of pairs in C++

```
int n = muons.size();
for (i=0; i < n; i++)
    for (j=i+1; j < n; j++) {
        Muon* m1 = muons[i];
        Muon* m2 = muons[j];
        double s = m1->pt + m2->pt;
        fill_histogram(s); }
```

# FPGA acceleration

- Parallel pipelines
  - Partition the input
  - Unroll loops
- Reconfigurable with partial reconfiguration

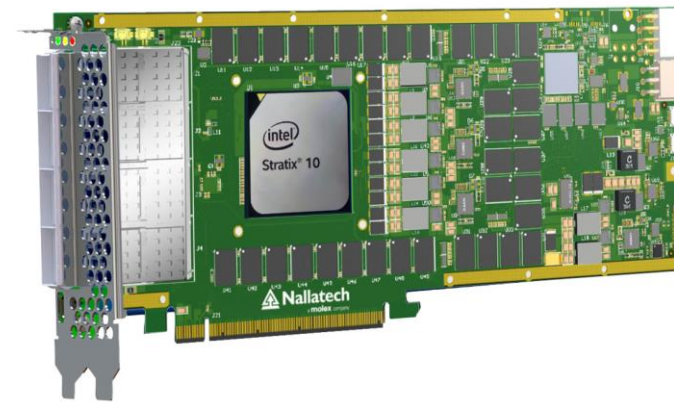


# FPGA vs GPU

- NVidia Tesla V100 GPU
  - 15 TFLOP single point
  - 60GFLOP per watt



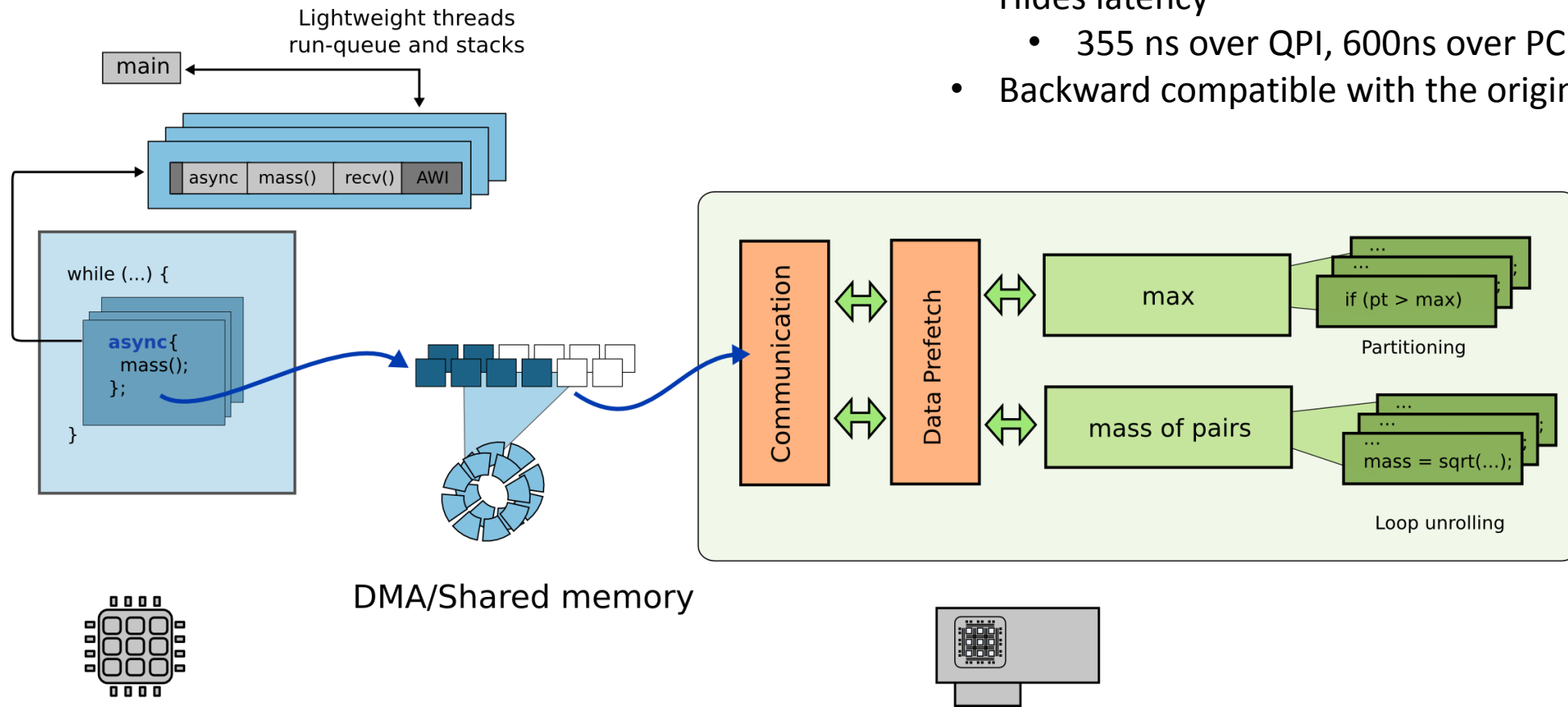
- Intel Stratix 10 FPGA
  - 10 TFLOP single point
  - 80GFLOP per watt



