



# Preparing for Extreme Heterogeneity in High Performance Computing

Jeffrey S. Vetter

*With many contributions from FTG Group and Colleagues*

Rocky Mountain Advanced Computing Consortium Symposium  
Boulder, Colorado  
22 May 2019



ORNL is managed by UT-Battelle, LLC for the US Department of Energy

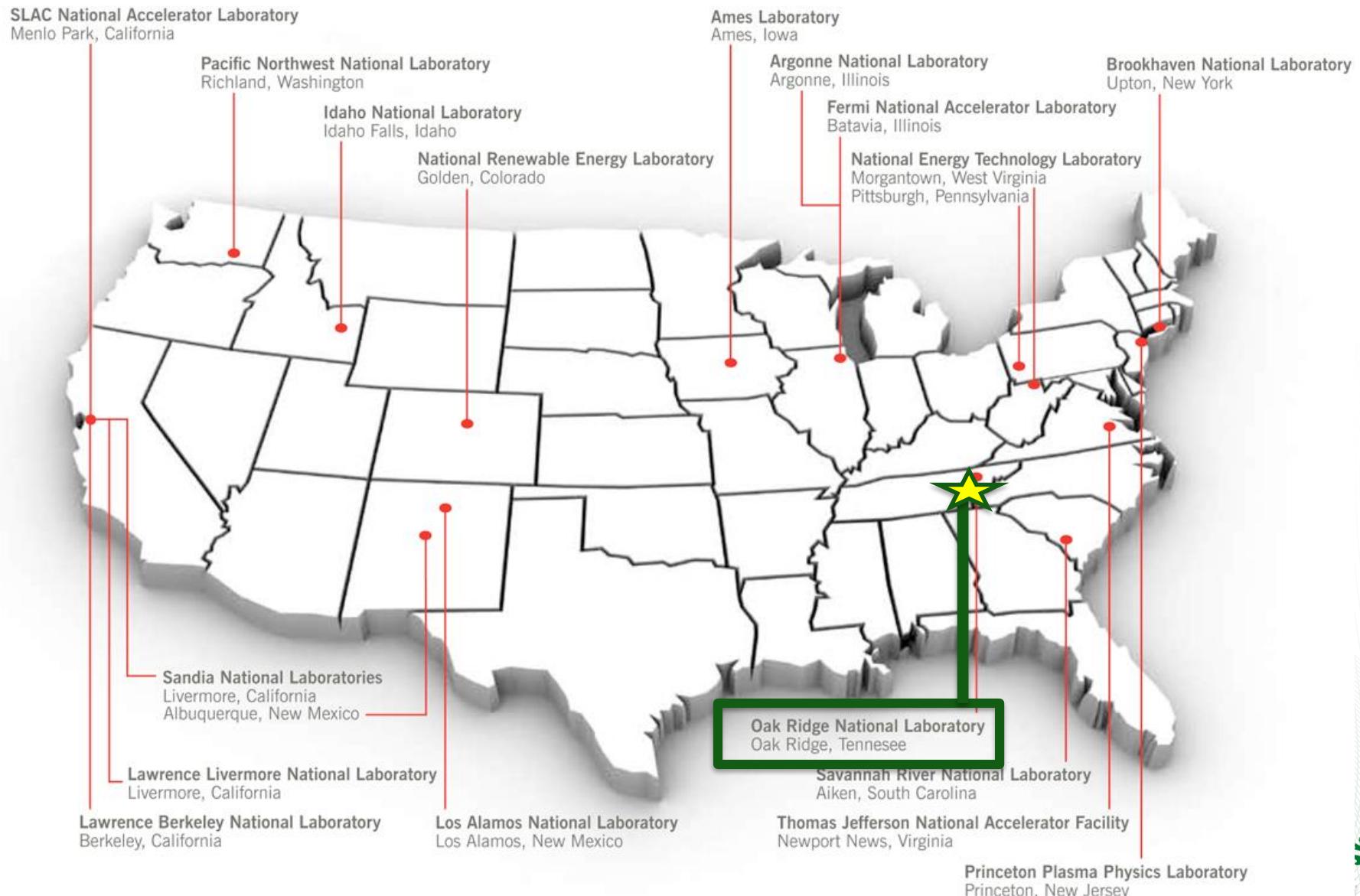


# Highlights

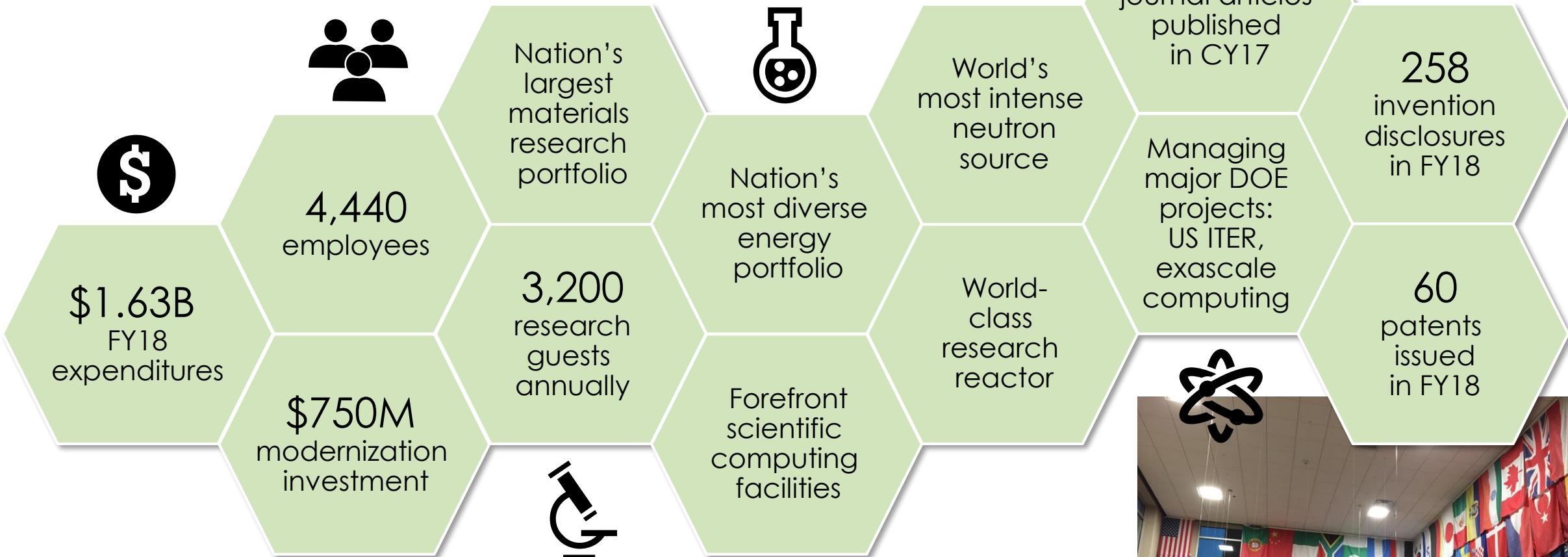
- Recent trends in extreme-scale HPC paint an ambiguous future
  - Contemporary systems provide evidence that power constraints are driving architectures to change rapidly
  - Multiple architectural dimensions are being (dramatically) redesigned: Processors, node design, memory systems, I/O
  - Complexity is our main challenge
- Applications and software systems are all reaching a state of crisis
  - Applications will not be functionally or performance portable across architectures
  - Programming and operating systems need major redesign to address these architectural changes
  - Procurements, acceptance testing, and operations of today's new platforms depend on performance prediction and benchmarking.
- We need portable programming models and performance prediction now more than ever!
  - Heterogeneous processing
    - OpenACC->FPGAs
    - Clacc – OpenACC support in LLVM (not covered today)
  - Emerging memory hierarchies (NVM)
    - DRAGON – transparent NVM access from GPUs
    - NVL-C – user management of nonvolatile memory in C
    - Papyrus – parallel aggregate persistent storage (not covered today)
- Performance prediction is critical for design and optimization (not covered today)

# Very Brief Introduction to ORNL

# Oak Ridge National Laboratory is the DOE Office of Science's Largest Lab



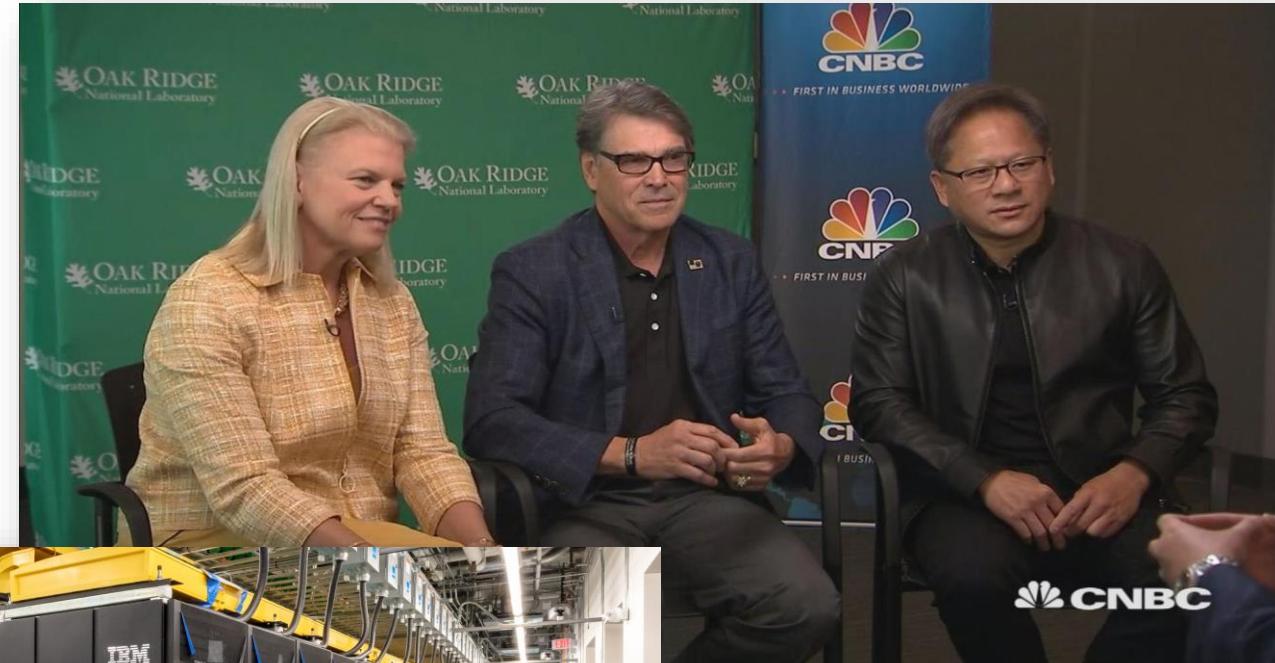
# Today, ORNL is a leading science and energy laboratory



# ORNL 75<sup>th</sup> Lab Day and Summit Unveiling – 8 June 2018

#1 on Top 500

<b>Application Performance</b>	<b>200 PF</b>
<b>Number of Nodes</b>	<b>4,608</b>
<b>Node performance</b>	<b>42 TF</b>
<b>Memory per Node</b>	<b>512 GB DDR4 + 96 GB HBM2</b>
<b>NV memory per Node</b>	<b>1600 GB</b>
<b>Total System Memory</b>	<b>&gt;10 PB DDR4 + HBM2 + Non-volatile</b>
<b>Processors</b>	<b>2 IBM POWER9™ 9,216 CPUs 6 NVIDIA Volta™ 27,648 GPUs</b>
<b>File System</b>	<b>250 PB, 2.5 TB/s, GPFS™</b>
<b>Power Consumption</b>	<b>13 MW</b>
<b>Interconnect</b>	<b>Mellanox EDR 100G InfiniBand</b>
<b>Operating System</b>	<b>Red Hat Enterprise Linux (RHEL) version 7.4</b>



# U.S. Department of Energy and Cray to Deliver Record-Setting Frontier Supercomputer at ORNL

Exascale system expected to be world's most powerful computer for science and innovation

Topic: Supercomputing

May 7, 2019



OAK RIDGE, Tenn., May 7, 2019—The U.S. Department of Energy today announced a contract with Cray Inc. to build the Frontier supercomputer at Oak Ridge National Laboratory, which is anticipated to debut in 2021 as the world's most powerful computer with a peak performance of greater than 1.5 exaflops.

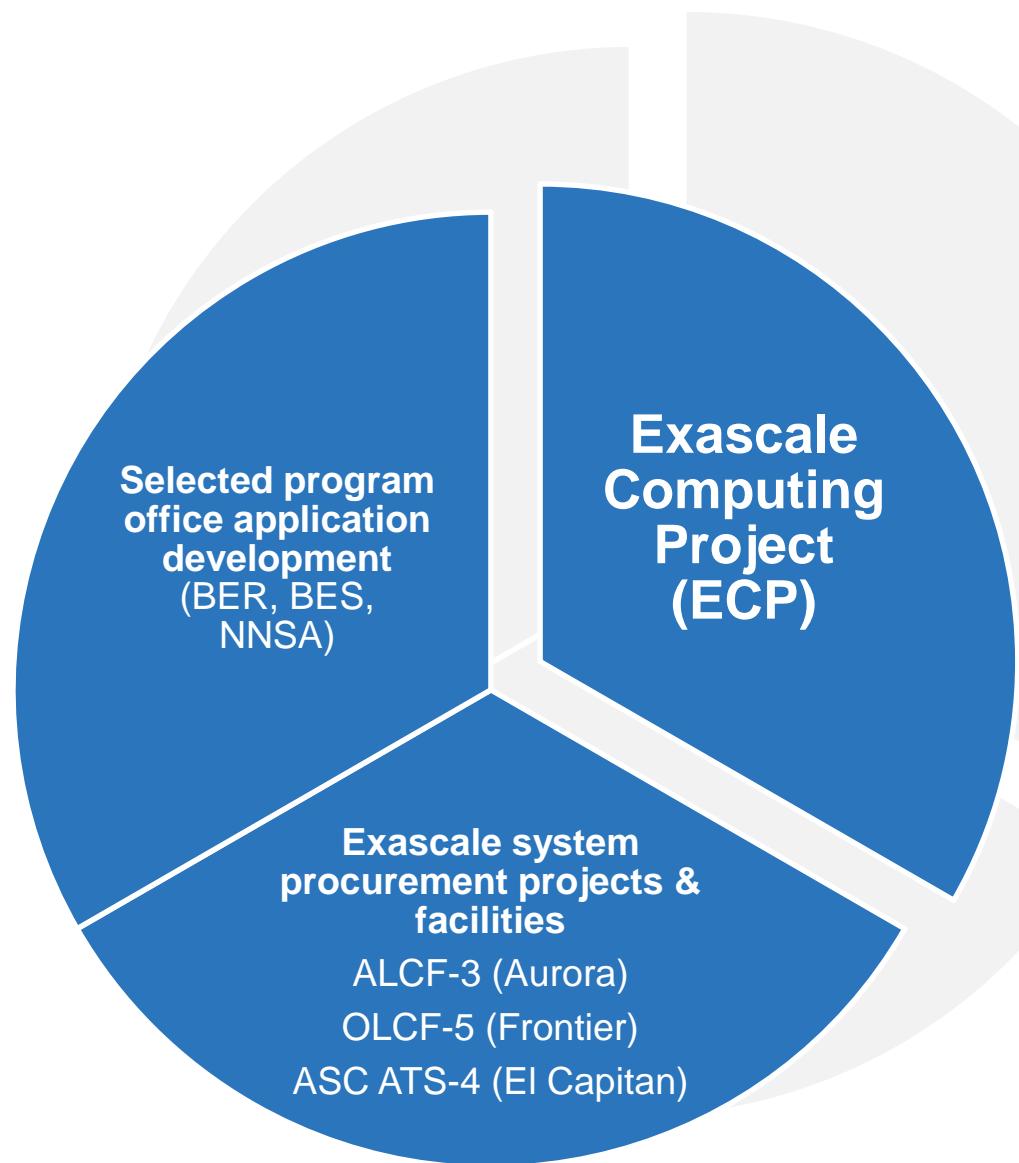
Scheduled for delivery in 2021, Frontier will accelerate innovation in science and technology and maintain U.S. leadership in high-performance computing and artificial intelligence. The total contract award is valued at more than \$600 million for the system and technology development. The system will be based on Cray's new Shasta architecture and Slingshot interconnect and will feature high-performance AMD EPYC CPU and AMD Radeon Instinct GPU technology.

Peak Performance	>1.5 EF
Footprint	> 100 cabinets
Node	1 HPC and AI Optimized AMD EPYC CPU 4 Purpose Built AMD Radeon Instinct GPU
CPU-GPU Interconnect	AMD Infinity Fabric Coherent memory across the node
System Interconnect	Multiple Slingshot NICs providing 100 GB/s network bandwidth Slingshot dragonfly network which provides adaptive routing, congestion management and quality of service.
Storage	2-4x performance and capacity of Summit's I/O subsystem. Frontier will have near node storage like Summit.



