# Preparing for Extreme Heterogeneity in High Performance Computing

Jeffrey S. Vetter
*With many contributions from FTG Group and Colleagues*
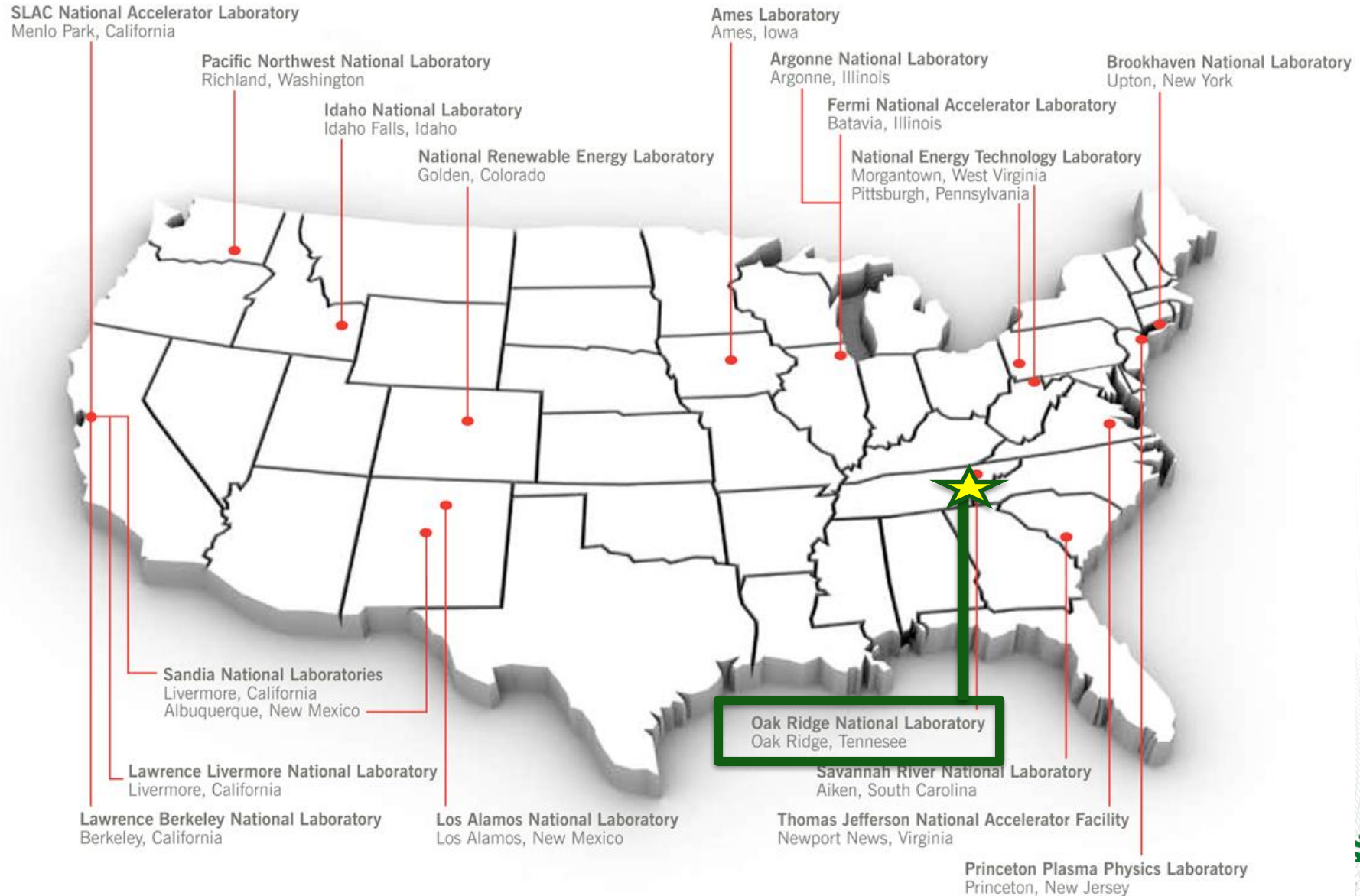