# More control

- Low-latency communication via DMA or shared memory with the main program
  - Simple ring-buffer optimized for the number of cache-coherence or PCIe transactions
- Data prefetching from the host (CPU) and device (FPGA) memories and even from NVMe
- Direct communication over the network and with NVMe

# Integration with existing programs: asynchronous runtime

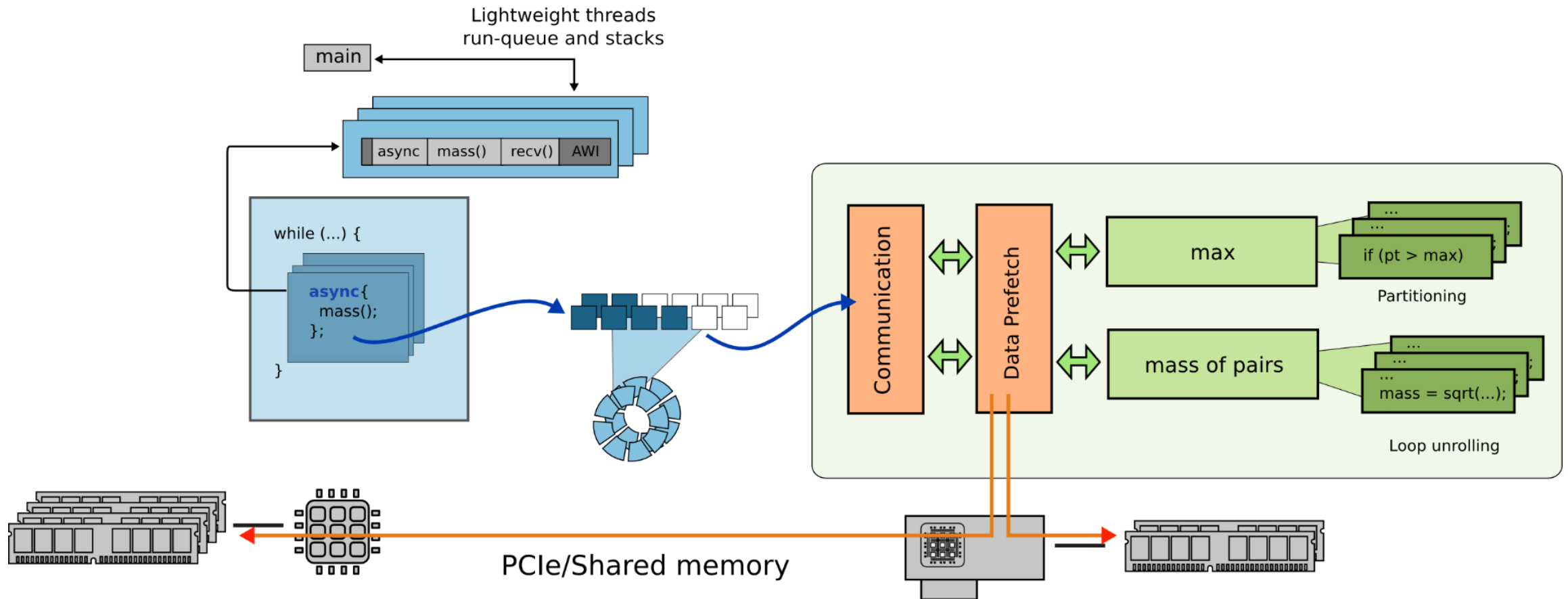


- Hides latency
  - 355 ns over QPI, 600ns over PCIe
- Backward compatible with the original code