# US Exascale Computing Project

# DOE Exascale Program: The Exascale Computing Initiative (ECI)



Three Major Components of the ECI

# ECP by the Numbers

7  
YEARS  
\$1.7B

A seven-year, \$1.7 B R&D effort that launched in 2016

6  
CORE DOE  
LABS

Six core DOE National Laboratories: Argonne, Lawrence Berkeley, Lawrence Livermore, Los Alamos, Oak Ridge, Sandia

- Staff from most of the 17 DOE national laboratories take part in the project

3  
TECHNICAL  
FOCUS  
AREAS

Three technical focus areas (Application Development, Software Technology, Hardware and Integration) supported by project management expertise in the ECP Project Office

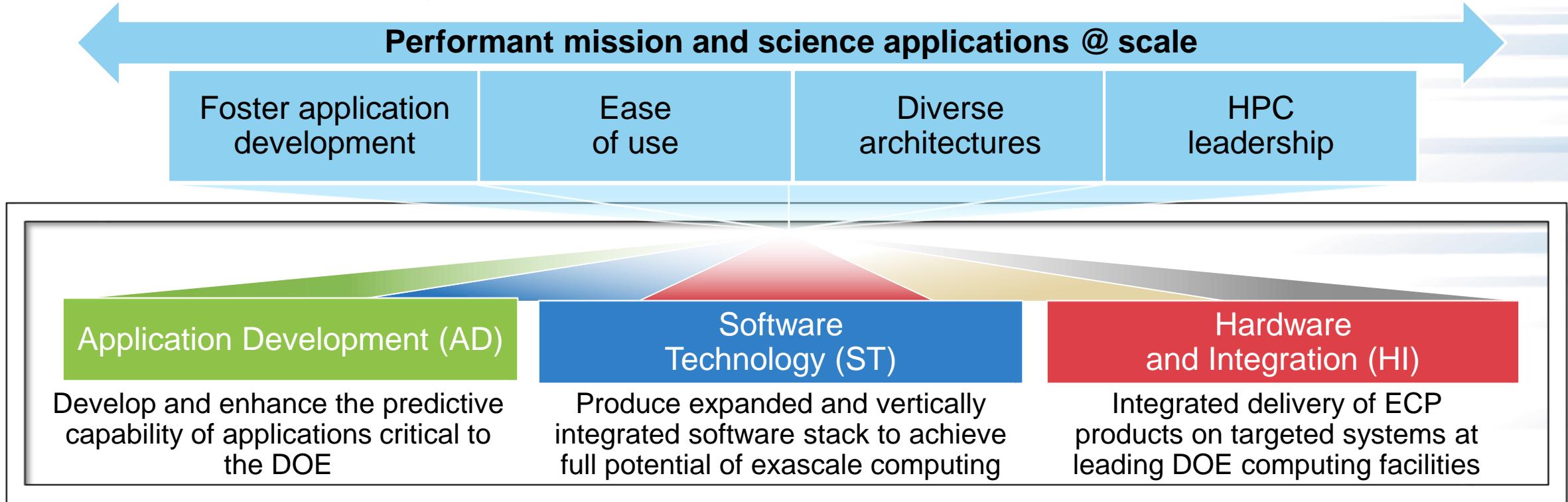
100  
R&D TEAMS  
1000  
RESEARCHERS

More than 100 top-notch R&D teams

ECP  
Project  
Office

Hundreds of consequential milestones delivered on schedule and within budget since project inception

# The three technical areas in ECP have the necessary components to meet national goals



25 applications ranging from national security, to energy, earth systems, economic security, materials, and data

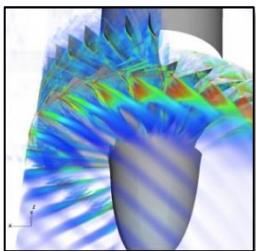
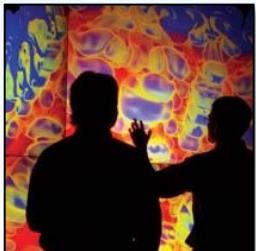
80+ unique software products spanning programming models and run times, math libraries, data and visualization

6 vendors supported by PathForward focused on memory, node, connectivity advancements; deployment to facilities

# ECP applications target national problems in 6 strategic areas

## National security

Stockpile stewardship  
Next-generation electromagnetics simulation of hostile environment and virtual flight testing for hypersonic re-entry vehicles

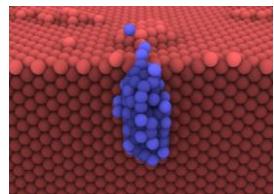


## Energy security

Turbine wind plant efficiency  
High-efficiency, low-emission combustion engine and gas turbine design  
Materials design for extreme environments of nuclear fission and fusion reactors  
Design and commercialization of Small Modular Reactors  
Subsurface use for carbon capture, petroleum extraction, waste disposal  
Scale-up of clean fossil fuel combustion  
Biofuel catalyst design

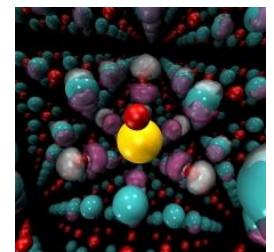
## Economic security

Additive manufacturing of qualifiable metal parts  
Reliable and efficient planning of the power grid  
Seismic hazard risk assessment  
Urban planning



## Scientific discovery

Find, predict, and control materials and properties  
Cosmological probe of the standard model of particle physics  
Validate fundamental laws of nature  
Demystify origin of chemical elements  
Light source-enabled analysis of protein and molecular structure and design  
Whole-device model of magnetically confined fusion plasmas



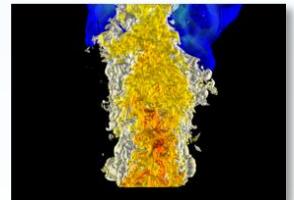
## Earth system

Accurate regional impact assessments in Earth system models  
Stress-resistant crop analysis and catalytic conversion of biomass-derived alcohols  
Metagenomics for analysis of biogeochemical cycles, climate change, environmental remediation



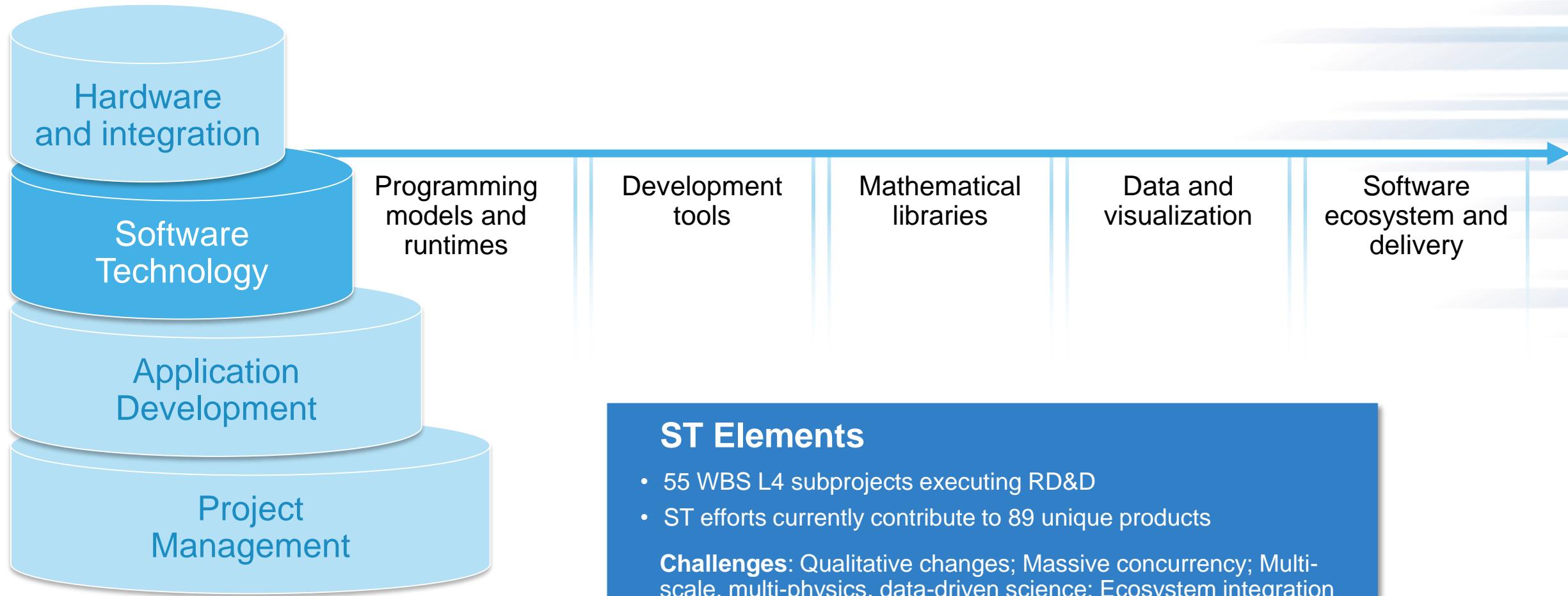
## Health care

Accelerate and translate cancer research



# Software Technology

Develop the exascale software stack and deliver using Software Development Kits (SDKs)



# Many ECP ST products are available (many github)

For example...

## Programming Models and Runtimes Products (16)

Legion  
ROSE  
Kokkos  
DARMA  
Global Arrays  
RAJA  
CHAI  
Umpire  
MPICH  
PaRSEC  
Open MPI  
Intel GEOPM  
LLVM OpenMP compiler  
OpenMP V&V Suite  
BOLT  
UPC++  
GASNet-EX  
Qthreads

<http://legion.stanford.edu>  
<https://github.com/rose-compiler>  
<https://github.com/kokkos>  
<https://github.com/darma-tasking>  
<http://hpc.pnl.gov/globalarrays/>  
<https://github.com/LLNL/RAJA>  
<https://github.com/LLNL/CHAI>

<http://www.xdk.info>  
<http://icl.utk.edu/exa-papi/>  
<http://www.hYPRE.org>  
<https://www.flecsi.org>  
<http://mfem.org/>  
[https://github.com/kokkos/kokkos-kernels/](https://github.com/kokkos/kokkos-kernels)  
<https://github.com/trilinos/Trilinos>  
<https://computation.llnl.gov/projects/sundials>  
<http://www.mcs.anl.gov/petsc>  
<https://github.com/Libensemble/libensemble>  
<http://portal.nersc.gov/project/sparse/strumpack/>  
<http://crd-legacy.lbl.gov/~xiaoye/SuperLU/>  
<https://trilinos.github.io/ForTrilinos/>  
<http://icl.utk.edu/slate/>  
<https://bitbucket.org/icl/magma>  
<https://github.com/ORNL-CEES/DataTransferKit>  
<http://tasmanian.ornl.gov/>

## Mathematical Libraries Products (16)

<https://xdk.info>  
<http://www.llnl.gov/casc/hypre>  
<http://www.flecsi.org>  
<http://mfem.org/>  
[https://github.com/kokkos/kokkos-kernels/](https://github.com/kokkos/kokkos-kernels)  
<https://github.com/trilinos/Trilinos>  
<https://computation.llnl.gov/projects/sundials>  
<http://www.mcs.anl.gov/petsc>  
<https://github.com/Libensemble/libensemble>  
<http://portal.nersc.gov/project/sparse/strumpack/>  
<http://crd-legacy.lbl.gov/~xiaoye/SuperLU/>  
<https://trilinos.github.io/ForTrilinos/>  
<http://icl.utk.edu/slate/>  
<https://bitbucket.org/icl/magma>  
<https://github.com/ORNL-CEES/DataTransferKit>  
<http://tasmanian.ornl.gov/>

etc...

## Development Tools (19)

SICM  
QUO  
Kitsune  
SCR  
Caliper  
mpiFileUtils  
Gotcha  
TriBITS  
Exascale Code Generation Toolkit  
PAPI  
CHILL Autotuning Compiler  
Search using Pandas  
<https://confluence.exascaleproject.org/display/STSS07>  
<https://github.com/lanl/libquo>  
<https://github.com/lanl/kitsune>  
<https://github.com/llnl/scr>  
<https://github.com/llnl/caliper>  
<https://github.com/llnl/mpifileutils>  
<https://github.com/llnl/gotcha>  
<https://tribits.org>  
<http://icl.utk.edu/exa-papi/>

org  
lyn.org  
[regon.edu/research/tau](http://regon.edu/research/tau)  
[/research/papyrus](http://research/papyrus)  
[/research/openarc](http://research/openarc)  
[regon.edu/research/pdt/home.php](http://regon.edu/research/pdt/home.php)



# Software Development Kits (SDKs): A Key ST Design Feature

An important delivery vehicle for software products with a direct line of sight to ECP applications

## ECP software projects

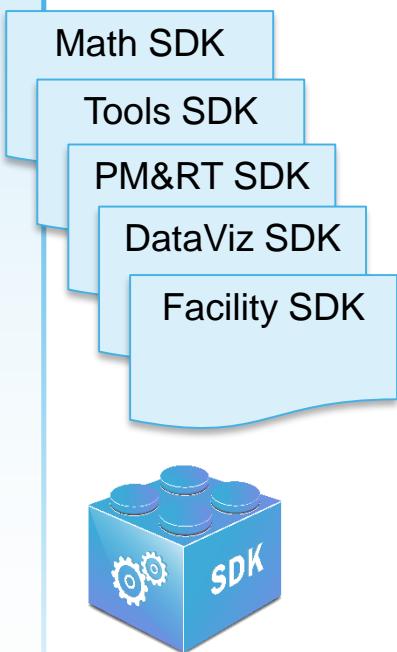
Each project to define (at least 2) release vectors

More projects

### SDKs

Reusable software libraries embedded in applications; cohesive/interdependent libraries released as sets modeled on xSDK

- Regular coordinated releases
- Hierarchical collection built on Spack
- Products may belong to >1 SDK based on dependences
- Establish community policies for library development
- Apply Continuous Integration and other robust testing practices



Fewer projects

### OpenHPC

Potential exit strategy for binary distributions

- Target similar software to existing OpenHPC stack
- Develop super-scalable release targeting higher end systems

Assume all releases are delivered as “build from source” via Spack – at least initially

Focus on ensuring that software compiles robustly on all platforms of interest to ECP (including testbeds)

### Direct2Facility

Platform-specific software in support of a specified 2021–2023 exascale system

- Software **exclusively** supporting a specific platform
- System software, some tools and runtimes



# Time for a short poll...

**Q: Think back 10 years. How many of you would have predicted that most of our top HPC systems would be GPU-based architectures?**

- a) Yes
- b) No
- c) Waffle ☺

# **Q: Think forward 10 years. How many of you predict that most of our top HPC systems will have the following architectural features?**

- a) X86 multicore CPU
- b) GPU
- c) FPGA/Reconfigurable processor
- d) Neuromorphic processor
- e) Deep learning processor
- f) Quantum processor
- g) Some new unknown processor
- h) All/some of the above in one SoC

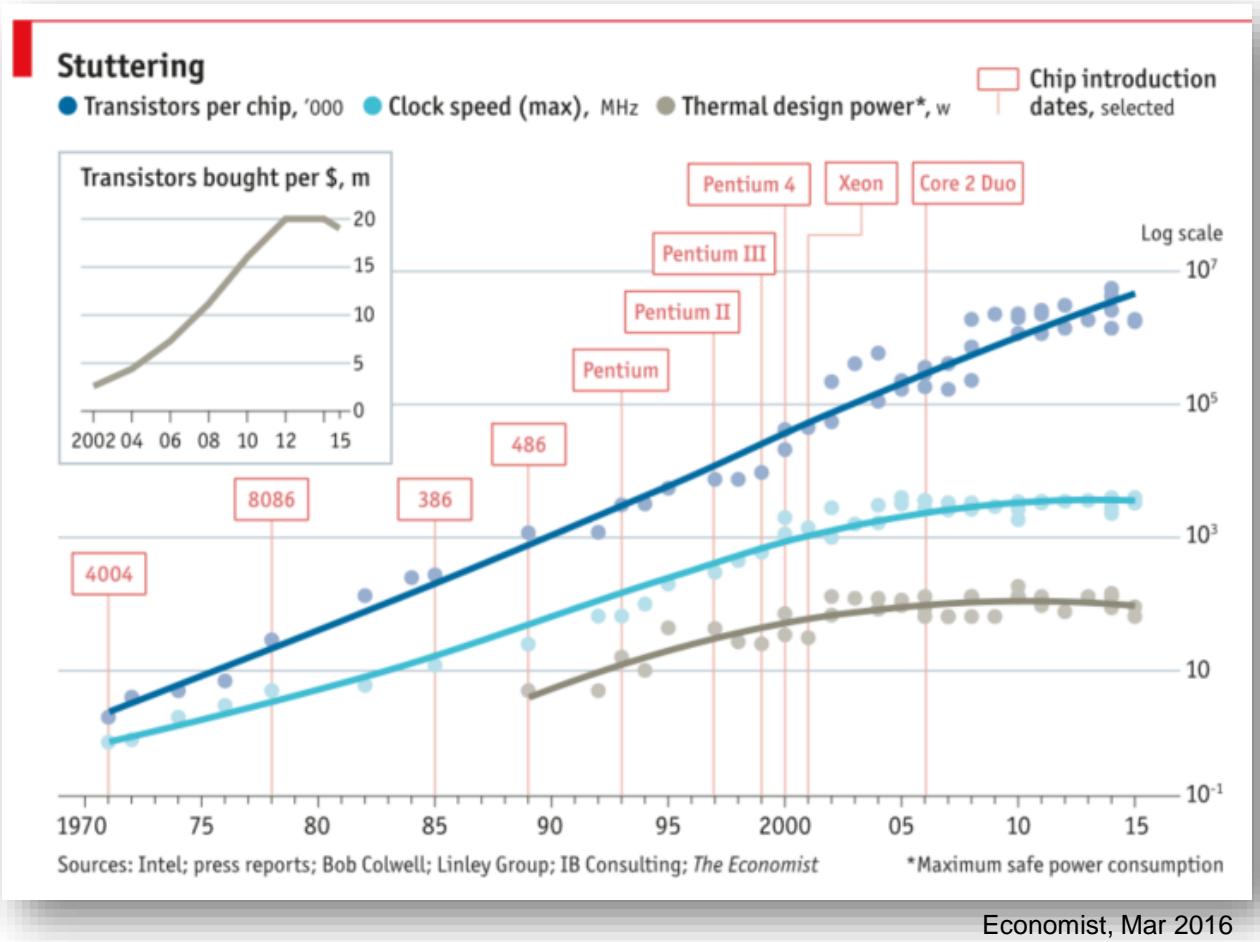
**Q: Now imagine you are building a new application with ~3M LOC and 20 team members over 10 years.**

**What on-node programming model/system do you use?**

- a) C, C++, Fortran
- b) C++ templates, policies, etc (e.g., AMP, Kokkos)
- c) CUDA, cu\*\*\*, HIP
- d) OpenCL, SYCL
- e) OpenMP or OpenACC
- f) A Domain Specific Language (e.g., Claw)
- g) A Domain Specific Framework (e.g., PetSc)
- h) Some new unknown programming approach
- i) All/some of the above

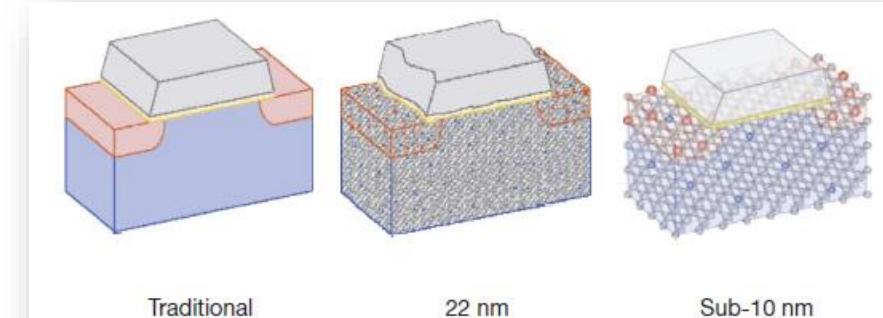
# Motivating Trends

# Contemporary devices are approaching fundamental limits

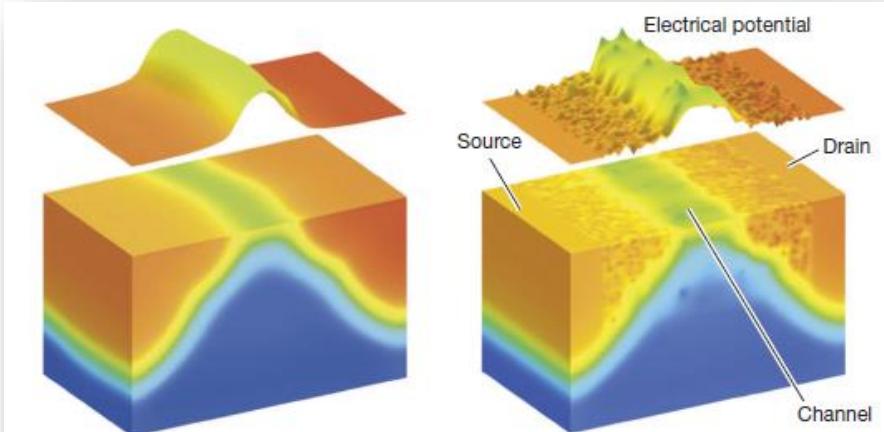


Dennard scaling has already ended. Dennard observed that voltage and current should be proportional to the linear dimensions of a transistor: 2x transistor count implies 40% faster and 50% more efficient.