Barcelona Supercomputing Center
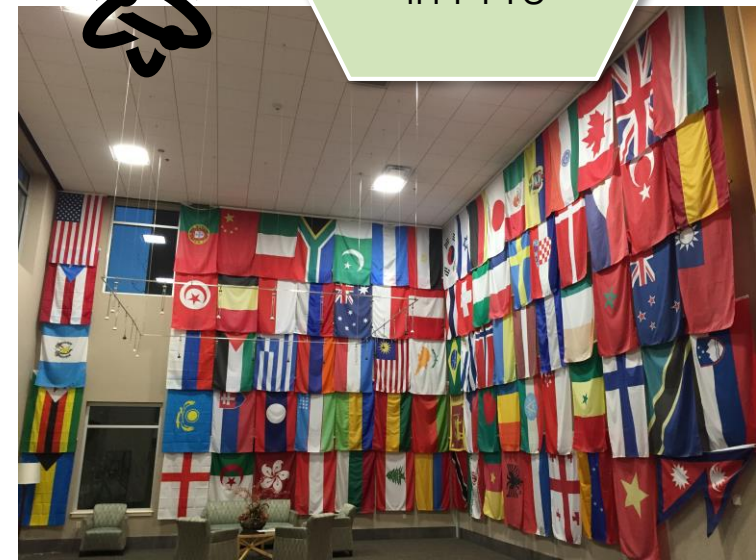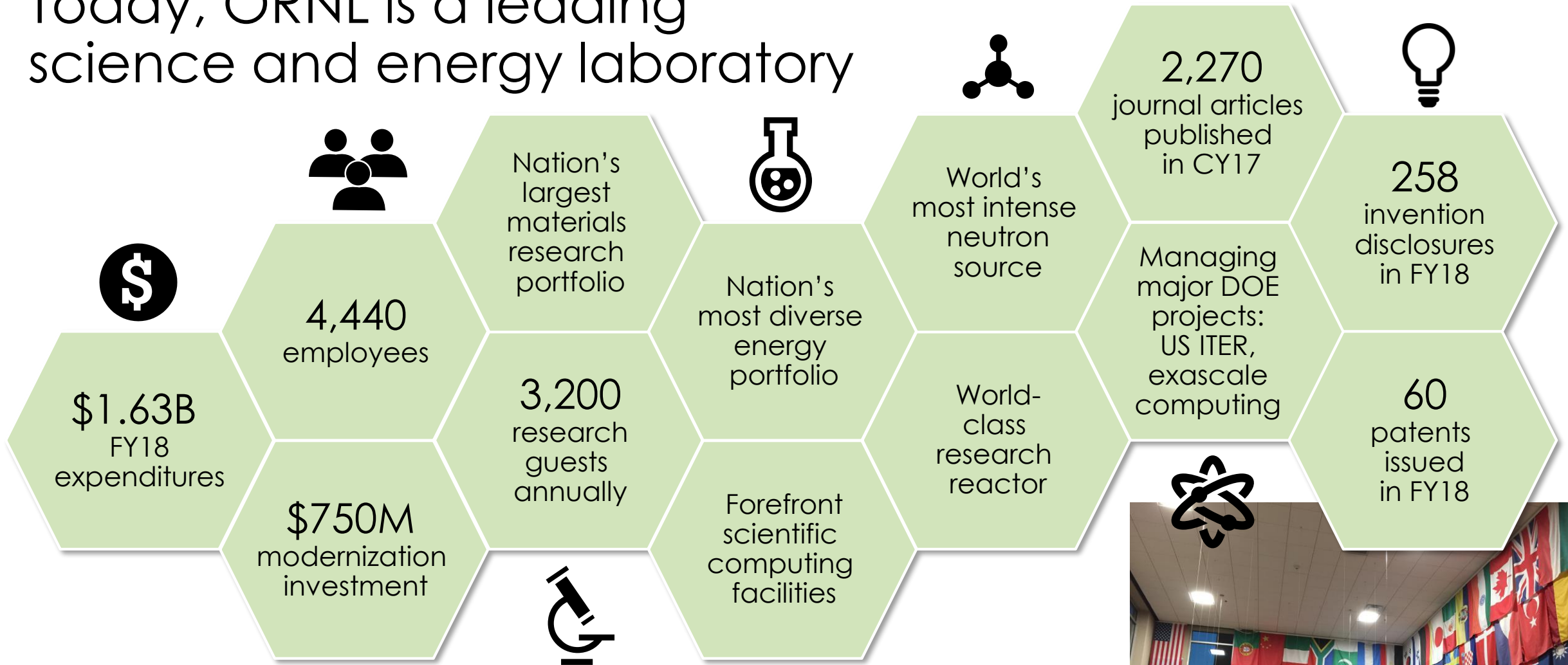Technical University of Catalonia (UPC)
8 May 2019