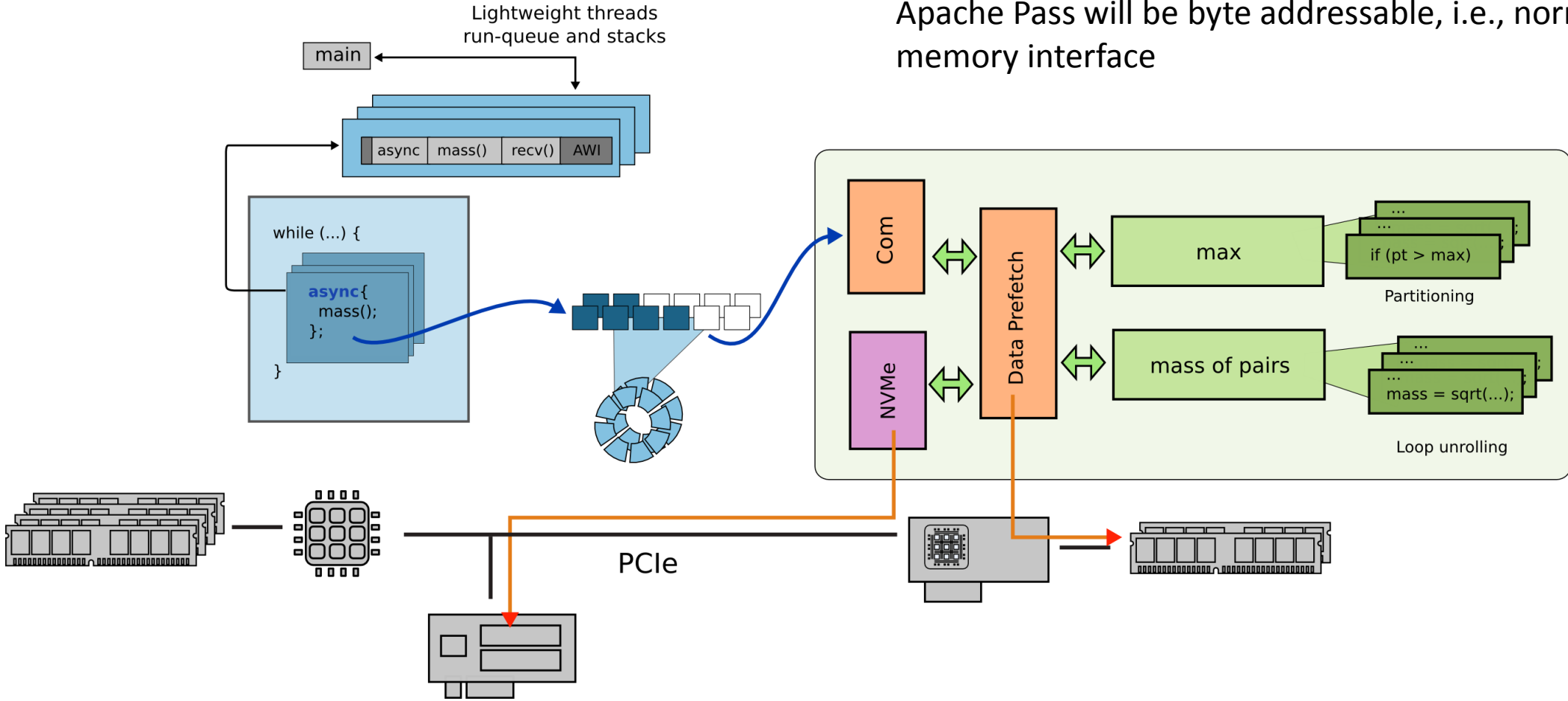
# Data prefetching

- FPGA has
  - 6MB of fast block RAM
  - 4GB of DRAM
- Program a custom prefetch logic that is aware of the data layout