R.H. Dennard, F.H. Gaenslen, V.L. Rideout, E. Bassous, and A.R. LeBlanc, "Design of ion-implanted MOSFET's with very small physical dimensions," *IEEE Journal of Solid-State Circuits*, 9(5):256-68, 1974.



**Figure 1 |** As a metal oxide–semiconductor field effect transistor (MOSFET) shrinks, the gate dielectric (yellow) thickness approaches several atoms (0.5 nm at the 22-nm technology node). Atomic spacing limits the



**Figure 2 |** As a MOSFET transistor shrinks, the shape of its electric field departs from basic rectilinear models, and the level curves become disconnected. Atomic-level manufacturing variations, especially for dopant

# Business climate reflects this uncertainty, cost, complexity, consolidation

designlines WIRELESS & NETWORKING

Blog

## IC Merger Mania Hits Fever Pitch

Dylan McGrath, Contributing Editor

12/2/2015 10:13 AM EST

1 comments post a comm

[Like](#) 10 [Tweet](#) 10

With the announcement  
PMC-Sierra, the total val  
acquisitions announced

The wave of consolidatio

SEMICONDUCTOR ENGINEERING

Home > Manufacturing, Design & Test > Uncertainty Grows For 5nm, 3nm

MANUFACTURING, DESIGN & TEST

## Uncertainty Grows For

5nm, 3nm

1,787 C 74

Nanosheets and nanowire FETs us  
costs are skyrocketing. New packa  
provide an alternative.

DECEMBER 19TH, 2016 - BY: MARK LAPEDUS

As several chipmakers ramp up their

designlines SoC

News & Analysis

## TSMC Grows Share of Foundry

Business

Repercussions of Samsung's b

Alan Patterson

10/13/2016 09:38 AM EDT

Post a comment

[Tweet](#) [Share](#) 20 [G+](#)

TAIPEI — Taiwan Semiconductor Manufac

increased its share of the foundry business  
on better than expected demand for smart  
quarter.

"In the third quarter, we gained market sha

Globalfoundries s

the bleeding edge

"In the third quarter, we gained market sha

## Intel to acquire Altera for \$54 a share

Monday, 1 Jun 2015 | 8:33



## Avago Agrees to Buy Broadcom for \$37 Billion

By MICHAEL J. de la MERCED and CHAD BRAY MAY 28, 2015



designlines AUTOMOTIVE

News & Analysis

## Foundries' Sales Show Hard Times Continuing

Peter Clarke

2016 09:33 PM EDT  
Comments

[Like](#) 6

[Tweet](#)

## Tech giant ARM Holdings sold to Japanese firm for £24bn

#BUSINESS NEWS NOVEMBER 19, 2017 / 7:57 PM / UPDATED 21 MINUTES AGO

EXCLUSIVE

## Marvell Technology to buy rival chipmaker Cavium for \$6 billion

By EXCLUSIVE

## SoftBank to sell 25% of Arm to Saudi-backed fund

Son puts stake worth \$8bn in UK's largest tech company into \$100bn Vision Fund



## Amazon Is Becoming an AI Chip Maker, Speeding Alexa Responses

By Aaron Tilley Feb. 12, 2018 7:00 AM PST • Comments by Yonatan Raz-Fridman and Mohammad Musa

nytimes.com

## Hewlett Packard Enterprise to Acquire Supercomputer Pioneer Cray

6-7 minutes

5-6 minutes

## NVIDIA Buys Mellanox To Bring HPC Scaling To Data Centers

to Acquire Supercomputer

\$1.4 billion to acquire

## SANDISK COMPLETES ACQUISITION OF FUSION IO

JUL 23, 2014

ACQUISITION TO BOOST SANDISK'S ENTERPRISE GROWTH

MILPITAS, Calif., July 23, 2014 - SanDisk Corporation (NASDAQ: SNDK), a global leader in flash storage solutions, today

announce

hardware

"I am deli

the Fusio

and chief

the indus

## Western Digital Now A Storage Powerhouse With SanDisk Acquisition

to-market talent of

ehrotra, preside

se flash solutions in



## Broadcom acquires Brocade in \$5.9 billion deal

Posted 1 hour ago by Ron Miller (@ron\_miller)

[Like](#) [Share](#) [Email](#) [Print](#)



## Toshiba to sell 'minority stake' in chip business to Western Digital

In 2016, Toshiba had a 20.4% share in global NAND flash memory market.

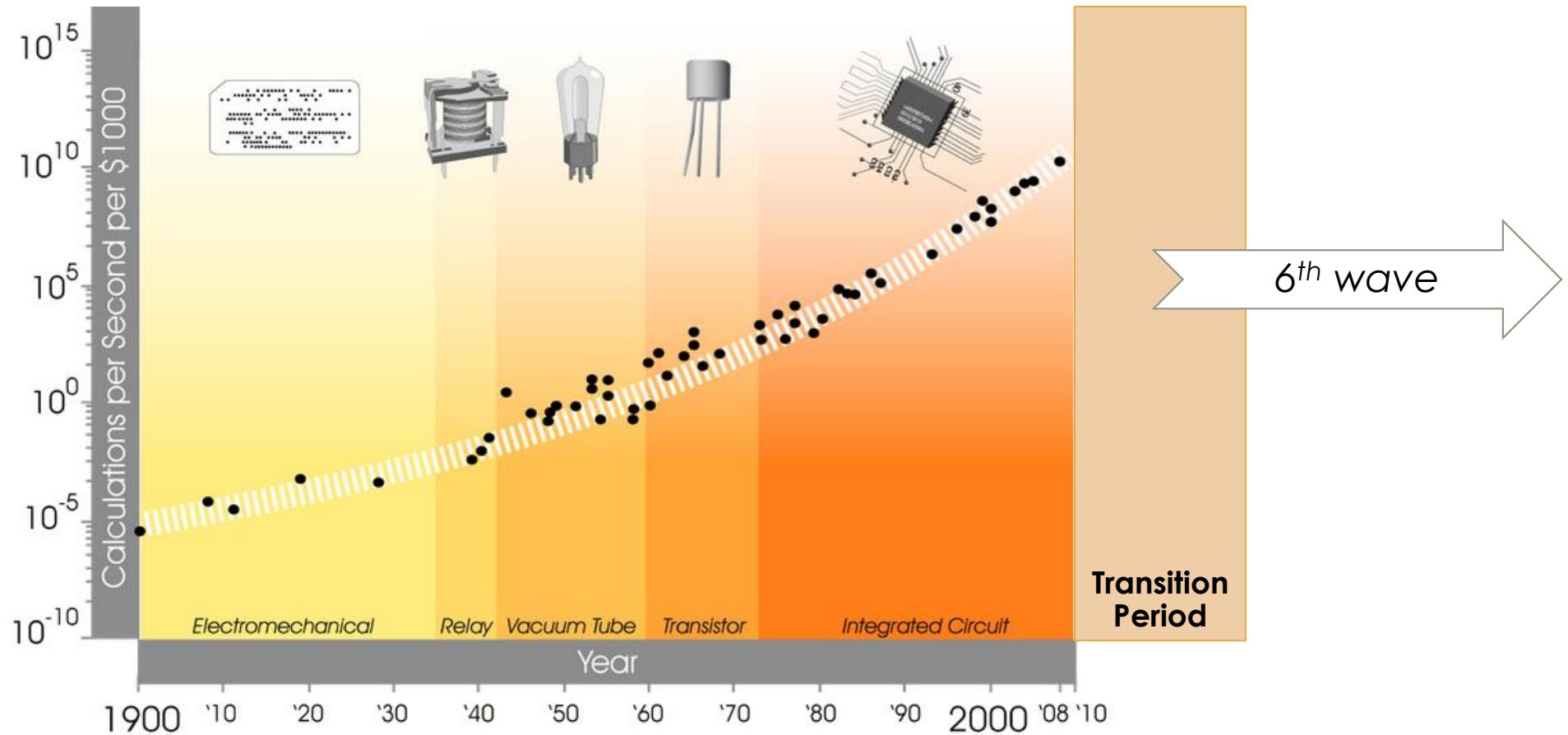
In Intel's Arduous Journey to 10 nm, Moore's Law Comes Up Short

Dairis Latimer, Technical Advisor, Red Oak Consulting | August 30, 2018 11:53 CEST

With a share price riding high and dominance in the datacentre market, it may seem perverse to state that Intel is a company facing a range of significant problems. So what caused the technology behemoth on the occasion of its 50th birthday to act so spectacularly on its back foot?

Andy Grove's famous maxim, "Success breeds complacency. Complacency breeds failure." Only the paranoid survive has proven accurate once again. Intel again finds itself at a classic Grove strategic inflection point. The problem is

# Sixth Wave of Computing



<http://www.kurzweilai.net/exponential-growth-of-computing>

# Predictions for Transition Period

## Optimize Software and Expose New Hierarchical Parallelism

- Redesign software to boost performance on upcoming architectures
- Exploit new levels of parallelism and efficient data movement

## Architectural Specialization and Integration

- Use CMOS more efficiently for our workloads
- Integrate components to boost performance and eliminate inefficiencies

## Emerging Technologies

- Investigate new computational paradigms
  - Quantum
  - Neuromorphic
  - Advanced Digital
  - Emerging Memory Devices

# Predictions for Transition Period

## Optimize Software and Expose New Hierarchical Parallelism

- Redesign software to boost performance on upcoming architectures
- Exploit new levels of parallelism and efficient data movement

## Architectural Specialization and Integration

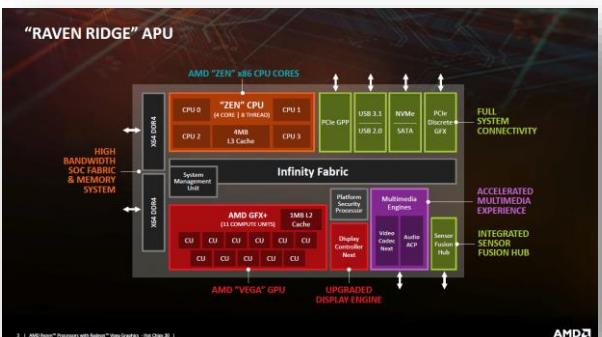
- Use CMOS more efficiently for our workloads
- Integrate components to boost performance and eliminate inefficiencies

## Emerging Technologies

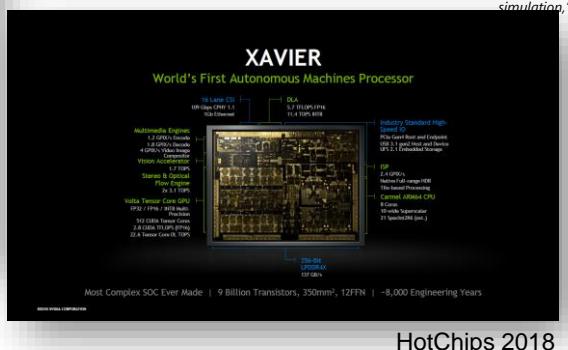
- Investigate new computational paradigms
  - Quantum
  - Neuromorphic
  - Advanced Digital
  - Emerging Memory Devices

# Pace of Architectural Specialization is Quickening

- Industry, lacking Moore's Law, will need to continue to differentiate products (to stay in business)
- Grant that advantage of better CMOS process stalls
- Use the same transistors differently to enhance performance
- Architectural design will become extremely important, critical
  - Dark Silicon
  - Address new parameters for benefits/curse of Moore's Law



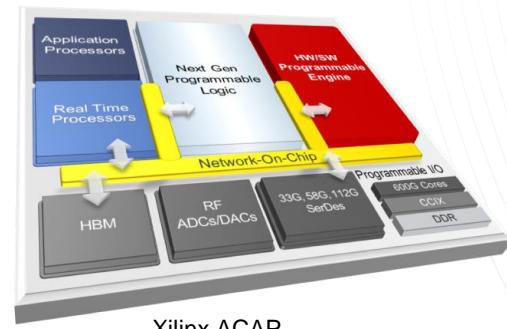
HotChips 2018



HotChips 2018



D.E. Shaw, M.M. Deneroff, R.O. Dror et al., "Anton, a special-purpose machine for molecular dynamics simulation," *Communications of the ACM*, 51(7):91-7, 2008.



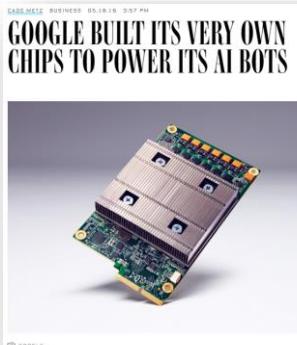
<https://fossbytes.com/nvidia-volta-gddr6-2018/>



<https://www.broadcastbridge.com/content/entry/1094/altera-announces-new-arria-10-fpgas-and-socs>

TOM SIMONITE BUSINESS 11.27.18 08:12 PM

**NEW AT AMAZON: ITS OWN CHIPS FOR CLOUD COMPUTING**



GOOGLE BUILT ITS VERY OWN CHIPS TO POWER ITS AI BOTS



GOOGLE HAS DESIGNED its own computer chip for driving deep neural networks, an AI technology that is reinventing the way Internet services operate.

This morning at Google I/O, the centerpiece of the company's year, CEO Sundar Pichai said that Google has designed an ASIC, or application-specific integrated circuit, that's specific to deep neural nets. These are networks of

<http://www.wired.com/2016/05/google-tpu-custom-chips/>

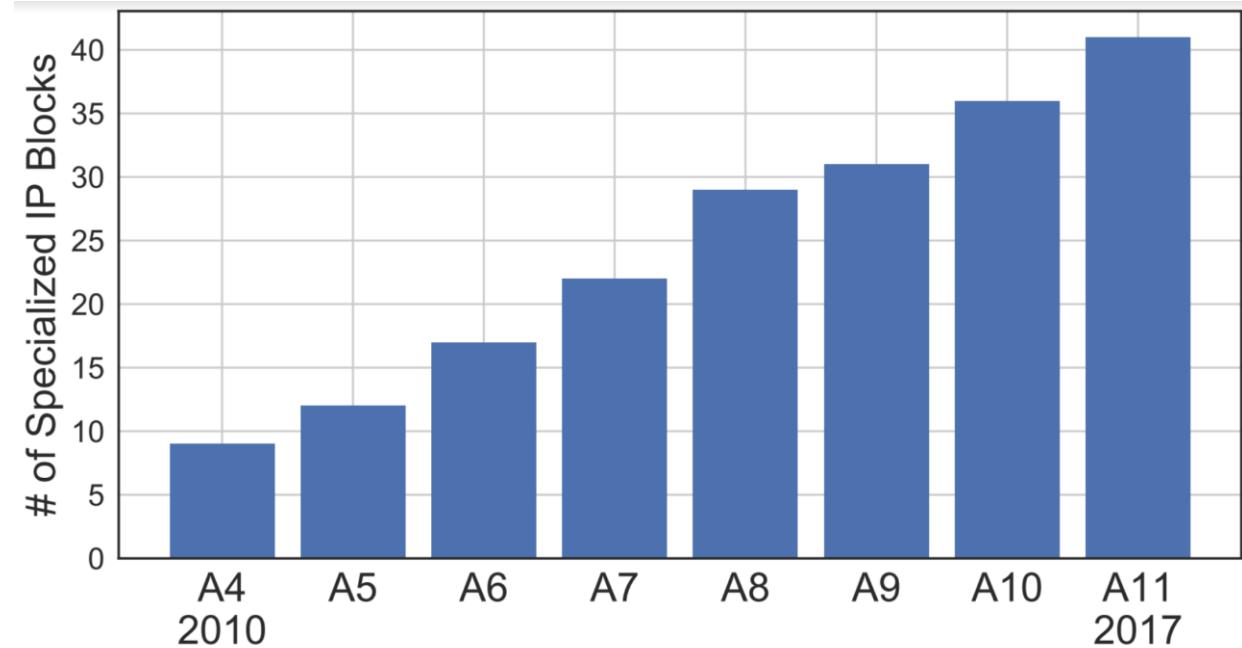
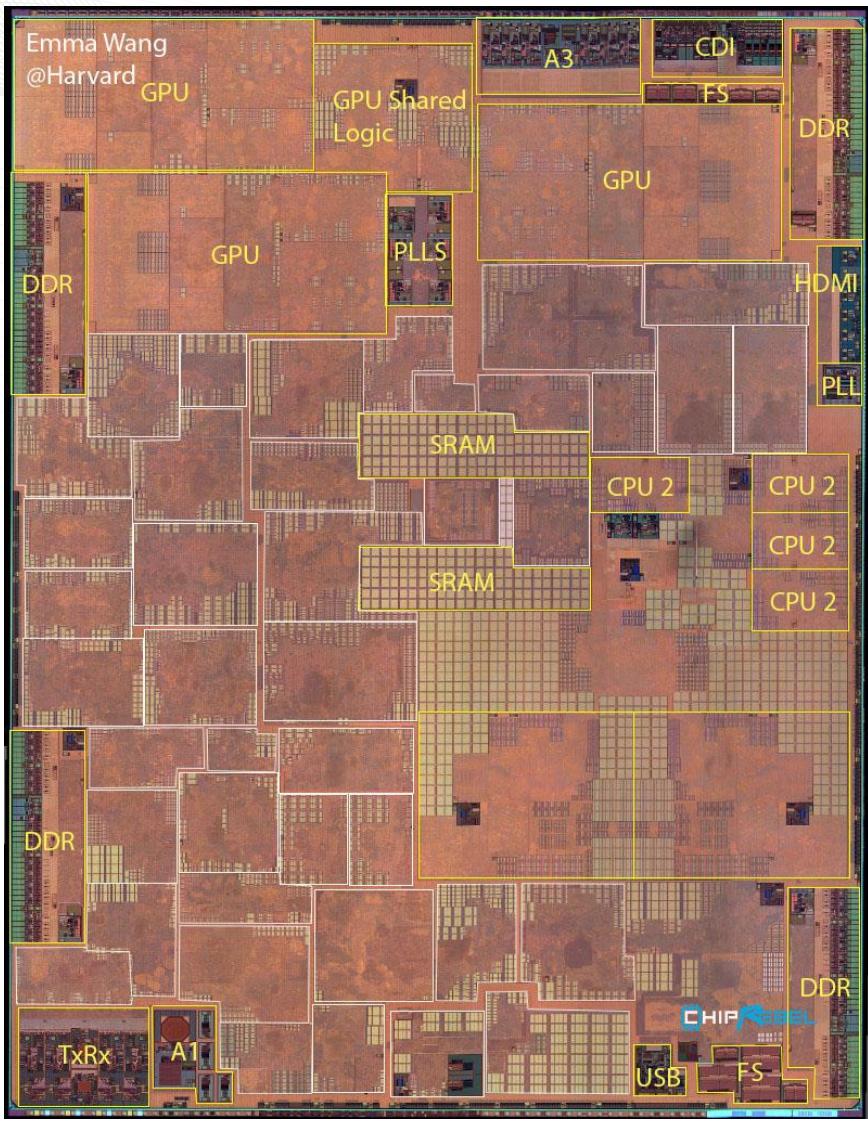


Amazon Web Services CEO Andy Jassy speaks at an event in San Francisco in 2017.