# Highlights

- Recent trends in extreme-scale HPC paint an uncertain 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!

- Programming systems must provide performance portability (beyond functional portability)!!
  - Heterogeneous processor
    - OpenACC->FGPAs
    - 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)

OAK RIDGE
National Laboratory

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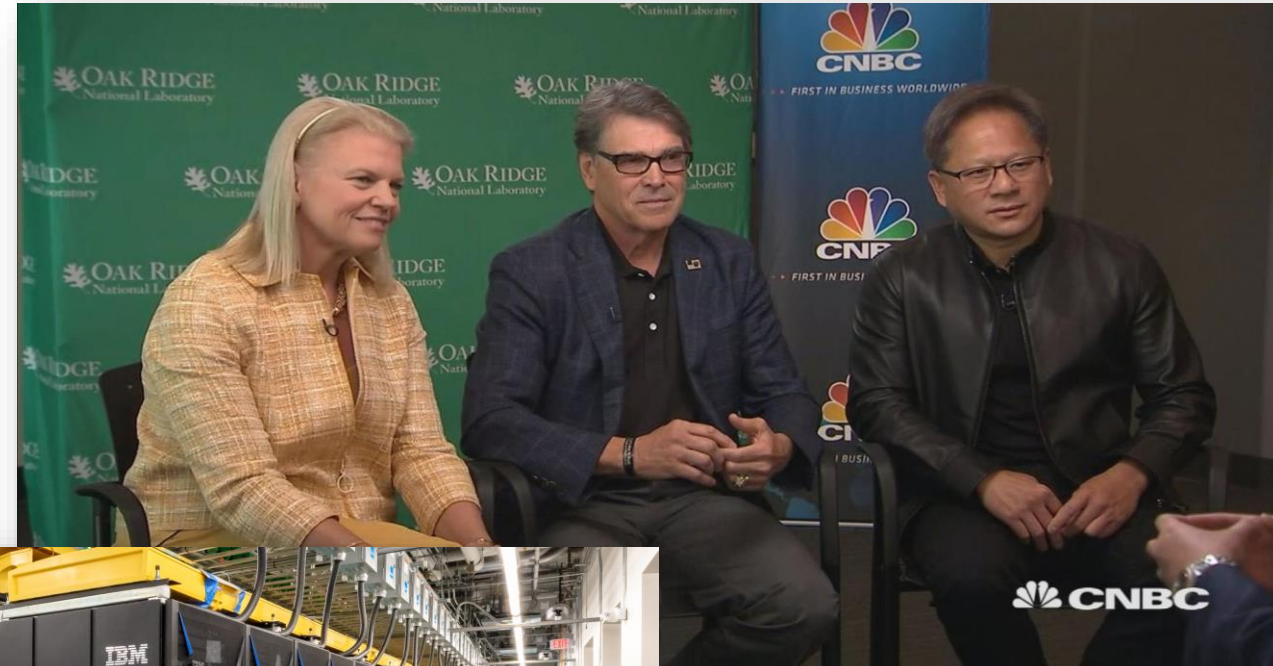


SLAC National Accelerator Laboratory
Menlo Park, California

Pacific Northwest National Laboratory
Richland, Washington

Idaho National Laboratory
Idaho Falls, Idaho

National Renewable Energy Laboratory
Golden, Colorado

Ames Laboratory
Ames, Iowa

Argonne National Laboratory
Argonne, Illinois

Fermi National Accelerator Laboratory
Batavia, Illinois

National Energy Technology Laboratory
Morgantown, West Virginia
Pittsburgh, Pennsylvania

Brookhaven National Laboratory
Upton, New York

Sandia National Laboratories
Livermore, California
Albuquerque, New Mexico

Oak Ridge National Laboratory
Oak Ridge, Tennesee

Savannah River National Laboratory
Aiken, South Carolina

Lawrence Livermore National Laboratory
Livermore, California

Lawrence Berkeley National Laboratory
Berkeley, California

Los Alamos National Laboratory
Los Alamos, New Mexico

Thomas Jefferson National Accelerator Facility
Newport News, Virginia

Princeton Plasma Physics Laboratory
Princeton, New Jersey

# Today, ORNL is a leading science and energy laboratory

- **$1.63B** FY18 expenditures
- **4,440** employees
- **$750M** modernization investment
- Nation's largest materials research portfolio
- **3,200** research guests annually
- Nation's most diverse energy portfolio
- Forefront scientific computing facilities
- World's most intense neutron source
- World-class research reactor
- **2,270** journal articles published in CY17
- Managing major DOE projects: US ITER, exascale computing
- **258** invention disclosures in FY18
- **60** patents issued in FY18

**OAK RIDGE**
National Laboratory

24

#1 on Top 500

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

**OAK RIDGE** National Laboratory

**OAK RIDGE** National Laboratory | **75** YEARS

# 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 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<br>4 Purpose Built AMD Radeon Instinct GPU |
| CPU-GPU Interconnect | AMD Infinity Fabric<br><br>Coherent memory across the node |
| System Interconnect | Multiple Slingshot NICs providing 100 GB/s network bandwidth<br>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)



**Selected program office application development** (BER, BES, NNSA)

**Exascale Computing Project (ECP)**

**Exascale system procurement projects & facilities**
ALCF-3 (Aurora)
OLCF-5 (Frontier)
ASC ATS-4 (El Capitan)

**ECI partners** — US DOE Office of Science (SC) and National Nuclear Security Administration (NNSA)

**ECI mission** — Accelerate R&D, acquisition, and deployment to deliver exascale computing capability to DOE national labs by the early- to mid-2020s

**ECI focus** — Delivery of an *enduring and capable exascale computing capability for use by a wide range of applications* of importance to DOE and the US

*Three Major Components of the ECI*

45

# 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

**ECP Project Office**

**100 R&D TEAMS 1000 RESEARCHERS**

More than 100 top-notch R&D teams

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

**Performant mission and science applications @ scale**

| Foster application development | Ease of use | Diverse architectures | HPC leadership |
|---|---|---|---|

**Application Development (AD)**

Develop and enhance the predictive capability of applications critical to the DOE

**Software Technology (ST)**

Produce expanded and vertically integrated software stack to achieve full potential of exascale computing

**Hardware and Integration (HI)**

Integrated delivery of ECP products on targeted systems at leading DOE computing facilities