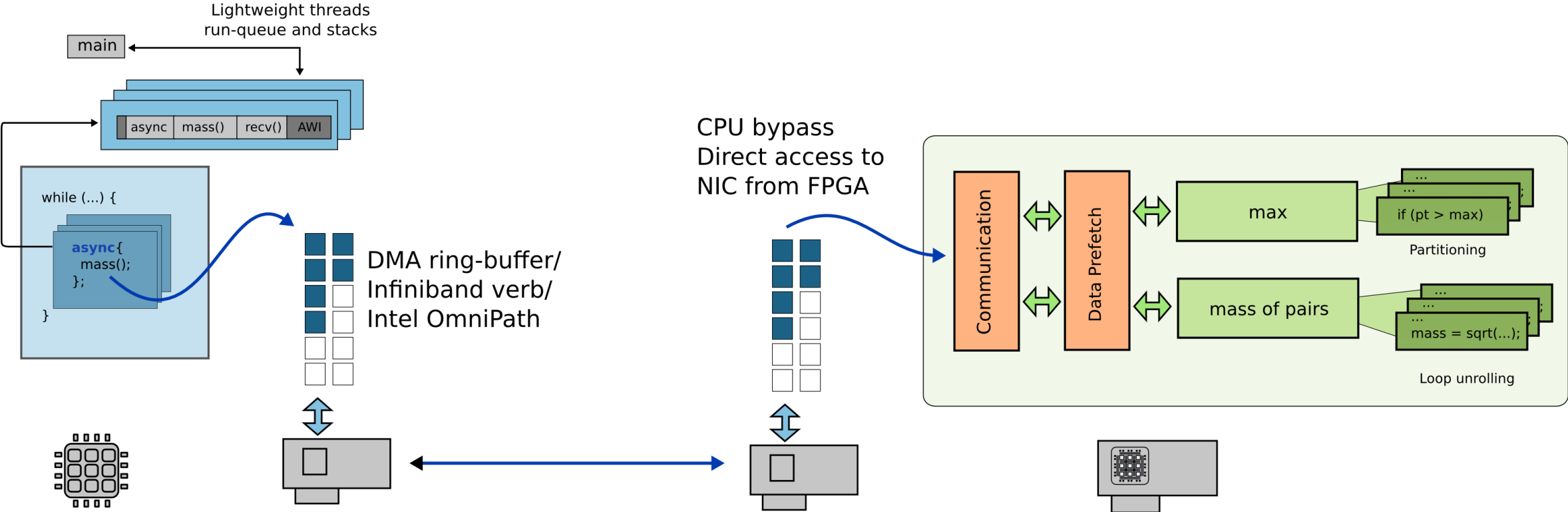


# Direct access to NVMe

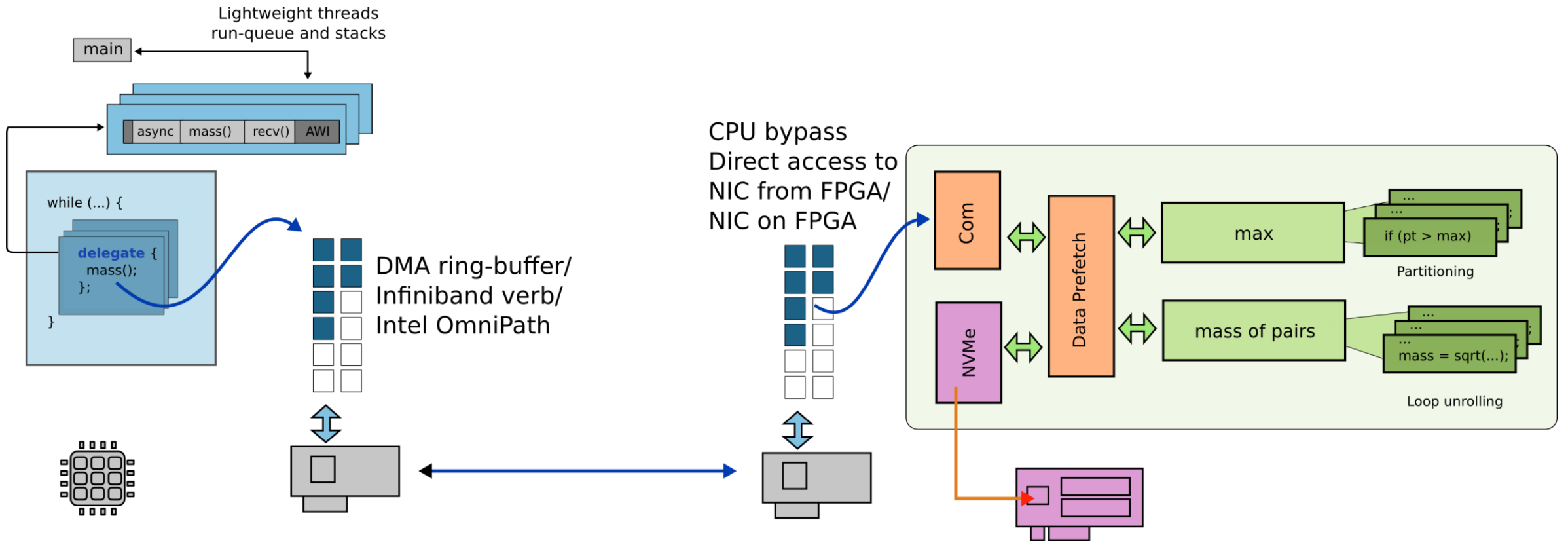
- Direct access to storage devices
  - NVMe is a simple ring-based protocol
  - Easy to program in FPGA
- Emerging non-volatile DIMMs, e.g., Intel 3D Xpoint Apache Pass will be byte addressable, i.e., normal memory interface