DAVID PAUL MORRIS/BLOOMBERG/GETTY IMAGES

**BIG SOFTWARE COMPANIES** don't just stick to software any more—they build computer chips. The latest proof comes from Amazon, which announced late Monday that its cloud computing division has created its own chips to power customers' websites and other services. The chips, dubbed Graviton, are built around the same technology that powers smartphones and tablets. That approach has been much discussed in the cloud industry but never

# Analysis of Apple A-\* SoCs



# DARPA ERI Programs Aiming for Agile (and Frequent) Chip Creation



## IDEA/POSH End State – A Universal Hardware Compiler

```
$ git clone https://github.com/darpa/idea  
$ git clone https://github.com/darpa/posh  
$ cd posh  
$ make soc42
```



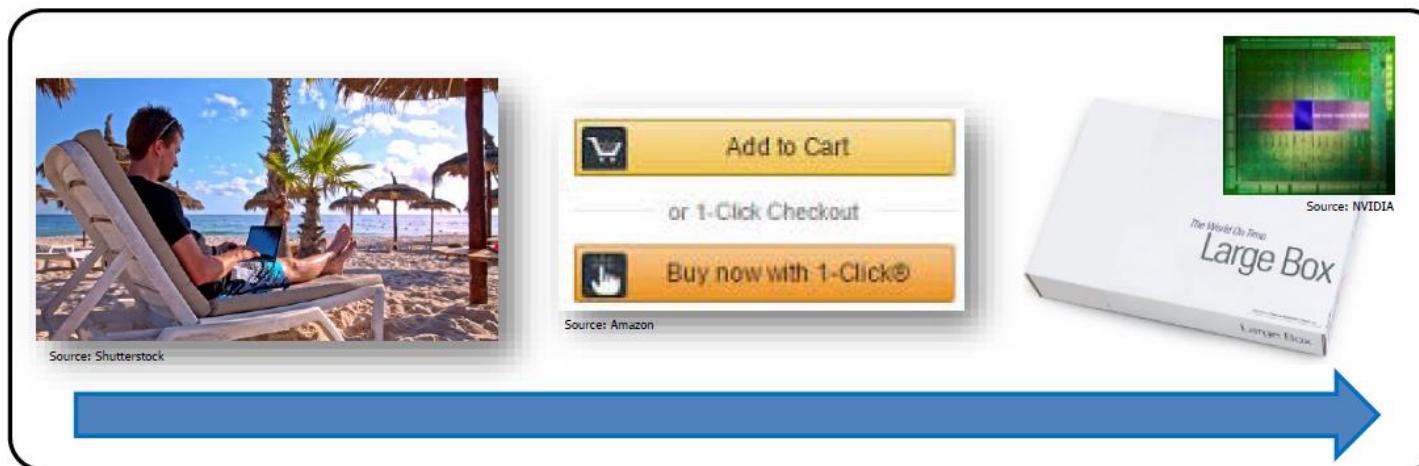
Source: Shutterstock



Source: Amazon



Source: NVIDIA



Distribution Statement "A" (Approved for Public Release, Distribution Unlimited)

23

A. Olofsson, 2018

# Growing Open Source Hardware Movement Enables Rapid Chip Design



## RISC-V Ecosystem

### Software

**Open-source software:**  
Gcc, binutils, glibc, Linux, BSD,  
LLVM, QEMU, FreeRTOS,  
ZephyrOS, LiteOS, SylixOS, ...

**Commercial software:**  
Lauterbach, Segger, Micrium,  
ExpressLogic, ...



ISA specification

Golden Model

Compliance

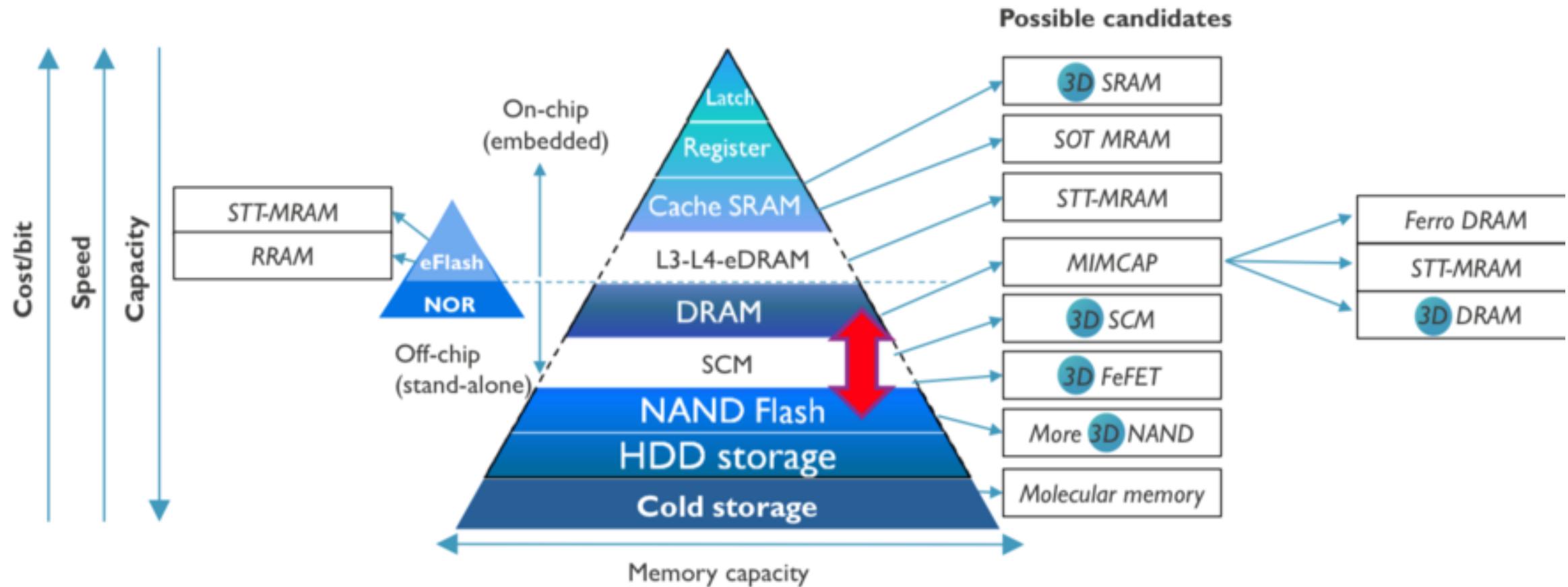
### Hardware

**Open-source cores:**  
Rocket, BOOM, RI5CY,  
Ariane, PicoRV32, Piccolo,  
SCR1, Hummingbird, ...

**Commercial core providers:**  
Andes, Bluespec, Cloudbear,  
Codasip, Cortus, C-Sky,  
Nuclei, SiFive, Syntacore, ...

**Inhouse cores:**  
Nvidia, +others

# Memory Hierarchy is Specializing too



# Transition Period will be Disruptive

- New devices and architectures may not be hidden in traditional levels of abstraction
  - A new type of CNT transistor may be completely hidden from higher levels
  - A new paradigm like quantum may require new architectures, programming models, and algorithmic approaches
- Solutions need a co-design framework to evaluate and mature specific technologies

Layer	Switch, 3D	NVM	Approximate	Neuro	Quantum
<i>Application</i>	1	1	2	2	3
<i>Algorithm</i>	1	1	2	3	3
<i>Language</i>	1	2	2	3	3
<i>API</i>	1	2	2	3	3
<i>Arch</i>	1	2	2	3	3
<i>ISA</i>	1	2	2	3	3
<i>Microarch</i>	2	3	2	3	3
<i>FU</i>	2	3	2	3	3
<i>Logic</i>	3	3	2	3	3
<i>Device</i>	3	3	2	3	3

Adapted from IEEE Rebooting Computing Chart

# Department of Energy (DOE) Roadmap to Exascale Systems

An impressive, productive lineup of *accelerated node* systems supporting DOE's mission

## Pre-Exascale Systems [Aggregate Linpack (Rmax) = 323 PF!]

2012



**Titan (9)**  
ORNL  
Cray/AMD/NVIDIA



**Mira (21)**  
ANL  
IBM BG/Q



**Sequoia (10)**  
LLNL  
IBM BG/Q

2016

### Heterogeneous Cores



**Summit (1)**  
ORNL  
IBM/NVIDIA

### Theta (24)

**ANL**  
Cray/Intel KNL



**Cori (12)**  
LBNL  
Cray/Intel Xeon/KNL

### Trinity (6)

**LANL/SNL**  
Cray/Intel Xeon/KNL

2018



**Sierra (2)**  
LLNL  
IBM/NVIDIA



**CROSSROADS**  
LANL/SNL  
TBD

2020



**Perlmutter**  
LBNL  
Cray/AMD/NVIDIA

## First U.S. Exascale Systems

2021-2023



**FRONTIER**  
ORNL  
AMD/Cray



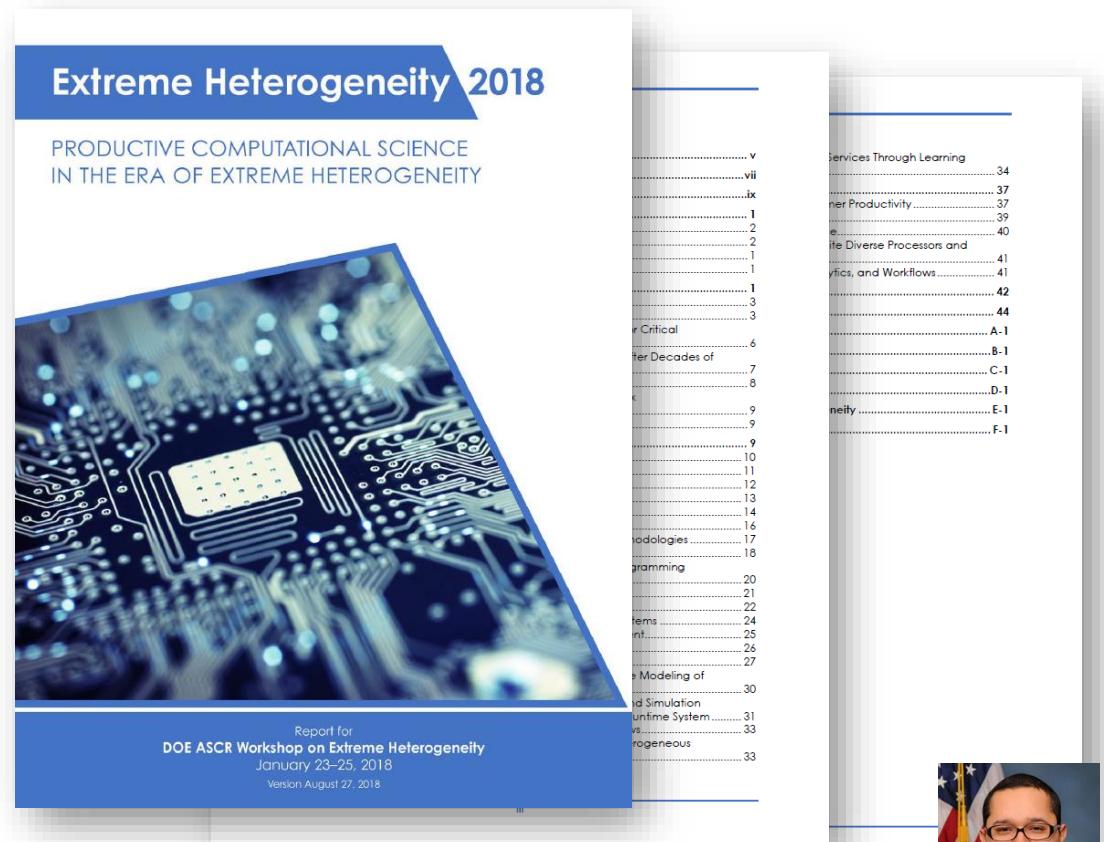
**Aurora**  
ANL  
Intel/Cray



**EL CAPITAN**  
LLNL  
TBD

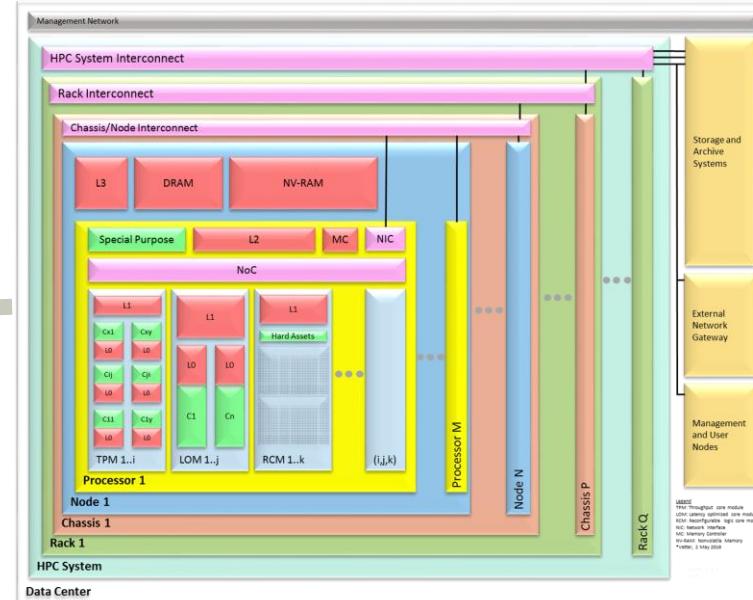
# Final Report on Workshop on Extreme Heterogeneity

1. Maintaining and improving programmer productivity
  - Flexible, expressive, programming models and languages
  - Intelligent, domain-aware compilers and tools
  - Composition of disparate software components
- Managing resources intelligently
  - Automated methods using introspection and machine learning
  - Optimize for performance, energy efficiency, and availability
- Modeling & predicting performance
  - Evaluate impact of potential system designs and application mappings
  - Model-automated optimization of applications
- Enabling reproducible science despite non-determinism & asynchrony
  - Methods for validation on non-deterministic architectures
  - Detection and mitigation of pervasive faults and errors
- Facilitating Data Management, Analytics, and Workflows
  - Mapping of science workflows to heterogeneous hardware and software services
  - Adapting workflows and services to meet facility-level objectives through learning approaches



# Programming Heterogeneous Systems

# Complex Architectures Yields Complex Programming Models



- This approach is not scalable, affordable, robust, elegant, etc.
- Not performance portable across different architectures