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 | Energy security | Economic security | Scientific discovery | Earth system | Health care |
|---|---|---|---|---|---|
| Stockpile stewardship | Turbine wind plant efficiency | Additive manufacturing of qualifiable metal parts | Find, predict, and control materials and properties | Accurate regional impact assessments in Earth system models | Accelerate and translate cancer research |
| Next-generation electromagnetics simulation of hostile environment and virtual flight testing for hypersonic re-entry vehicles | High-efficiency, low-emission combustion engine and gas turbine design | Reliable and efficient planning of the power grid | Cosmological probe of the standard model of particle physics | Stress-resistant crop analysis and catalytic conversion of biomass-derived alcohols | |
| | Materials design for extreme environments of nuclear fission and fusion reactors | Seismic hazard risk assessment | Validate fundamental laws of nature | Metagenomics for analysis of biogeochemical cycles, climate change, environmental remediation | |
| | Design and commercialization of Small Modular Reactors | Urban planning | Demystify origin of chemical elements | | |
| | Subsurface use for carbon capture, petroleum extraction, waste disposal | | Light source-enabled analysis of protein and molecular structure and design | | |
| | Scale-up of clean fossil fuel combustion | | Whole-device model of magnetically confined fusion plasmas | | |
| | Biofuel catalyst design | | | | |

ECP EXASCALE COMPUTING PROJECT

# ECP SW Stack: Strategic Alignment & Synergies



Data & Visualization

Applications | Co-Design

Programming Models Runtimes | Mathematical Libraries | Embedded Data & Visualization | Development Tools

Software Ecosystem & Delivery

Hardware interface

# Many ECP ST products are available (many github)

**For example…**

## Programming Models and Runtimes Products

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

## Development Tools (19)

SICM
QUO
Kitsune
SCR
Caliper
mpiFileUtils
Gotcha
TriBITS
Exascale Code Geneneration Toolkit
PAPI
CHiLL Autotuning Compiler
Search using Rand

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
http://github.com/hpc/mpifileutils
https://tribits.org

http://icl.utk.edu/exa-papi/

## Mathematical Libraries Products (16)

xSDK
hypre
FleCSI
MFEM
Kokkoskernels
Trilinos
SUNDIALS
PETSc/TAO
libEnsemble
STRUMPACK
SuperLU
ForTrilinos
SLATE
MAGMA-sparse
DTK
Tasmanian

https://xsdk.info
http://www.llnl.gov/casc/hypre
http://www.flecsi.org
http://mfem.org/
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…**

EXASCALE COMPUTING PROJECT

# 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 → Fewer 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

Math SDK
Tools SDK
PM&RT SDK
DataViz SDK
Facility SDK

### OpenHPC
Potential exit strategy for binary distributions

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

### 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

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)

http://e4s.io

ECP EXASCALE COMPUTING PROJECT

# Major Trends in Computing

# Contemporary devices are approaching fundamental limits



Economist, Mar 2016



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

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. Gaensslen, 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,*

I.L. Markov, "Limits on fundamental limits to computation," *Nature, 512(7513):147-54, 2014, doi:10.1038/nature13570.*

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

# 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

**OAK RIDGE**
National Laboratory

# Transition Period Predictions

| Optimize Software and Expose New Hierarchical Parallelism | Architectural Specialization and Integration | Emerging Technologies |
|---|---|---|
| • Redesign software to boost performance on upcoming architectures<br>• Exploit new levels of parallelism and efficient data movement | • Use CMOS more efficiently for our workloads<br>• Integrate components to boost performance and eliminate inefficiencies | • Investigate new computational paradigms<br>  • Quantum<br>  • Neuromorphic<br>  • Advanced Digital<br>  • Emerging Memory Devices |

OAK RIDGE
National Laboratory

# 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



Intel's Nervana AI platform takes aim at Nvidia's GPU techology

http://www.theinquirer.net/inquirer/news/2477796/intels-nervana-ai-platform-takes-aim-at-nvidias-gpu-techology

GOOGLE BUILT ITS VERY OWN CHIPS TO POWER ITS AI BOTS

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

NEW AT AMAZON: ITS OWN CHIPS FOR CLOUD COMPUTING

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/

HotChips 2018

HotChips 2018

Xilinx ACAP

https://www.thebroadcastbridge.com/content/entry/1094/altera-announced-...

143

# Analysis of Apple A-* SoCs

http://vlsiarch.eecs.harvard.edu/accelerators/die-photo-analysis

# Memory Hierarchy is Specializing, Expanding, and Diversifying



Image Source: IMEC

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

# 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 |
|---|---|---|---|---|---|
| Application | 1 | 1 | 2 | 2 | 3 |
| Algorithm | 1 | 1 | 2 | 3 | 3 |
| Language | 1 | 2 | 2 | 3 | 3 |
| API | 1 | 2 | 2 | 3 | 3 |
| Arch | 1 | 2 | 2 | 3 | 3 |
| ISA | 1 | 2 | 2 | 3 | 3 |
| Microarch | 2 | 3 | 2 | 3 | 3 |
| FU | 2 | 3 | 2 | 3 | 3 |
| Logic | 3 | 3 | 2 | 3 | 3 |
| Device | 3 | 3 | 2 | 3 | 3 |