# Remote access over the network

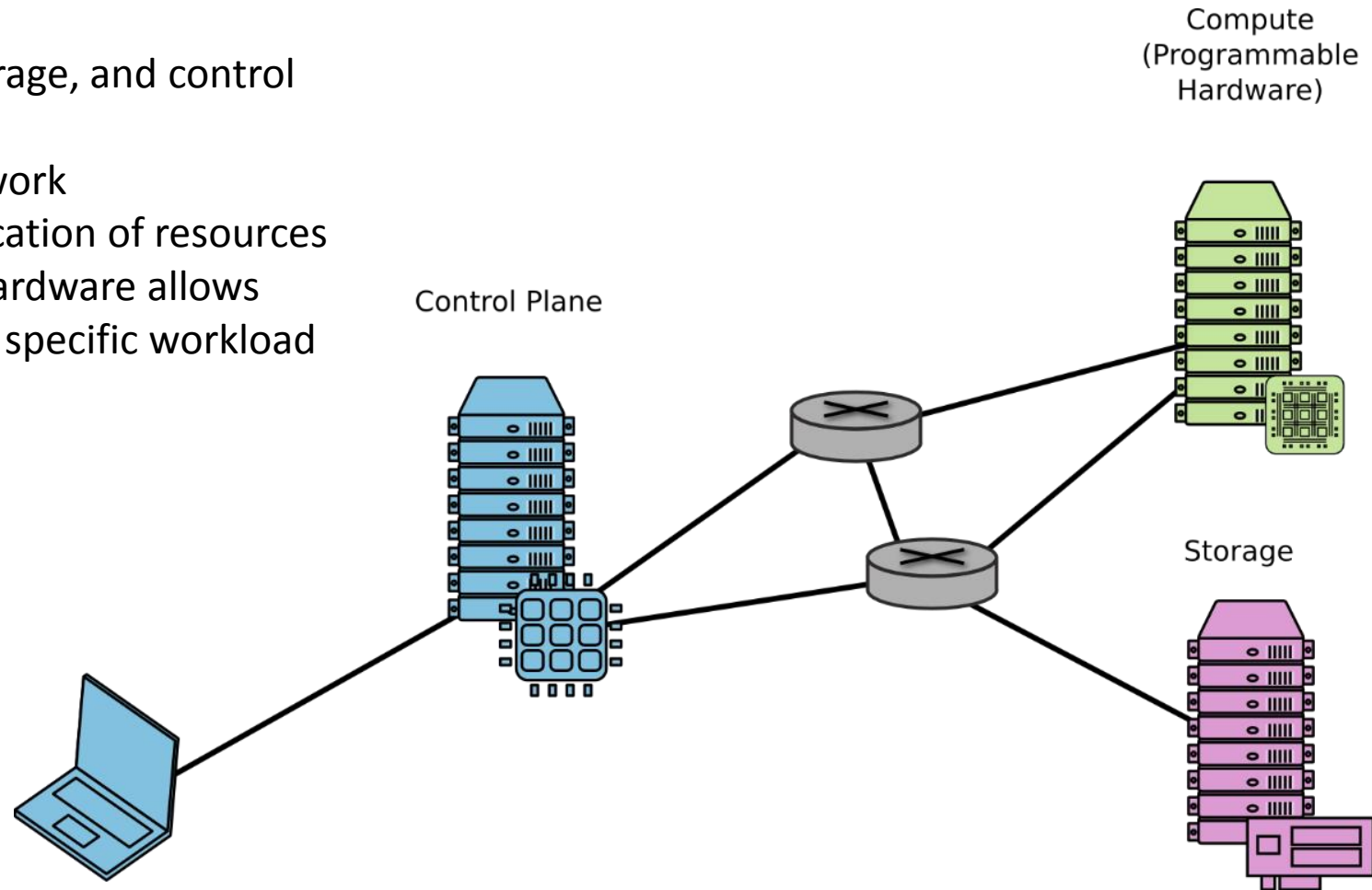


# Collocating compute and storage



# Disaggregated programmable datacenter

- Pools of compute, storage, and control plane servers
  - Low-latency network
- Flexible, dynamic allocation of resources
  - Programmable hardware allows optimization of a specific workload





Example applications

# Parallelized Kalman-Filter-Based Reconstruction of Particle Tracks on Many-Core Processors and GPUs

Giuseppe **Cerati**<sup>4,a</sup>, Peter **Elmer**<sup>2,b</sup>, Slava **Krutelyov**<sup>1,c</sup>, Steven **Lantz**<sup>3,d</sup>, Matthieu **Lefebvre**<sup>2,e</sup>, Mario **Masciovecchio**<sup>1,f</sup>, Kevin **McDermott**<sup>3,g</sup>, Daniel **Riley**<sup>3,h</sup>, Matevž **Tadel**<sup>1,i</sup>, Peter **Wittich**<sup>3,j</sup>, Frank **Würthwein**<sup>1,k</sup>, and Avi **Yagil**<sup>1,l</sup>

<sup>1</sup> *UC San Diego, La Jolla, CA, United States of America 92093*

<sup>2</sup> *Princeton University, Princeton, NJ, United States of America 08544*