**System:** MPI, Legion, HPX, Charm++, etc

Low overhead

Resource contention

Locality

**Node:** OpenMP, Pthreads, U-threads, etc

SIMD

NUMA, HBM

**Cores:** OpenACC, CUDA, OpenCL, OpenMP4, SYCL, Kokkos...

Memory use,  
coalescing

Data  
orchestration

Fine grained  
parallelism

Hardware  
features

# **Directive-based Solutions for FPGA Computing**

# Challenges in FPGA Computing

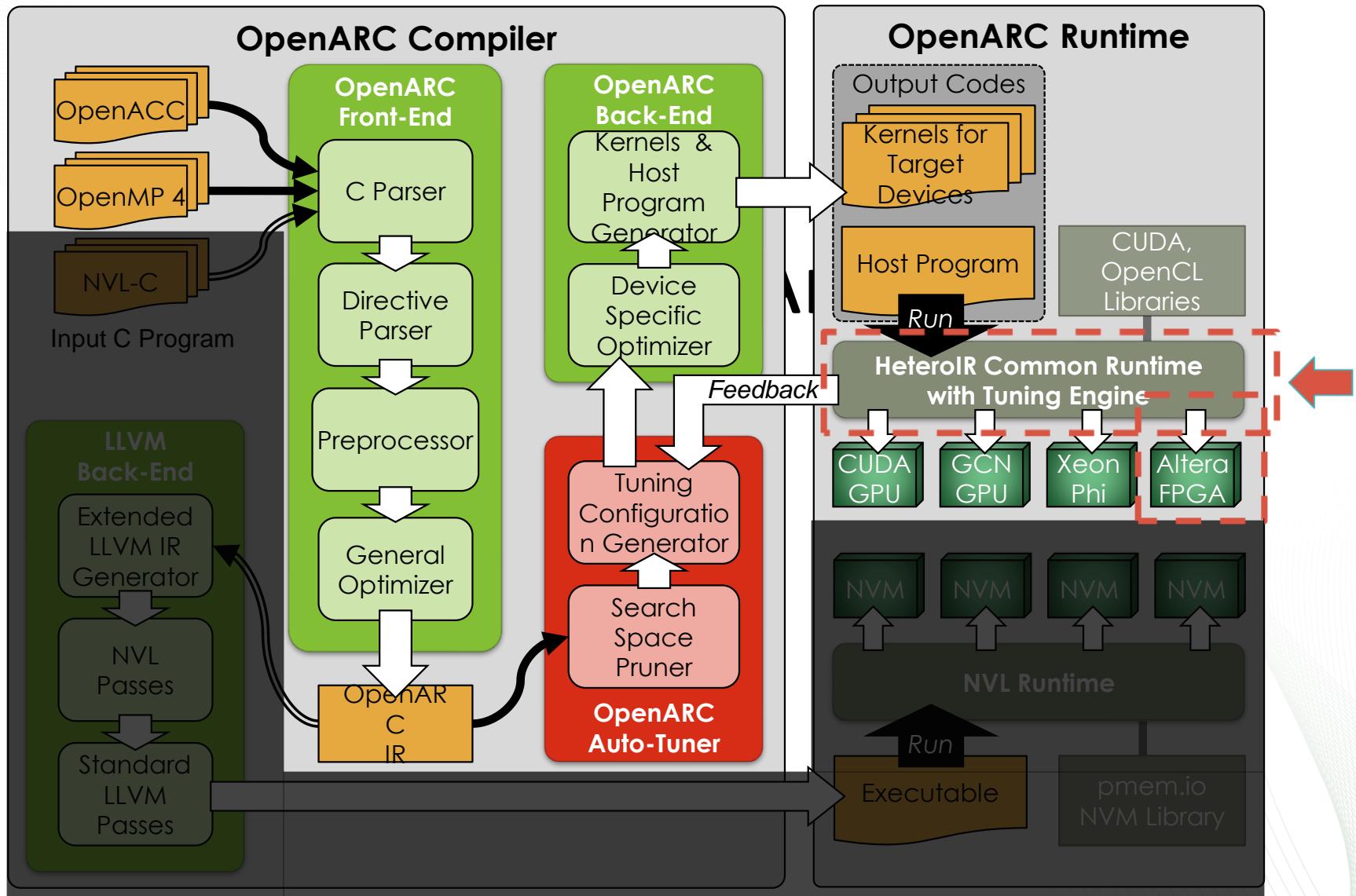
- Programmability and Portability Issues
  - Best performance for FPGAs requires writing Hardware Description Languages (HDLs) such as VHDL and Verilog; too complex and low-level
    - HDL requires substantial knowledge on hardware (digital circuits).
    - Programmers must think in terms of a state machine.
    - HDL programming is a kind of digital circuit design.
  - High-Level Synthesis (HLS) to provide better FPGA programmability
    - SRC platforms, Handel-C, Impulse C-to-FPGA compiler, Xilinx Vivado (AutoPilot), FCUDA, etc.
    - None of these use a portable, open standard.

# Standard, Portable Programming Models for Heterogeneous Computing

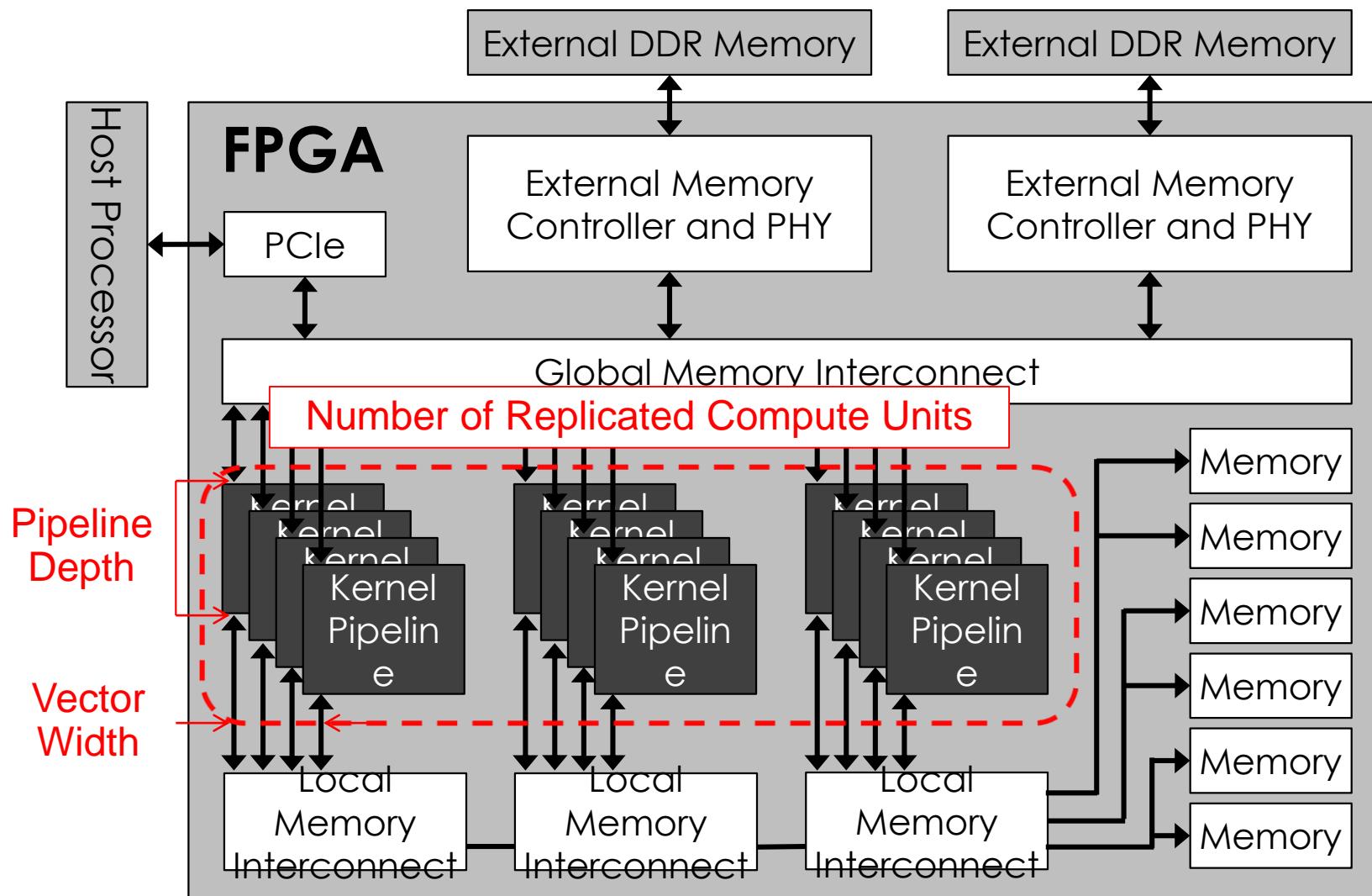
- OpenCL
  - Open standard portable across diverse heterogeneous platforms (e.g., CPUs, GPUs, DSPs, Xeon Phis, FPGAs, etc.)
  - Much higher than HDL, but still complex for typical programmers.
- Directive-based accelerator programming models
  - OpenACC, OpenMP4, etc.
  - Provide higher abstraction than OpenCL.
  - Most of existing OpenACC/OpenMP4 compilers target only specific architectures; none supports FPGAs.

# FPGAs | Approach

- Design and implement an OpenACC-to-FPGA translation framework, which is the first work to use a standard and portable directive-based, high-level programming system for FPGAs.
- Propose FPGA-specific optimizations and novel pragma extensions to improve performance.
- Evaluate the functional and performance portability of the framework across diverse architectures (Altera FPGA, NVIDIA GPU, AMD GPU, and Intel Xeon Phi).

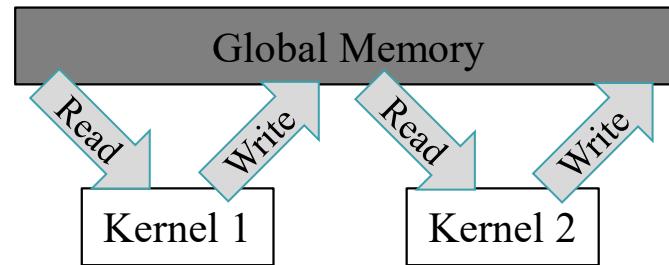


# FPGA OpenCL Architecture

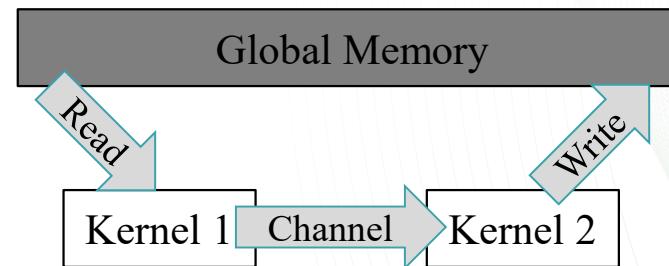


# Kernel-Pipelining Transformation Optimization

- Kernel execution model in OpenACC
  - Device kernels can communicate with each other only through the device global memory.
  - Synchronizations between kernels are at the granularity of a kernel execution.
- Altera OpenCL channels
  - Allows passing data between kernels and synchronizing kernels with high efficiency and low latency



Kernel communications through global memory in OpenACC

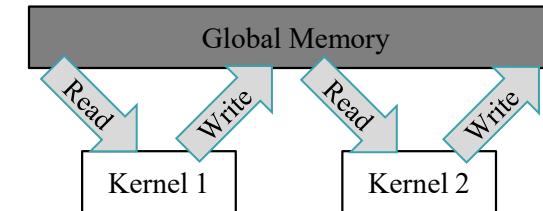


Kernel communications with Altera channels

# Kernel-Pipelining Transformation Optimization (2)

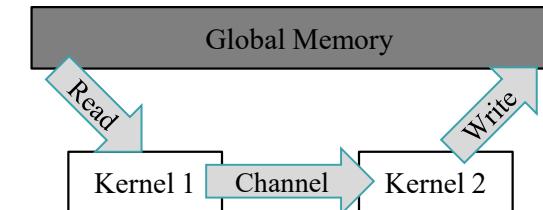
(a) Input OpenACC code

```
#pragma acc data copyin (a) create (b) copyout (c)
{
    #pragma acc kernels loop gang worker present (a, b)
    for(i=0; i<N; i++) { b[i] = a[i]*a[i]; }
    #pragma acc kernels loop gang worker present (b, c)
    for(i=0; i<N; i++) {c[i] = b[i]; }
}
```



(b) Altera OpenCL code with channels

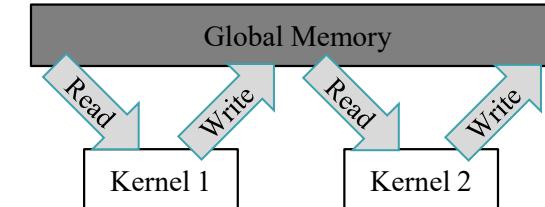
```
channel float pipe_b;
__kernel void kernel1(__global float* a) {
    int i = get_global_id(0);
    write_channel_altera(pipe_b, a[i]*a[i]);
}
__kernel void kernel2(__global float* c) {
    int i = get_global_id(0);
    c[i] = read_channel_altera(pipe_b);
}
```



# Kernel-Pipelining Transformation Optimization (3)

(a) Input OpenACC code

```
#pragma acc data copyin (a) create (b) copyout (c)
{
    #pragma acc kernels loop gang worker present (a, b)
    for(i=0; i<N; i++) { b[i] = a[i]*a[i]; }
    #pragma acc kernels loop gang worker present (b, c)
    for(i=0; i<N; i++) {c[i] = b[i]; }
}
```



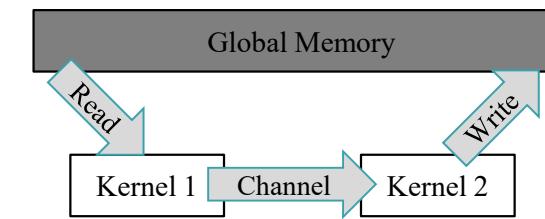
Kernel-pipelining transformation

Valid under specific conditions



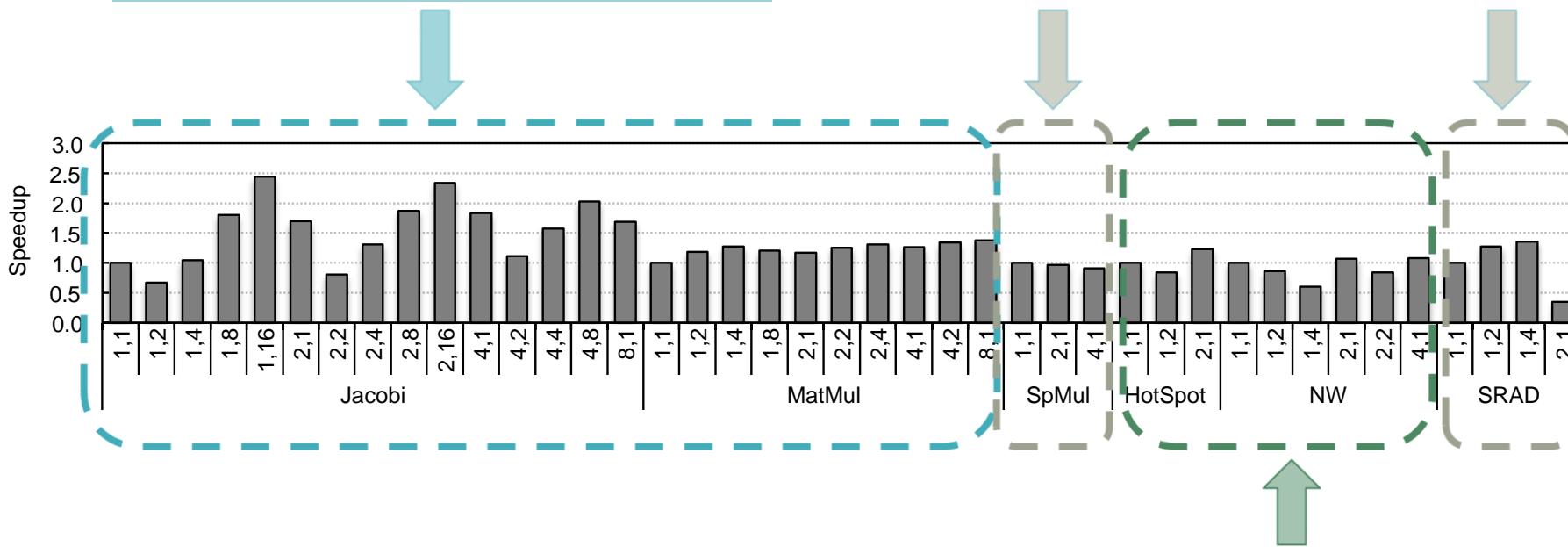
(c) Modified OpenACC code for kernel-pipelining

```
#pragma acc data copyin (a) pipe (b) copyout (c)
{
    #pragma acc kernels loop gang worker pipeout (b) present (a)
    For(i=0; i<N; i++) { b[i] = a[i]*a[i]; }
    #pragma acc kernels loop gang worker pipein (b) present (c)
    For(i=0; i<N; i++) {c[i] = b[i]; }
}
```



# Speedup over CU, SIMD (1,1)

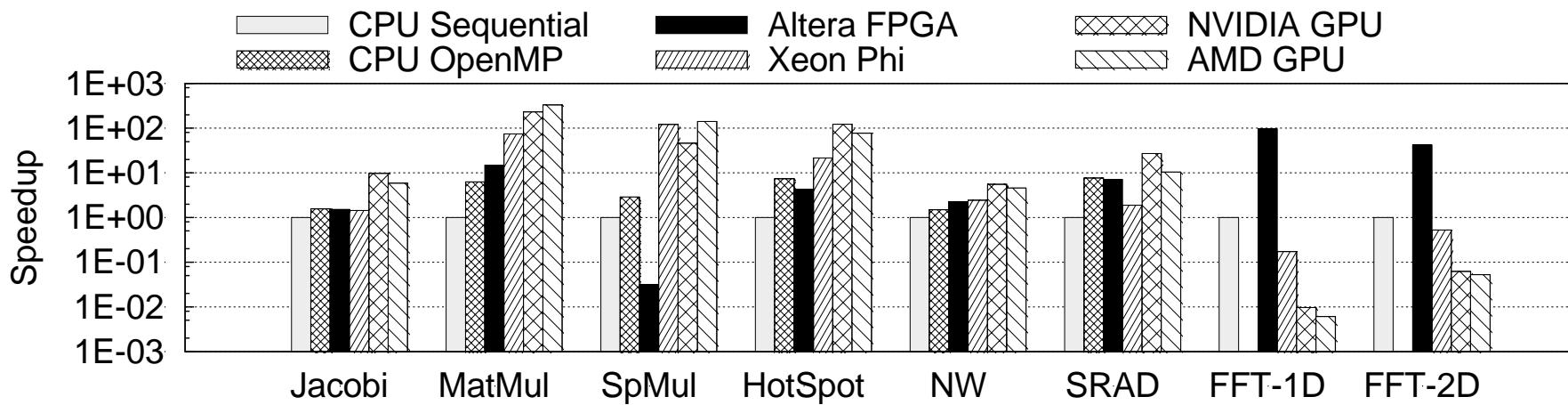
Jacobi and MatMul show better performance with increase in CU and SIMD, thanks to regular memory accesses.



SpMul and SRAD perform worse with multiple CUs, mainly due to memory contention.

Performance of HotSpot and NW increases with multiple CUs, but decreases with vectorization.

# Overall Performance



FPGAs prefer applications with deep execution pipelines (e.g., FFT-1D and FFT-2D), performing much higher than other accelerators.

For traditional HPC applications with abundant parallel floating-point operations, it seems to be difficult for FPGAs to beat the performance of other accelerators, even though FPGAs can be much more power-efficient.

- Tested FPGA does not contain dedicated, embedded floating-point cores, while others have fully-optimized floating-point computation units.

Current and upcoming high-end FPGAs are equipped with hardened floating-point operators, whose performance will be comparable to other accelerators, while remaining power-efficient.