Adapted from IEEE Rebooting Computing Chart

OAK RIDGE
National Laboratory

# 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!]　　　　　　**First U.S. Exascale Systems**

| 2012 | 2016 | 2018 | 2020 | 2021-2023 |
|---|---|---|---|---|

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

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

**FRONTIER**
**ORNL**
TBD

Heterogeneous Cores

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

**Theta (24)**
**ANL**
Cray/Intel KNL

Deep Memory incl NVM

**Aurora**
**ANL**
Intel/Cray

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

**Perlmutter**
**LBNL**
Cray/AMD/NVIDIA

Plateauing I/O Performance

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

**Trinity (6)**
**LANL/SNL**
Cray/Intel Xeon/KNL

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

**CROSSROADS**
**LANL/SNL**
TBD

**EL CAPITAN**
**LLNL**
TBD

**OAK RIDGE**
National Laboratory

Jan 2018

# 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

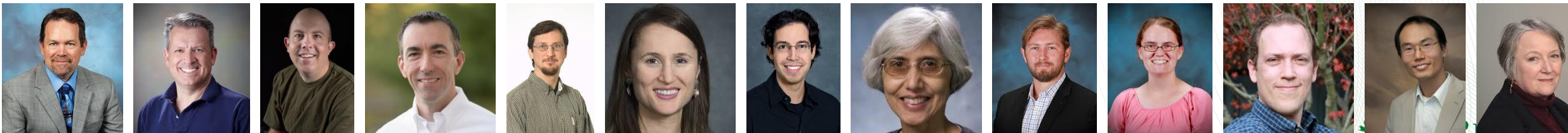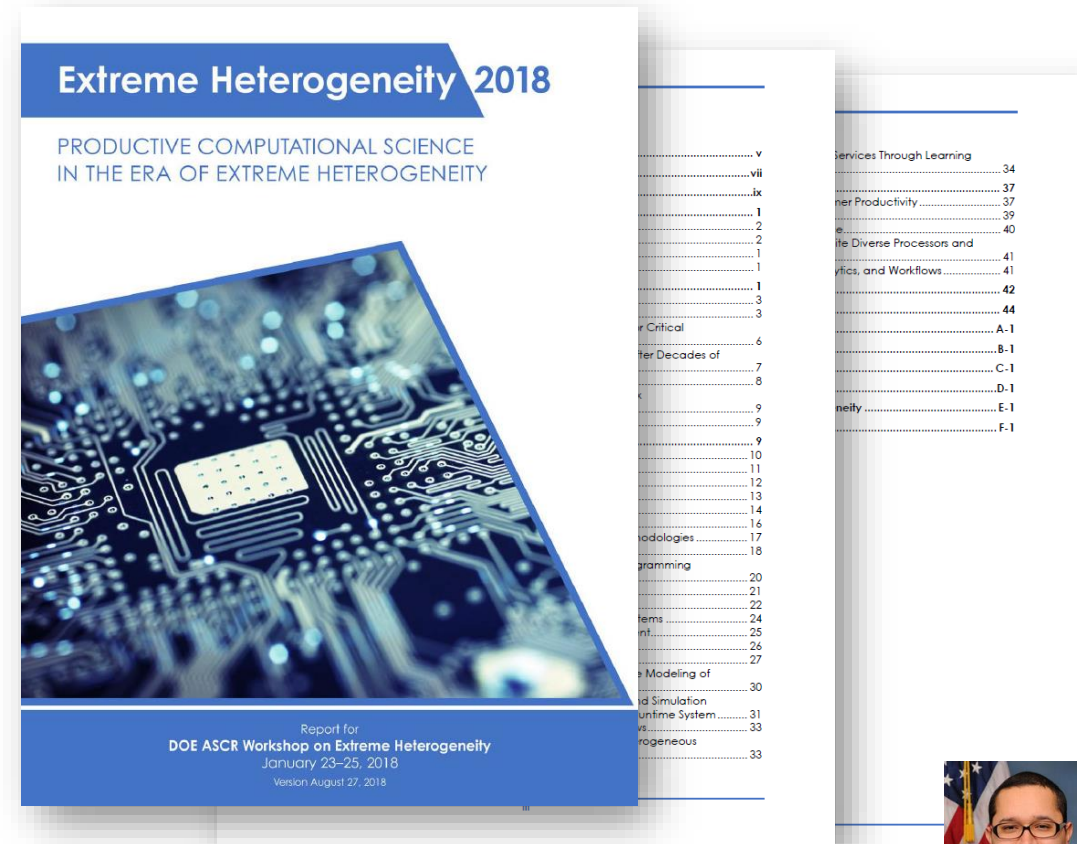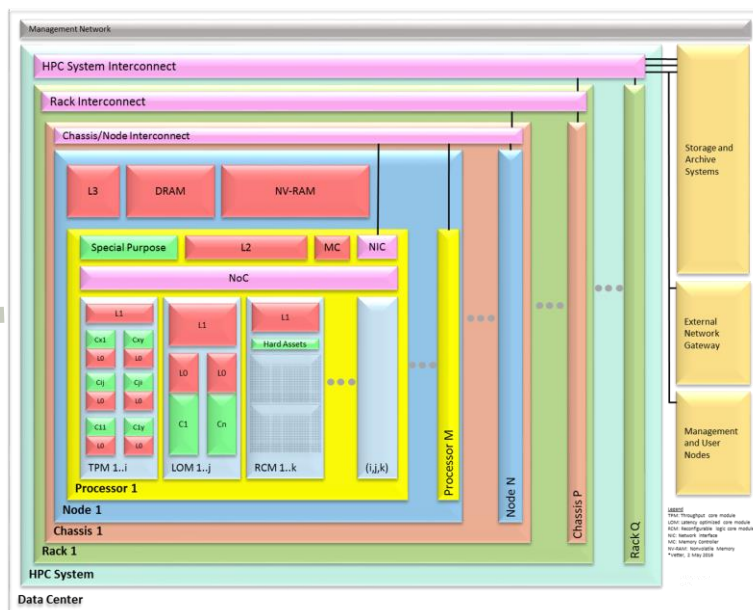- Enabling reproducible science despite non-determinism & asynchrony
   - Methods for validation on non-deterministic architectures
   - Detection and mitigation of pervasive faults and errors