<sup>3</sup> *Cornell University, Ithaca, NY, United States of America 14853*

<sup>4</sup> *Fermilab, Batavia, IL, United States of America 60510-5011*

**Abstract.** For over a decade now, physical and energy constraints have limited clock speed improvements in commodity microprocessors. Instead, chipmakers have been pushed into producing lower-power, multi-core processors such as Graphical Processing Units (GPU), ARM CPUs, and Intel MICs. Broad-based efforts from manufacturers and developers have been devoted to making these processors user-friendly enough to



ELSEVIER

Available online at [www.sciencedirect.com](http://www.sciencedirect.com)

SCIENCE @ DIRECT®

Nuclear Instruments and Methods in Physics Research A 506 (2003) 250–303

**NUCLEAR  
INSTRUMENTS  
& METHODS  
IN PHYSICS  
RESEARCH**  
Section A

[www.elsevier.com/locate/nima](http://www.elsevier.com/locate/nima)

## GEANT4—a simulation toolkit

S. Agostinelli<sup>ae</sup>, J. Allison<sup>as,\*</sup>, K. Amako<sup>e</sup>, J. Apostolakis<sup>a</sup>, H. Araujo<sup>aj</sup>,  
P. Arce<sup>l,m,x,a</sup>, M. Asai<sup>g,ai</sup>, D. Axen<sup>i,t</sup>, S. Banerjee<sup>bi,l</sup>, G. Barrand<sup>an</sup>, F. Behner<sup>l</sup>,  
L. Bellagamba<sup>c</sup>, J. Boudreau<sup>bd</sup>, L. Broglia<sup>ar</sup>, A. Brunengo<sup>c</sup>, H. Burkhardt<sup>a</sup>,  
S. Chauvie<sup>bj,bl</sup>, J. Chuma<sup>h</sup>, R. Chytrcek<sup>a</sup>, G. Cooperman<sup>az</sup>, G. Cosmo<sup>a</sup>,  
P. Degtyarenko<sup>d</sup>, A. Dell'Acqua<sup>a,i</sup>, G. Depaola<sup>y</sup>, D. Dietrich<sup>af</sup>, R. Enami<sup>ab</sup>,  
A. Feliciello<sup>bj</sup>, C. Ferguson<sup>bh</sup>, H. Fesefeldt<sup>l,o</sup>, G. Folger<sup>a</sup>, F. Foppiano<sup>ac</sup>,  
A. Forti<sup>as</sup>, S. Garelli<sup>ac</sup>, S. Giani<sup>a</sup>, R. Giannitrapani<sup>bo</sup>, D. Gibin<sup>m,bc</sup>, J.J. Gómez

# Cling – The New Interactive Interpreter for ROOT 6

V Vasilev<sup>1</sup>, Ph Canal<sup>2</sup>, A Naumann<sup>1</sup>, P Russo<sup>2</sup>

<sup>1</sup> CERN, PH-SFT, Geneva, Switzerland

<sup>2</sup> FermiLab, P.O. Box 500 Batavia, IL, US

E-mail: vasil.georgiev.vasilev@cern.ch

**Abstract.** Cling is an interactive C++ interpreter, built on top of Clang and LLVM compiler infrastructure. Like its predecessor Cint, Cling realizes the *read-print-evaluate-loop* concept, in order to leverage rapid application development. Implemented as a small extension to LLVM and Clang, the interpreter reuses their strengths such as the praised concise and expressive compiler diagnostics. We show how to match the interpreter concept to the compiler library and generalize common set of requirements for building up an interactive interpreter. We reason the design and implementation decisions as solution to the challenge of implementing interpreter behaviour as an extension of the compiler library. We present the new features, e.g. how C++11 will come to Cling and how Cint-specific extensions are being adopted. We clarify the state of integration in the ROOT framework and the induced change set. We explain how ROOT dictionaries are simplified due to the new interpreter.

# Discussion

- We need help with understanding
  - Sizes of HEP datasets
  - Shape of the computation, e.g., similar to mass of pairs, but for Kalman Filter and Monte Carlo