# Hardware Resource Utilization (%)

Hardware resource utilization (%) depending on the number of the replicated compute units (CUs) and SIMD width in the kernel vectorization

App	Number of the replicated CUs, SIMD width in the kernel vectorization																	
	1,1	1,2	1,4	1,8	1,16	2,1	2,2	2,4	2,8	2,16	4,1	4,2	4,4	4,8	4,16	8,1	8,2	
Jacobi	29	33	37	41	49	36	43	51	59	74	48	62	78	95	124	73	101	
MatMul	28	34	45	67	109	35	46	68	110	195	48	69	112	197	367	72	115	
SpMul	35	-	-	-	-	46	-	-	-	-	69	-	-	-	-	114	-	
HotSpot	56	79	124	214	443	89	134	224	445	863	154	245	467	866	1704	285	518	
NW	35	46	68	112	200	46	68	112	200	377	69	113	201	377	730	115	202	
SRAD	54	65	80	110	170	84	106	136	197	317	145	189	249	370	621	266	354	
FFT-1D	80	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	
FFT-2D	56	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	-	

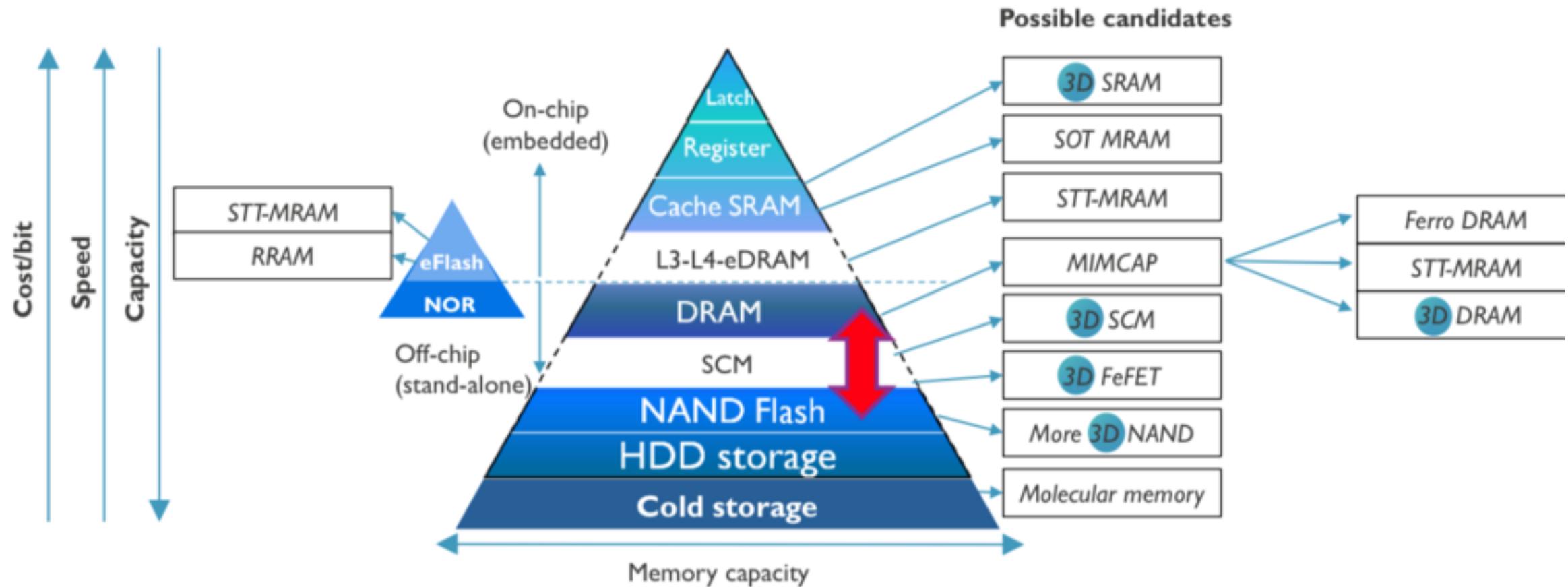
# of CU affects the resource utilization more than the SIMD width.

If a resource utilization is larger than 100%, the compiler cannot generate kernel execution file.

# Emerging Memory Systems

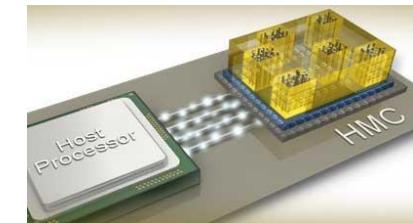


# Memory Hierarchy is Specializing too



# Memory Systems Started Diversifying Several Years Ago

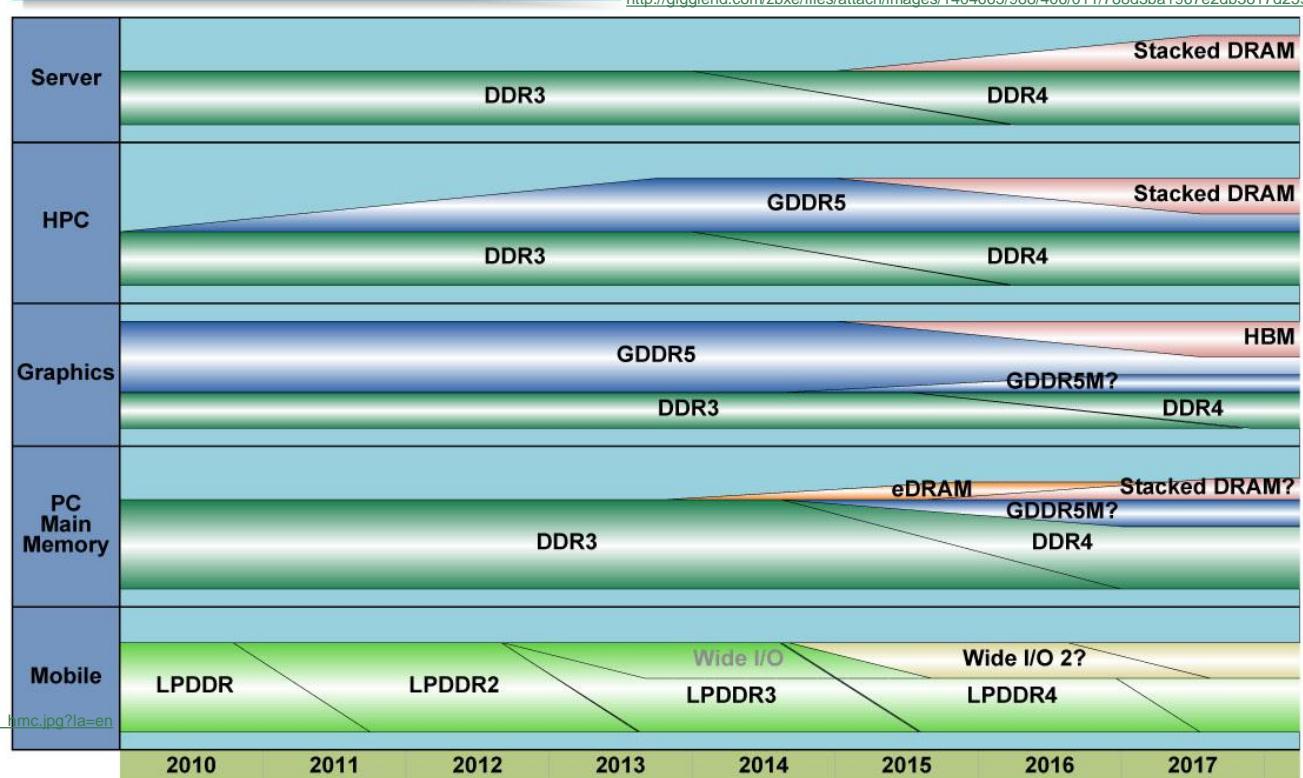
- Architectures
  - HMC, HBM/2/3, LPDDR4, GDDR5X, WIDEIO; etc
  - 2.5D, 3D Stacking
- Configurations
  - Unified memory
  - Scratchpads
  - Write through, write back, etc
  - Consistency and coherence protocols
  - Virtual v. Physical, paging strategies
- New devices
  - ReRAM, PCRAM, STT-MRAM, 3D-Xpoint



[https://www.micron.com/~media/track-2-images/content-images/content\\_image\\_hmc.jpg?la=en](https://www.micron.com/~media/track-2-images/content-images/content_image_hmc.jpg?la=en)

## DRAM Transition

<http://gigglehd.com/zixe/files/attach/images/1404665/988/406/011/788d3ba1967e2db3817d259d2e>



Copyright (c) 2014 Hiroshige Goto All rights reserved.

	SRAM	DRAM	eDRAM	2D NAND Flash	3D NAND Flash	PCRAM	STTRAM	2D ReRAM	3D ReRAM
Data Retention	N	N	N	Y	Y	Y	Y	Y	Y
Cell Size ( $\text{F}^2$ )	60-200	4-6	19-26	2-5	<1	4-10	8-40	4	<1
Minimum F demonstrated (nm)	14	25	22	16	64	20	28	27	24
Read Time (ns)	<1	30	5	$10^9$	$10^8$	10-50	3-10	10-50	10-50
Write Time (ns)	<1	50	5	$10^9$	$10^8$	100-300	3-10	10-50	10-50
Number of Rewrites	$10^{10}$	$10^{10}$	$10^{10}$	$10^9-10^8$	$10^8-10^7$	$10^{10}$	$10^9-10^{10}$	$10^9-10^{10}$	
Read Power	Low	Low	Low	High	High	Low	Medium	Medium	Medium
Write Power	Low	Low	Low	High	High	High	Medium	Medium	Medium
Power (other than R/W)	Leakage	Refresh	Refresh	None	None	None	Sleep	Sleep	
Maturity									

J.S. Vetter and S. Mittal, "Opportunities for Nonvolatile Memory Systems in Extreme-Scale High Performance Computing," CISE, 17(2):73-82, 2015.

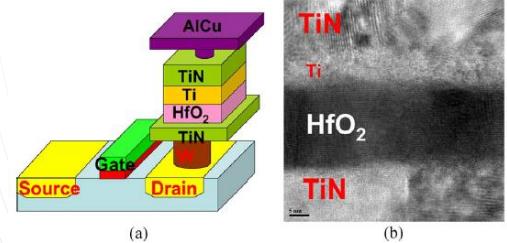
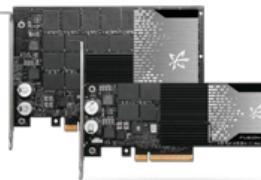


Fig. 4. (a) A typical 1T1R structure of RRAM with  $\text{HfO}_2$ ; (b) HR-TEM image of the  $\text{TiN}/\text{Ti}/\text{HfO}_2/\text{TiN}$  stacked layer; the thickness of the  $\text{HfO}_2$  is 20 nm.

H.S.P. Wong, H.Y. Lee, S. Yu et al., "Metal-oxide RRAM," Proceedings of the IEEE, 100(6), 1911-1922, 2012.

# NVRAM Technology Continues to Improve – Driven by Broad Market Forces



designlines MEMORY

Blog

## First Look at Samsung's 48L 3D V-NAND Flash

Kevin Gibb, Product Line Manager  
TechInsights

4/6/2016 04:40 PM EDT

9 comments post a comment

tom'sHARDWARE

PRODUCT REVIEWS NEWS DEALS FORUM

## Samsung's 10-Year Plan Starts With 128TB QLC SSD, 960 Successor

The highly anticipated Samsung memory is out in the market, first look.

Samsung had announced its 256Gb K9AFY8S0M 3D V-NAND as would be used in a variety of solid state drives (SSD), and would be on the market in early 2016. True to their word, we managed to find them in their 2 TB capacity, mSATA, T3 port.

Figure 1.

16

22  
COMMENTS

designlines WIRELESS & NETWORKING

Slideshow

## Facebook Likes Intel's 3D XPoint

Google joins open hardware effort

Rick Merritt

May 18, 2016

## IBM Puts 3D XPoint on Notice with 3 Bits/Cell PCM Breakthrough

Tiffany Trader



designlines MEMORY

News & Analysis

## 3D NAND Flash at 2 Cents per GB

BeSang wants to lower barrier to 3D NAND flash

R. Colin Johnson

7/18/2016 07:10 PM EDT

14 comments

14 14 4

NO RATINGS 1 saves LOGIN TO RATE

Original URL: [http://www.theregister.co.uk/2013/11/01/hp\\_memristor\\_2018/](http://www.theregister.co.uk/2013/11/01/hp_memristor_2018/)

HP 100TB Memristor drives by 2018 – if you're lucky, admits tech titan  
Universal memory slow in coming

By Chris Mellor

Posted in Storage, 1st November 2013 02:28 GMT

Blocks and Files HP has warned *E! Reg* not to get its hopes up too high after the tech titan's CTO Martin Fink suggested StoreServ arrays could be packed with 100TB Memristor drives come 2018.

designlines MEMORY

News & Analysis

## Samsung Debuts 3D XPoint Killer

3D NAND variant stakes out high-end SSDs

Rick Merritt

8/11/2016 00:01 AM EDT

5 comments

56 12 212 4

NO RATINGS 1 saves LOGIN TO RATE

SANTA CLARA, Calif. – Samsung lobbed a new variant of its 3D

## Memory Forecast to Account for 53% of Semiconductor Capex

By Dylan McGrath, 08/29/18 00:00

Share Post Share on Facebook Share on Twitter in

SAN FRANCISCO — Capital spending for memory chips is expected to account for 53% of the semiconductor industry capex of \$102 billion this year, nearly twice the percentage that memory accounted for five years ago, according to market research firm IC Insights.

With all NAND flash vendors ramping up 3D NAND capacity, NAND-related capital expenditures are forecast to total more than \$31 billion, 31% of the semiconductor industry total, according to the latest edition of IC Insights' McClean Report. The total for NAND capex would represent an increase of 13% over 2017, when NAND flash capex grew by 91%.

Meanwhile, the report forecasts that capital spending for DRAM and SRAM will increase 41% in 2018, while spending for any other industry segment, growing 41% in 2018 after an 82% increase last year. DRAM capex is expected to total \$22.9 billion, 22% of the industry-wide total, according to the report.

Memory's Growing Share of Capital Spending



ars TECHNICA

BIZ & IT TECH SCIENCE POLICY CARS GAMING & CULTURE

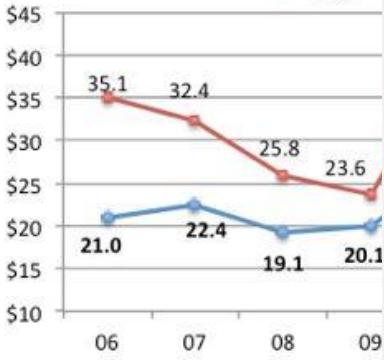
THE REVOLUTION IS HERE —

## Intel at last announces Optane memory: DDR4 that never forgets

New memory offers huge capacities and persistence, but fits in a DDR4 slot.

PETER BRIGHT - 5/30/2018, 8:45 PM

Flash



SanDisk 1TB  
Extreme UHS-I ...  
\$449.99  
B&H Photo

ES, FI, 10, Be superlatively treated 3D tec os with the the chip technology seems to have since lecture for NAND implementation of its door to offering to price that reduces about 2¢ per

46 PM 7,391 VIEWS

and Micron Jointly Announce  
Changing 3D XPoint Memory  
Technology

209

# Many Memory Architecture Options under Consideration...

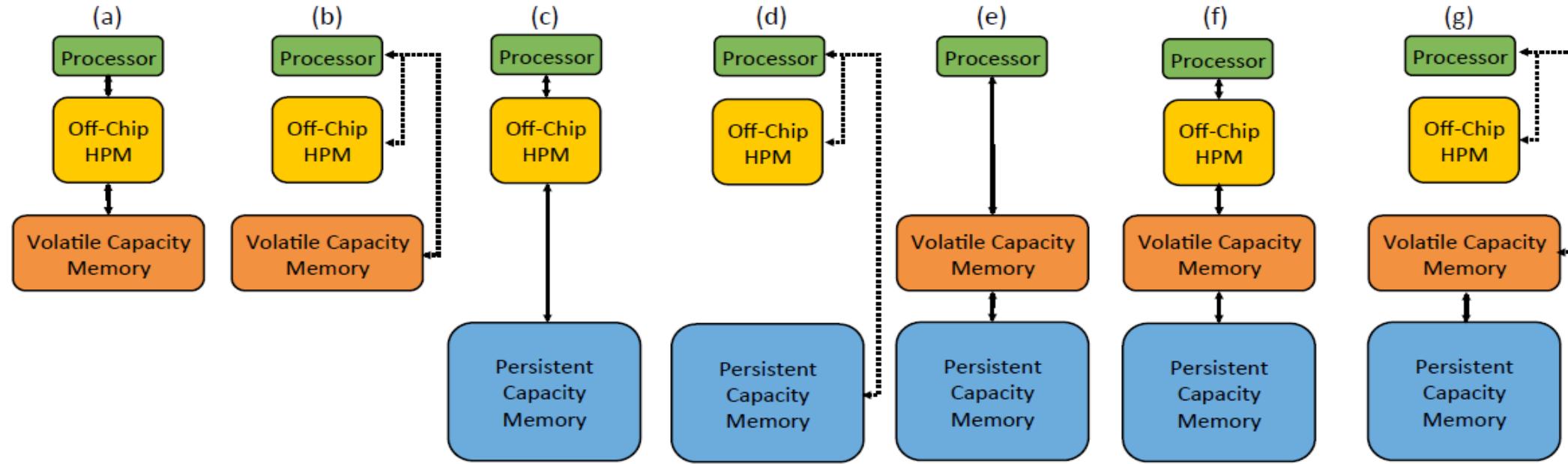


TABLE I: Comparison of four tiers of recent memory technologies [9], [11], [17], [18], [22]–[25], [28], [30], [35], [39], [40], [47]–[49].