- 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

Extreme Heterogeneity 2018

PRODUCTIVE COMPUTATIONAL SCIENCE
IN THE ERA OF EXTREME HETEROGENEITY

Report for
DOE ASCR Workshop on Extreme Heterogeneity
January 23–25, 2018
Version August 27, 2018

https://orau.gov/exheterogeneity2018/

https://doi.org/10.2172/1473756

OAK RIDGE
National Laboratory

# 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

OAK RIDGE
National Laboratory

# Directive-based Solutions for FPGA Computing

OAK RIDGE
National Laboratory

# 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).

OAK RIDGE
National Laboratory

OpenARC Compiler

OpenARC Runtime

OpenACC

OpenMP 4

NVL-C

Input C Program

OpenARC Front-End

C Parser

Directive Parser

Preprocessor

General Optimizer

OpenARC IR

OpenARC Back-End

Kernels & Host Program Generator

Device Specific Optimizer

OpenARC Auto-Tuner

Tuning Configuration Generator

Search Space Pruner

Feedback

LLVM Back-End

Extended LLVM IR Generator

NVL Passes

Standard LLVM Passes

Output Codes

Kernels for Target Devices

Host Program

Run

CUDA, OpenCL Libraries

HeteroIR Common Runtime with Tuning Engine

CUDA GPU

GCN GPU

Xeon Phi

Altera FPGA

NVM

NVM

NVM

NVM

NVL Runtime

Run

Executable

pmem.io NVM Library

214

OAK RIDGE
National Laboratory
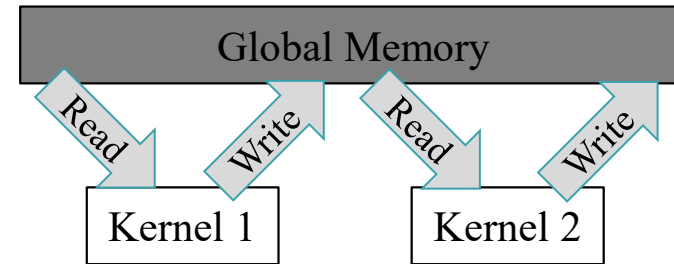
# Baseline Translation of OpenACC-to-FPGA

- Use OpenCL as the output model and the Altera Offline Compiler (AOC) as its backend compiler.

- Translates the input OpenACC program into a host code containing HeteroIR constructs and device-specific kernel codes.
  - Use the same HeteroIR runtime system of the existing OpenCL backends, except for the device initialization.
  - Reuse most of compiler passes for kernel generation.
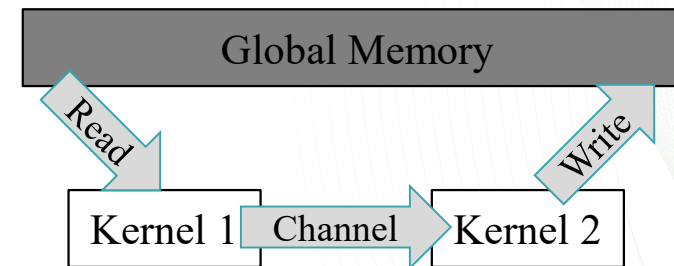
# FPGA OpenCL Architecture

# 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
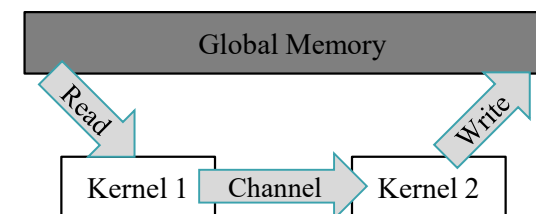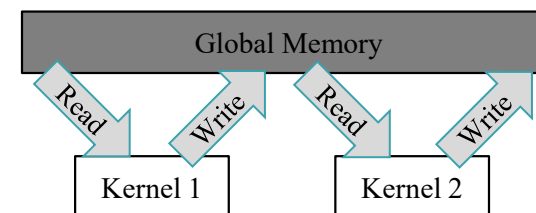
OAK RIDGE
National Laboratory