	Volatile	Density (GB)	BW (GB/s)	Est. Cost	Speed	Latency
HMC2.0	✓	4-8	320	3x	30 Gbps	~100s ns
HBM2	✓	2-8	256	2x	2 Gbps	~100s ns
GDDR6	✓	8-16	72	2x	18 Gbps	~100s ns
WIO2	✓	8-32	68	2x	1,066 MT/s	~100s ns
DDR4	✓	2-16	25.6	1x	3,200 MT/s	20-50 ns
STT-MRAM	✗	0.5	-	1x	1,600 MT/s	10-50 ns
PCM	✗	1	3.5	1x	3M IOPS	50-100 ns
3D-Xpoint	✗	750	2.4	0.5x	550K IOPS	10 $\mu$ s
Z-NAND	✗	800	3.2	0.5x	750K IOPS	12-20 $\mu$ s
NAND Flash	✗	>1,000	<3	0.1x	50K IOPS	25-125 $\mu$ s

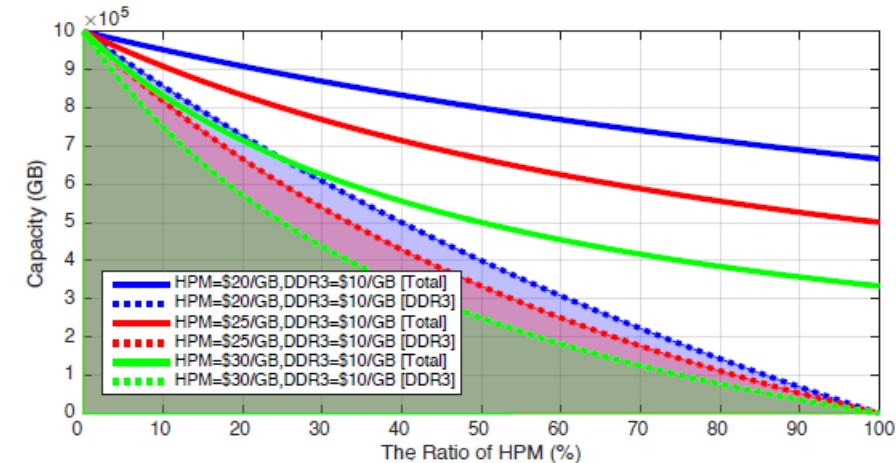
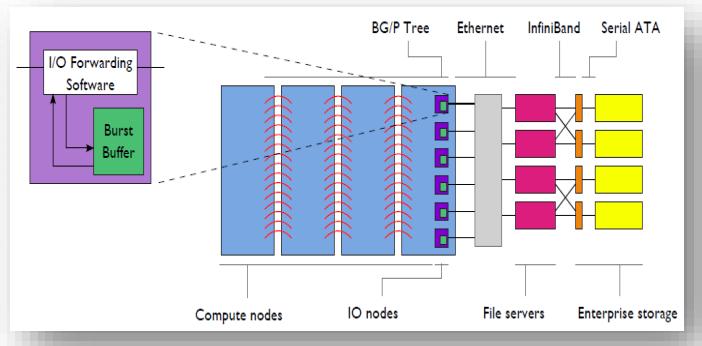


Fig. 1: Possible configurations of a memory system using DDR3 and HPM of different costs under a fixed budget.

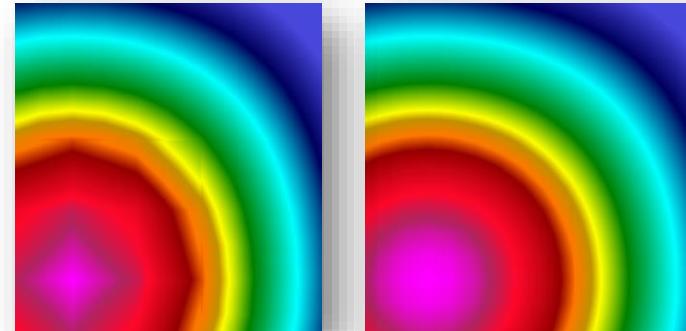
# Considerations for Programming NVM Systems

# NVM Opportunities in Applications

- Burst Buffers, C/R [Liu, et al., MSST 2012]



- In situ visualization and analytics



<http://ft.ornl.gov/eavl>

- Persistent data structures like materials tables

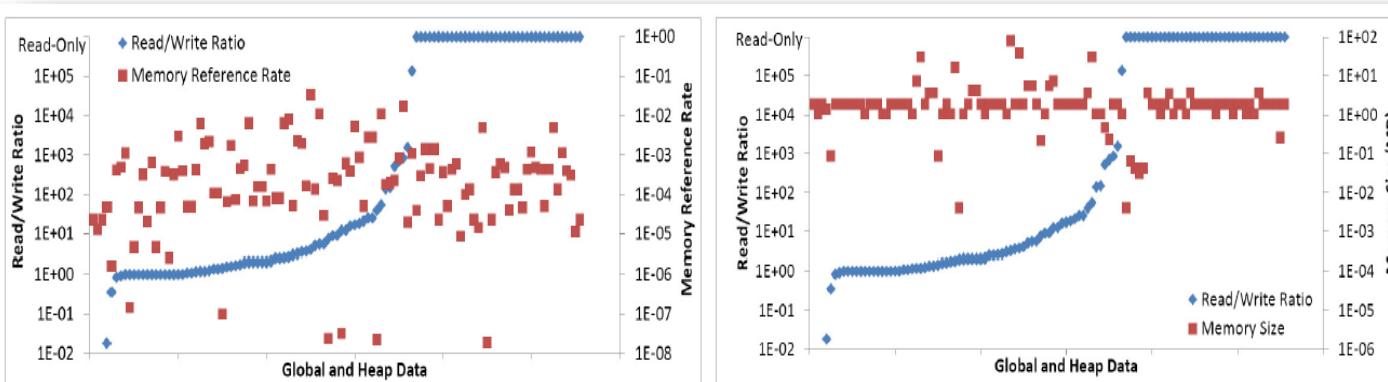


Figure 3: Read/write ratios, memory reference rates and memory object sizes for memory objects in Nek5000

Empirical results show many reasons...

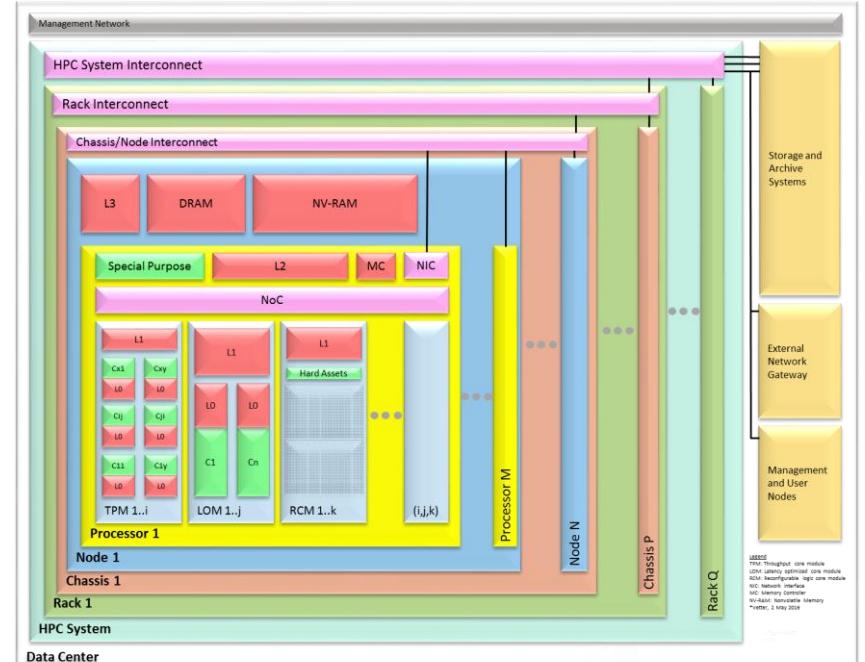
- Lookup, index, and permutation tables
- Inverted and 'element-lagged' mass matrices
- Geometry arrays for grids
- Thermal conductivity for soils
- Strain and conductivity rates
- Boundary condition data
- Constants for transforms, interpolation
- MC Tally tables, cross-section materials tables...

# NVM Design Choices

- Dimensions
  - Integration point
  - Exploit persistence
    - ACID?
  - Scalability
  - Programming model

- Our Approaches

- Transparent access to NVM from GPU
- NVL-C: expose NVM to user/applications
- Papyrus: parallel aggregate persistent memory
- Many others (See S. Mittal and J. S. Vetter, "A Survey of Software Techniques for Using Non-Volatile Memories for Storage and Main Memory Systems," in IEEE TPDS 27:5, pp. 1537-1550, 2016)



<http://j.mp/nvm-sw-survey>

IEEE TRANSACTIONS ON PARALLEL AND DISTRIBUTING SYSTEMS

## A Survey of Software Techniques for Using Non-Volatile Memories for Storage and Main Memory Systems

Sparsh Mittal, Member, IEEE, and Jeffrey S. Vetter, Senior Member, IEEE

**Abstract**—Non-volatile memory (NVM) devices, such as Flash, phase change RAM, spin transfer torque RAM, and resistive RAM, offer several advantages and challenges when compared to conventional memory technologies, such as DRAM and magnetic hard disk drives (HDDs). In this paper, we present a survey of software techniques that have been proposed to exploit the advantages and mitigate the disadvantages of NVMs when used for designing memory systems, and, in particular, secondary storage (e.g., solid state drive) and main memory. We classify these software techniques along several dimensions to highlight their similarities and differences. Given that NVMs are growing in popularity, we believe that this survey will motivate further research in the field of software technology for NVMs.

**Index Terms**—Review, classification, non-volatile memory (NVM) (NVRAM), flash memory, phase change RAM (PCM) (PCRAM), spin transfer torque RAM (STT-RAM) (STT-MRAM), resistive RAM (ReRAM) (RRAM), storage class memory (SCM), Solid State Drive (SSD).

# **Transparent Runtime Support for NVM from GPUs**

# DRAGON : Expanding the memory capacity of GPUs

- GPUs have limited memory capacity
- Recent GPUs have added paging support to host memory
- Recent datasets have grown larger than host memory
- Extend GPUs to NVM
  - Support for massive data
  - Support for temporary data
  - Support for read-only data
- Good performance (including surprises)

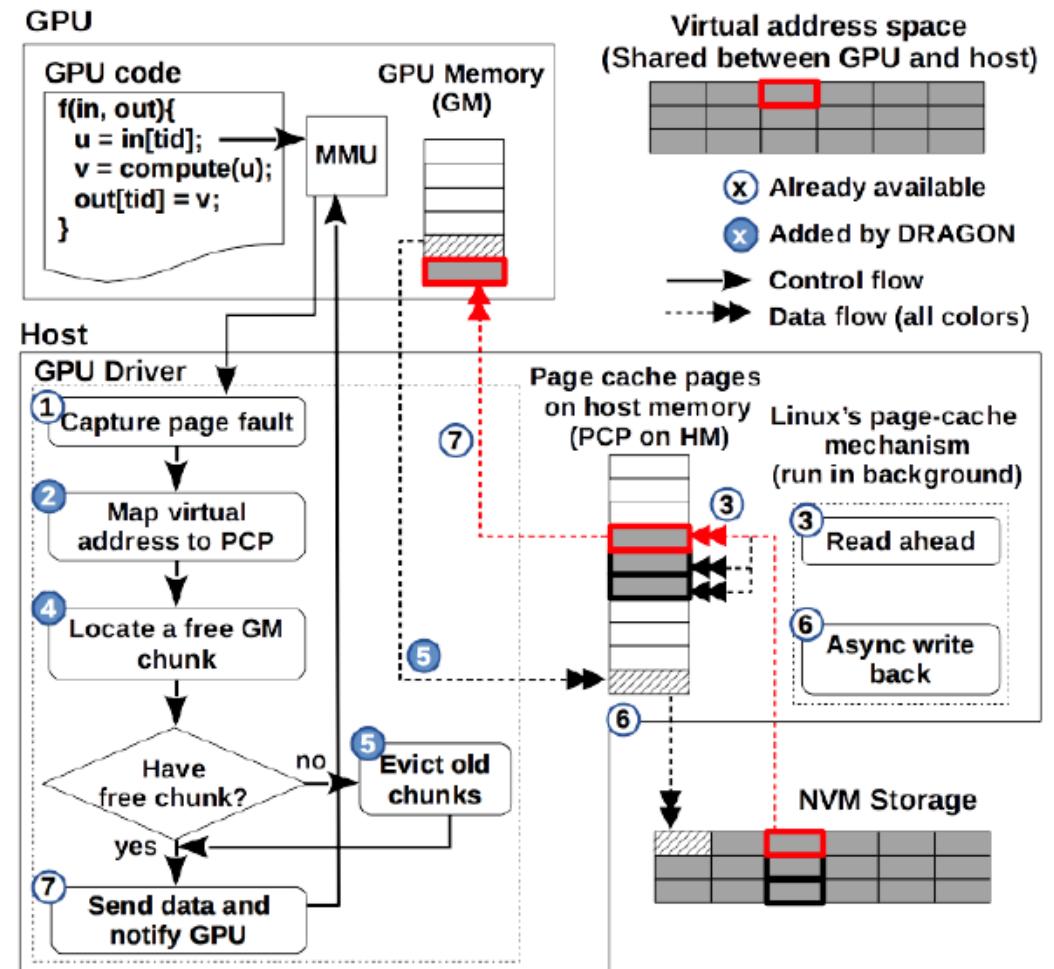


Fig. 1: DRAGON driver operation

# DRAGON: API and Integration

## Out-of-Core using CUDA

```
// Allocate host & device memory
h_buf = malloc(size);
cudaMalloc(&g_buf, size);

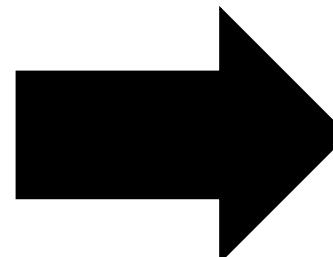
while() { // go over all chunks
    // Read-in data
    f = fopen(filepath, "r");
    fread(h_buf, size, 1, f);

    // H2D Transfer
    cudaMemcpy(g_buf, h_buf, H2D);

    // GPU compute
    compute_on_gpu(g_buf);

    // Transfer back to host
    cudaMemcpy(h_buf, g_buf, D2H);
    compute_on_host(h_buf);

    // Write out result
    fwrite(h_buf, size, 1, f);
}
```



## DRAGON

```
// mmap data to host and GPU
dragon_map(filepath, size,
           D_READ | D_WRITE, &g_buf);

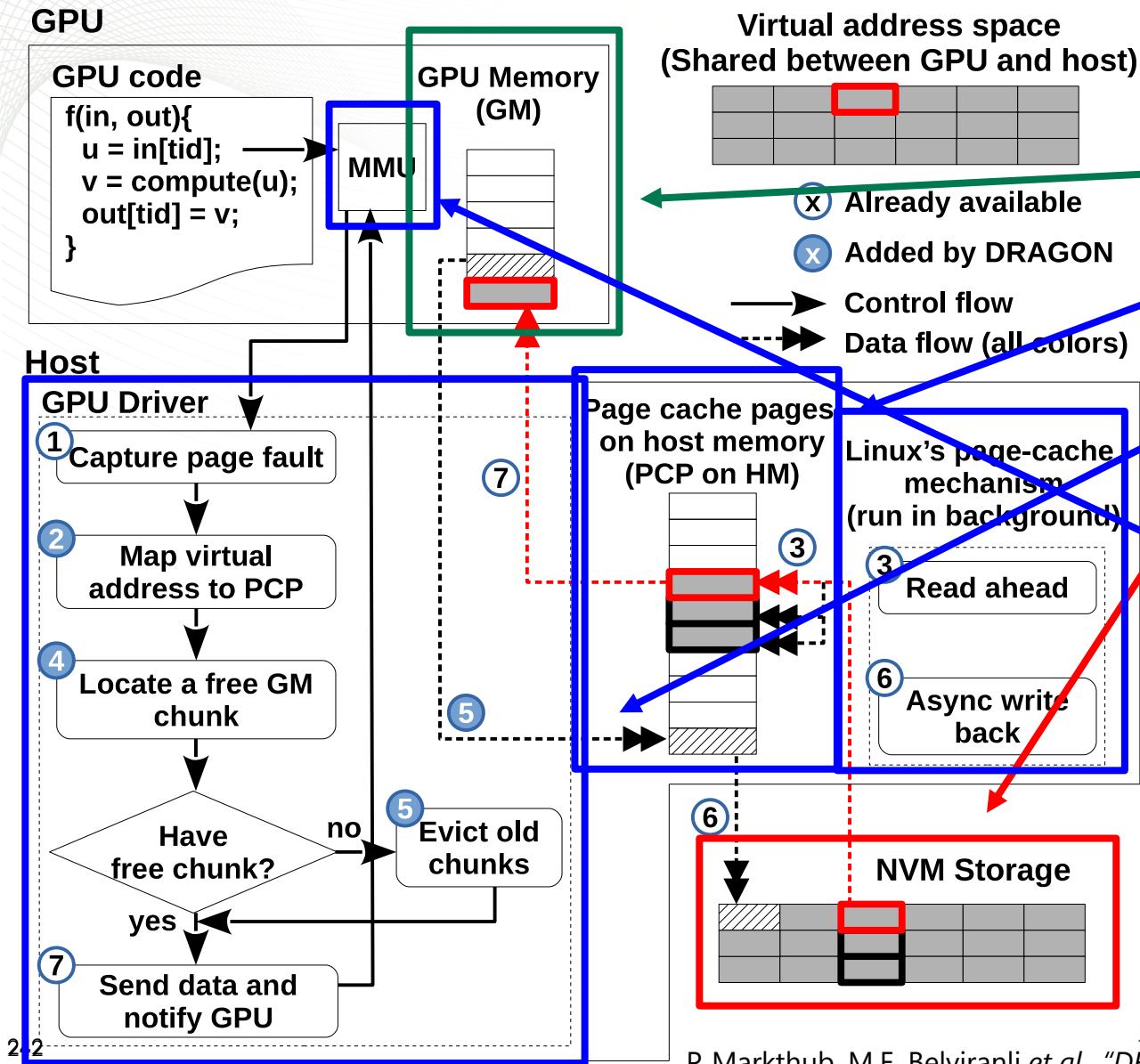
// Accessible on both host and GPU
compute_on_gpu(g_buf);
compute_on_host(g_buf);

// Implicitly called when program
// exits
dragon_sync(g_buf);
dragon_unmap(g_buf);
```

## Notes

- Similar to NVIDIA's Unified Memory (UM)
- Enable access to large memory on NVM
  - **UM is limited by host memory**

# DRAGON Operations: Key Components



- **Three memory spaces:**
  - **GPU Mem (GM)** as 1<sup>st</sup> level cache
  - **Host Mem (HM)** as 2<sup>nd</sup> level cache
  - **NVM** as primary storage
- **Modified GPU driver**
  - Manage data movement & coherency
- **GPU MMU with HW Page Fault**
  - Manage GPU virtual memory mapping
- **Page cache**
  - Buffer & accelerate data access

# Results with Caffe

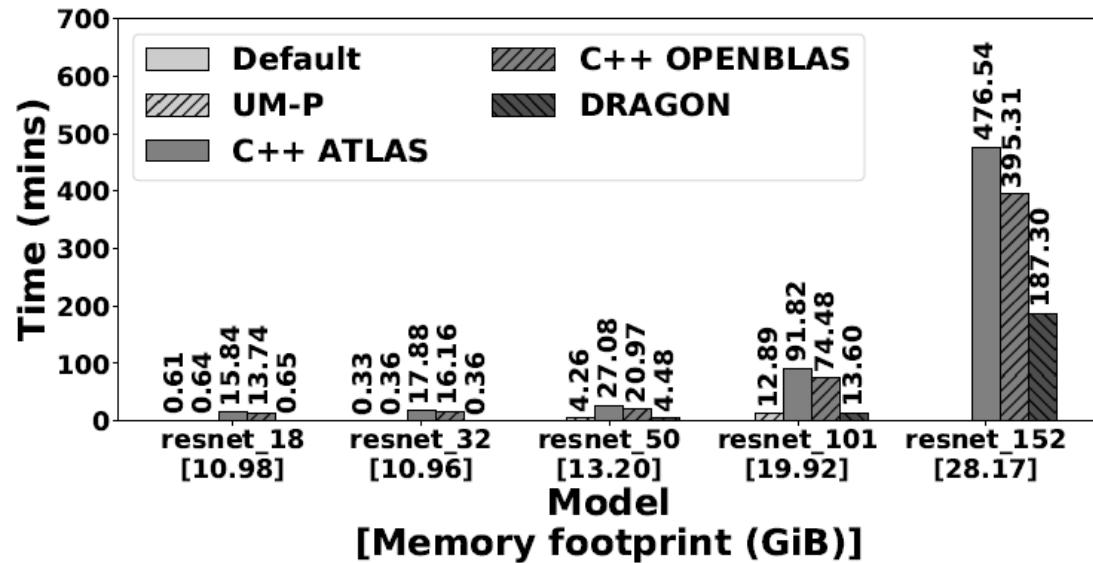


Figure 6: Comparison of ResNet execution times on Caffe.

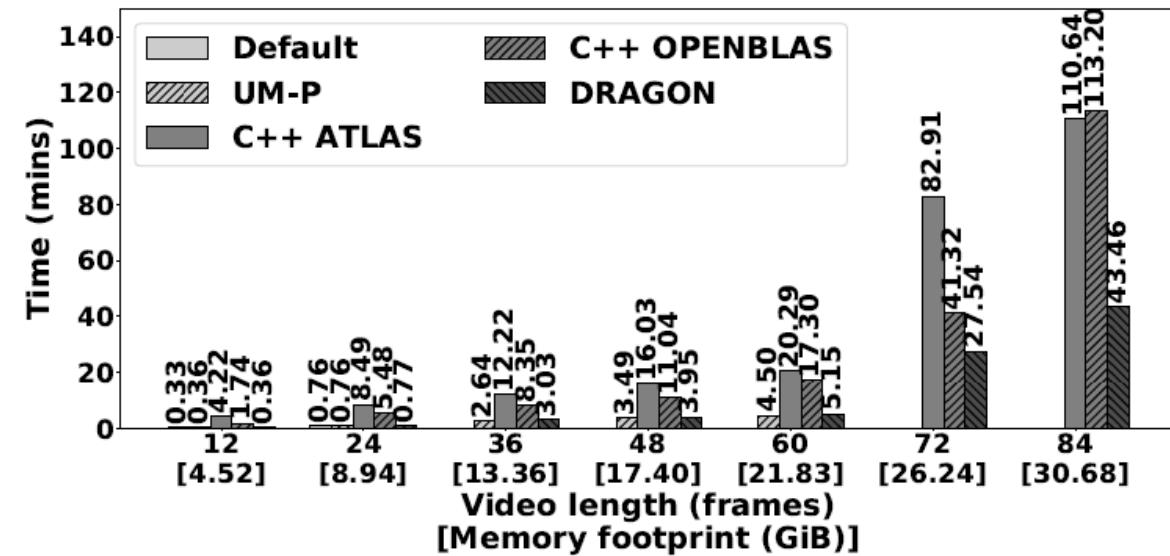


Figure 7: Comparison of C3D the execution times on Caffe.

- Improves capability and productivity
  - Larger problem sizes transparently
  - Handles irregularity easily
  - Surprising performance on applications

# **Language support for NVM: NVL-C - extending C to support NVM**

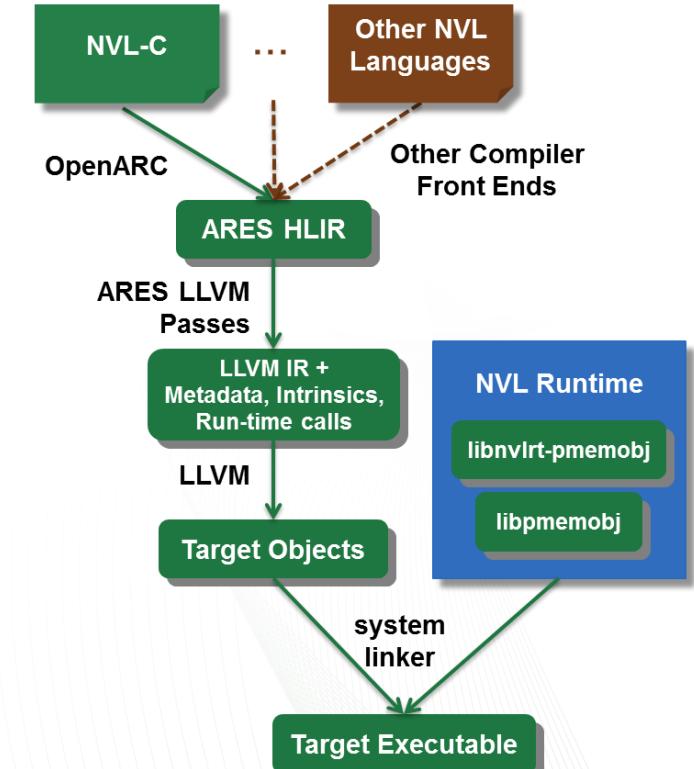
# NVL-C: Portable Programming for NVMM

- Minimal, familiar, programming interface:
  - Minimal C language extensions.
  - App can still use DRAM.
- Pointer safety:
  - Persistence creates new categories of pointer bugs.
  - Best to enforce pointer safety constraints at compile time rather than run time.
- Transactions:
  - Prevent corruption of persistent memory in case of application or system failure.
- Language extensions enable:
  - Compile-time safety constraints.
  - NVM-related compiler analyses and optimizations.
- LLVM-based:
  - Core of compiler can be reused for other front ends and languages.
  - Can take advantage of LLVM ecosystem.

```
#include <nvl.h>
struct list {
    int value;
    nvl struct list *next;
};
void remove(int k) {
    nvl_heap_t *heap
        = nvl_open("foo.nvl");
    nvl struct list *a
        = nvl_get_root(heap, struct list);
    #pragma nvl atomic
    while (a->next != NULL) {
        if (a->next->value == k)
            a->next = a->next->next;
        else
            a = a->next;
    }
    nvl_close(heap);
}
```

Pointer Class	Permitted
NV-to-V	no
V-to-NV	yes
intra-heap NV-to-NV	yes
inter-heap NV-to-NV	no

Table 1: Pointer Classes



# Design Goals: Avoiding persistent data corruption

- New categories of pointer bugs:
  - Caused by multiple memory types:
    - E.g., pointer from NVM to volatile memory will become dangling pointer
  - Prevented at compile time or run time
- Automatic reference counting:
  - No need to manually free
  - Avoids leaks and dangling pointers
- Transactions:
  - Avoids persistent data corruption across software and hardware failures
- High performance:
  - Performance penalty from memory management, pointer safety, and transactions
  - Compiler-based optimizations
  - Programmer-specified hints

# Programming Model: NVM Pointers

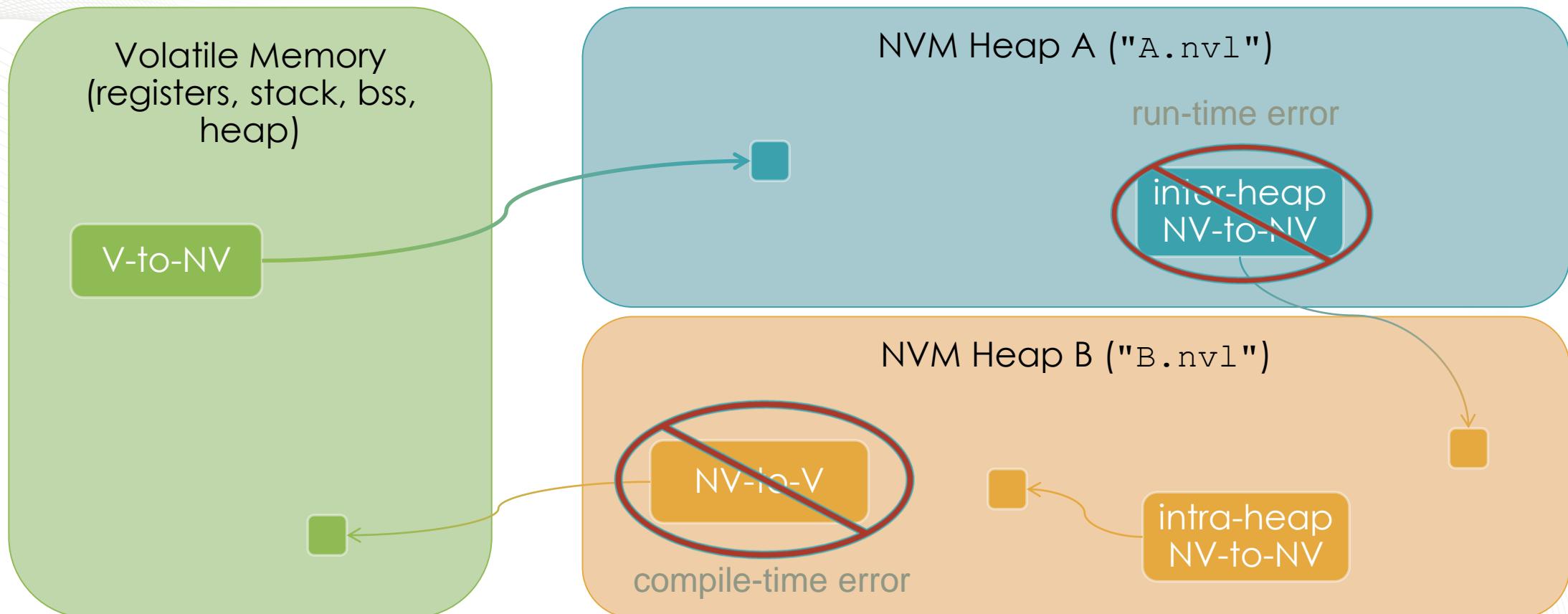
```
#include <nvl.h>
struct list {
    int value;
    nvl struct list *next;
};

void add(int k, nvl struct list *after) {
    struct list *node
        = malloc(sizeof(struct list));
    node->value = k;
    node->next = after->next;
    after->next = node;
}
```

*compile-time error  
explicit cast won't help*

- **nvl type qualifier:**
  - Indicates NVM storage
  - On target type, declares NVM pointer
  - No NVM-stored local or global variable
- **Stricter type safety for NVM pointers:**
  - Does not affect other C types
  - Avoids persistent data corruption
  - Facilitates compiler analysis
  - Needed for automatic reference counting
  - E.g., pointer conversions involving NVM pointers are strictly prohibited

# Programming Model: Pointer types (like Coburn et al.)



avoids dangling pointers when  
memory segments close

# Programming Model: Transactions: Purpose

- Ensures data consistency
- Handles unexpected application termination:
  - Hardware failure (e.g., power loss)
  - Application or OS failure (e.g., segmentation fault)
  - NVL-C safety constraint violation (e.g., inter-heap NV-to-NV pointer)
- Does not handle concurrent access to NVM:
  - Future work
  - Concurrency is still possible
  - Programmer must safeguard NVM data from concurrent access

# Programming Model: Transactions: MATMUL Example

```
#include <nvl.h>
void matmul (nvl float a[I][J],
              nvl float b[I][K],
              nvl float c[K][J])
{
    for (int i=0; i<I; ++i) {
        for (int j=0; j<J; ++j) {
            float sum = 0.0;
            for (int k=0; k<K; ++k)
                sum += b[i][k] * c[k][j];
            a[i][j] = sum;
        }
    }
}
```

- All three arrays are stored in NVM
- Progress is not recorded in NVM
- Problem: if power failure before computation is complete, must start over

# Programming Model: Transactions: MATMUL Example

```
#include <nvl.h>
void matmul(nvl float a[I][J],
            nvl float b[I][K],
            nvl float c[K][J],
            nvl int *i)
{
    for (; *i<I; ++*i) {
        for (int j=0; j<J; ++j) {
            float sum = 0.0;
            for (int k=0; k<K; ++k)
                sum += b[*i][k] * c[k][j];
            a[*i][j] = sum;
        }
    }
}
```

- Store  $i$  in NVM
- Caller initializes  $*i$  to 0 when allocated
- To recover after failure, matmul resumes at old  $*i$
- Problem: failure might have occurred before all of  $a[*i-1]$  became durable in NVM due to buffering and caching

# Programming Model: Transactions: MATMUL Example

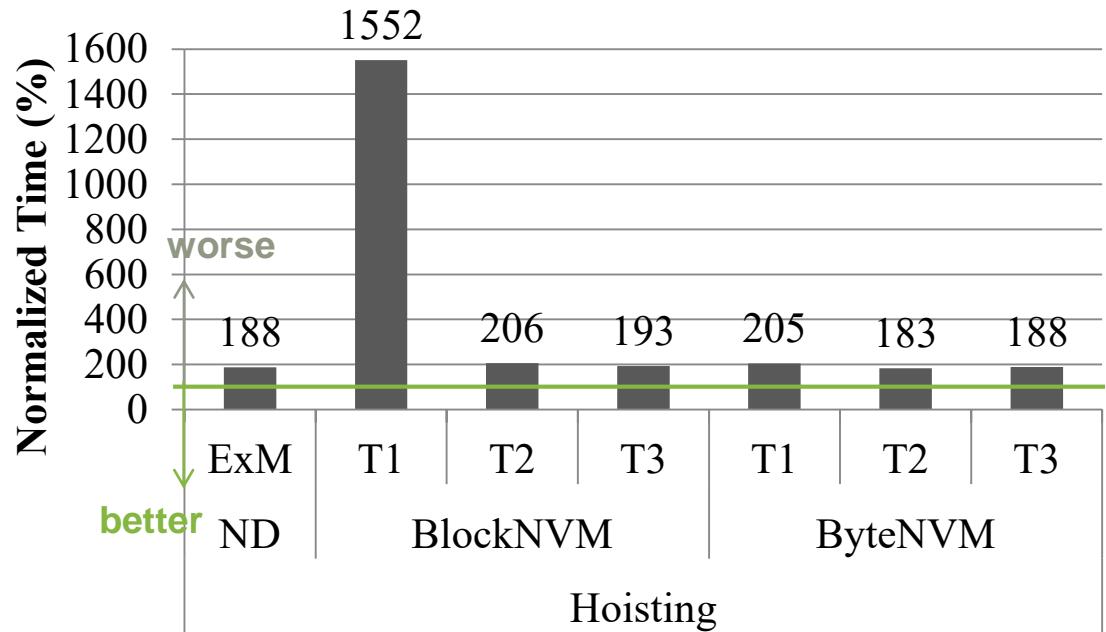
```
#include <nvl.h>
void matmul(nvl float a[I][J],
            nvl float b[I][K],
            nvl float c[K][J],
            nvl int *i)
{
    while (*i<I) {
        #pragma nvl atomic heap(heap)
        {
            for (int j=0; j<J; ++j) {
                float sum = 0.0;
                for (int k=0; k<K; ++k)
                    sum += b[*i][k] * c[k][j];
                a[*i][j] = sum;
            }
            ++*i;
        }
    }
}
```

- **nvl atomic** pragma specifies explicit transaction that computes one row of a
- Transaction guarantees atomicity: both  $*i$  is incremented and one row of  $a$  is written durably, or neither
- Incomplete transaction rolled back after failure

# Programming Model: Transactions: ACID

- Atomicity:
  - Incomplete transaction rolled back next time NVM heap is accessed
- Consistency:
  - Transactions begin and end with NVM data is in a consistent state
  - Implicit transactions: specify NVL-C internal data consistency
  - Explicit transactions: specify application data consistency
- Isolation (handles concurrent access):
  - Not guaranteed yet
- Durability:
  - All NVM writes are durable when transaction commits

# Evaluation: MATMUL



- ExM = use SSD as extended DRAM
- T1 = BSR + transactions
- T2 = T1 + backup clauses
- T3 = T1 + clobber clauses
- BlockNVM = msync included
- ByteNVM = msync suppressed

- Log aggregation (backup) is important for performance
- msync is the culprit
- Skipping undo logs (clobber) has little to improve upon
- NVL-C has minimal overhead

# NVL-C Summary

- Motivated a new programming model for NVM as persistent memory
- Introduced NVL-C, a new programming system for this purpose
  - First class language construct
  - Transactions
- Described several performance optimizations for NVL-C
- Showed performance results for these optimizations on an SSD
- Working on Optane DIMMs now

# Final Report on Workshop on Extreme Heterogeneity

1. Maintaining and improving programmer productivity
  - Flexible, expressive, programming models and languages
  - Intelligent, domain-aware compilers and tools
  - Composition of disparate software components
- Managing resources intelligently
  - Automated methods using introspection and machine learning
  - Optimize for performance, energy efficiency, and availability
- Modeling & predicting performance
  - Evaluate impact of potential system designs and application mappings
  - Model-automated optimization of applications
- Enabling reproducible science despite non-determinism & asynchrony
  - Methods for validation on non-deterministic architectures
  - Detection and mitigation of pervasive faults and errors
- Facilitating Data Management, Analytics, and Workflows
  - Mapping of science workflows to heterogeneous hardware and software services
  - Adapting workflows and services to meet facility-level objectives through learning approaches



# Recap

- Recent trends in extreme-scale HPC paint an ambiguous future
- Complexity is the next major hurdle
  - Heterogeneous compute
  - Deep memory with NVM
- New software solutions
  - Programming
    - Memory
      - DRAGON
      - NVL-C
      - Papyrus
    - Heterogeneity
      - OpenACC->FPGAs
      - Clacc for LLVM
- These changes will have a substantial impact on both software and application design

- Visit us
  - We host interns and other visitors year round
- Jobs in FTG
  - Postdoctoral Research Associate in Computer Science
  - Software Engineer
  - Computer Scientist
  - Visit <http://jobs.ornl.gov>
- Contact me [vetter@ornl.gov](mailto:vetter@ornl.gov)