# 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]; }
}
```



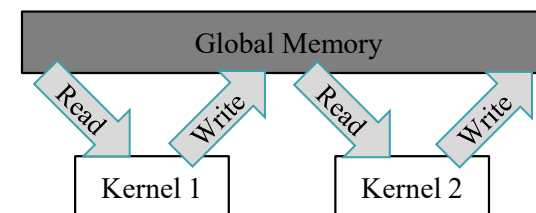(b) Altera OpenCL code with channels

```
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]; }
}
```
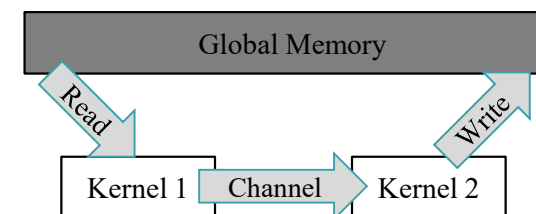


Kernel-pipelining transformation

Valid under specific conditions

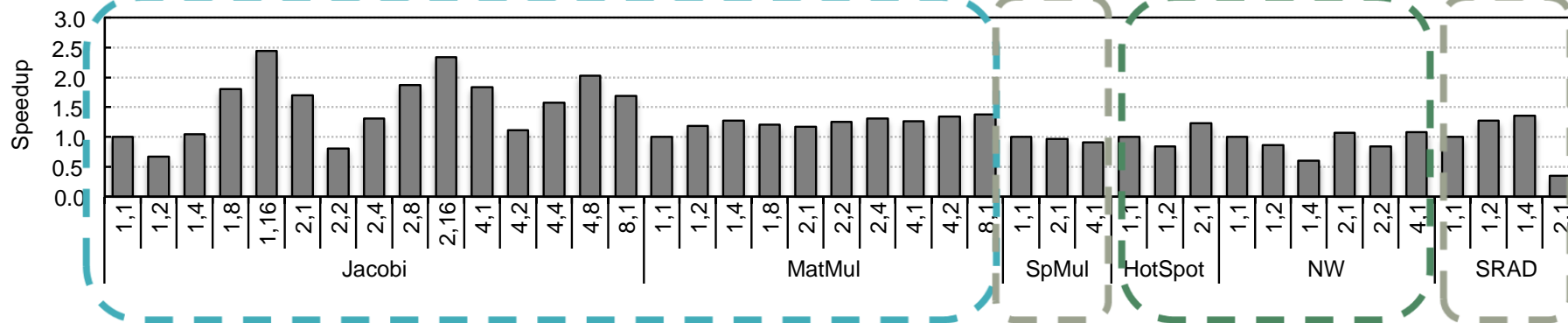(c) Modified OpenACC code for kernel-pipelining

```
#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];}
}
```
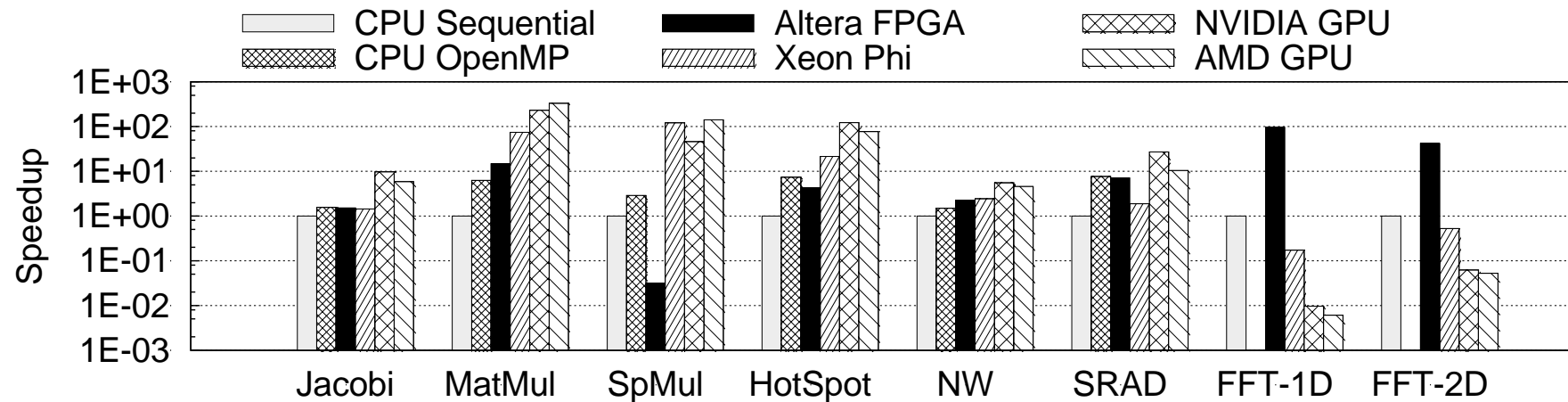
OAK RIDGE
National Laboratory

# Speedup over CU, SIMD (1,1)

# 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.

OAK RIDGE
National Laboratory

# Emerging Memory Systems

# Memory Systems Started Diversifying Several Years Ago
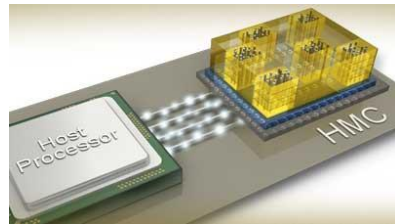
- ## Architectures
  - HMC, HBM/2/3, LPDDR4, GDDR5X, WIDEIO2, etc
  - 2.5D, 3D Stacking

- ## Configurations
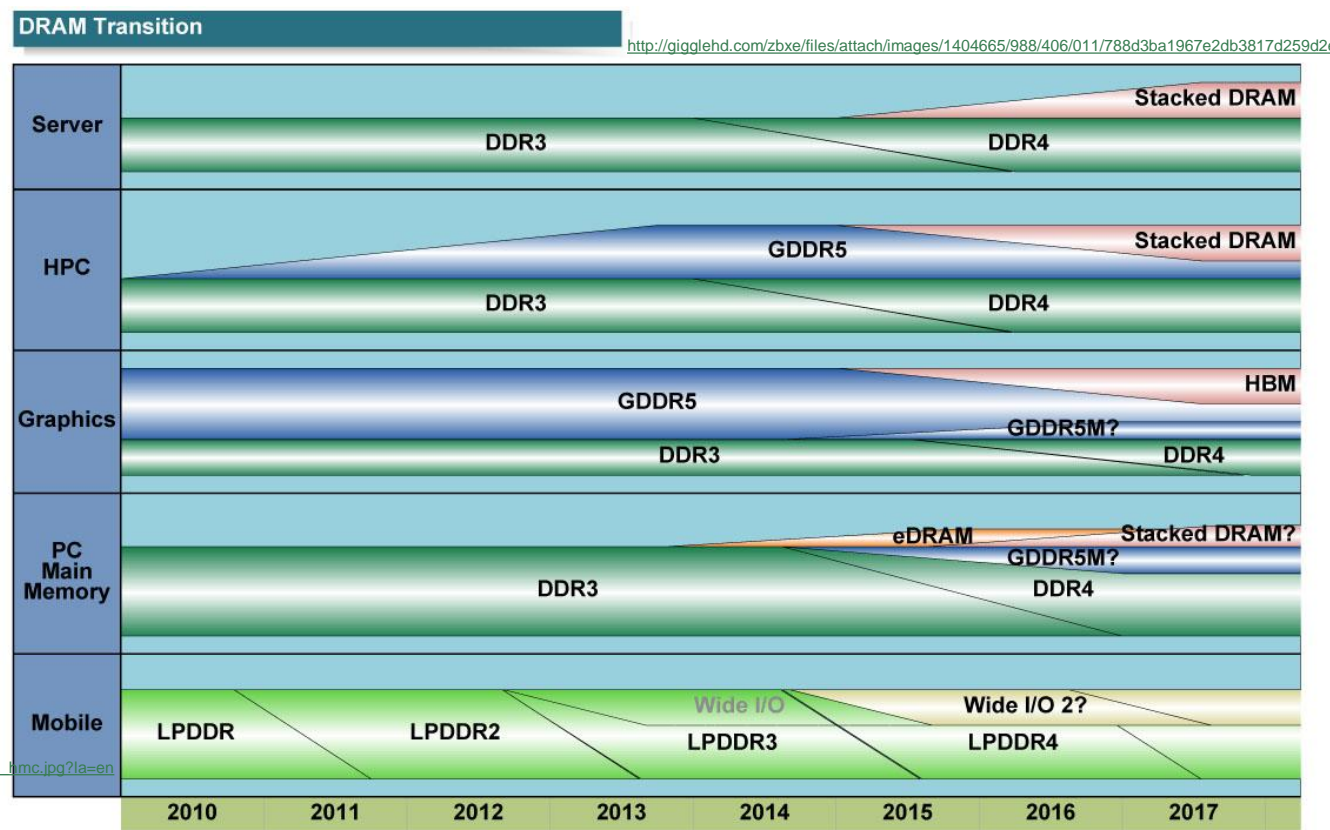  - Unified memory
  - Scratchpads
  - Write through, write back, etc
  - Consistency and coherence protocols
  - Virtual v. Physical, paging strategies

- ## New devices
  - ReRAM, PCRAM, STT-MRAM, 3D-Xpoint

https://www.micron.com/~/media/track-2-images/content-images/content_image_hmc.jpg?la=en

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

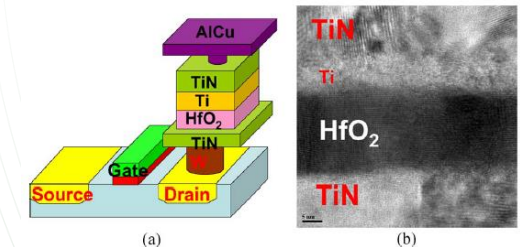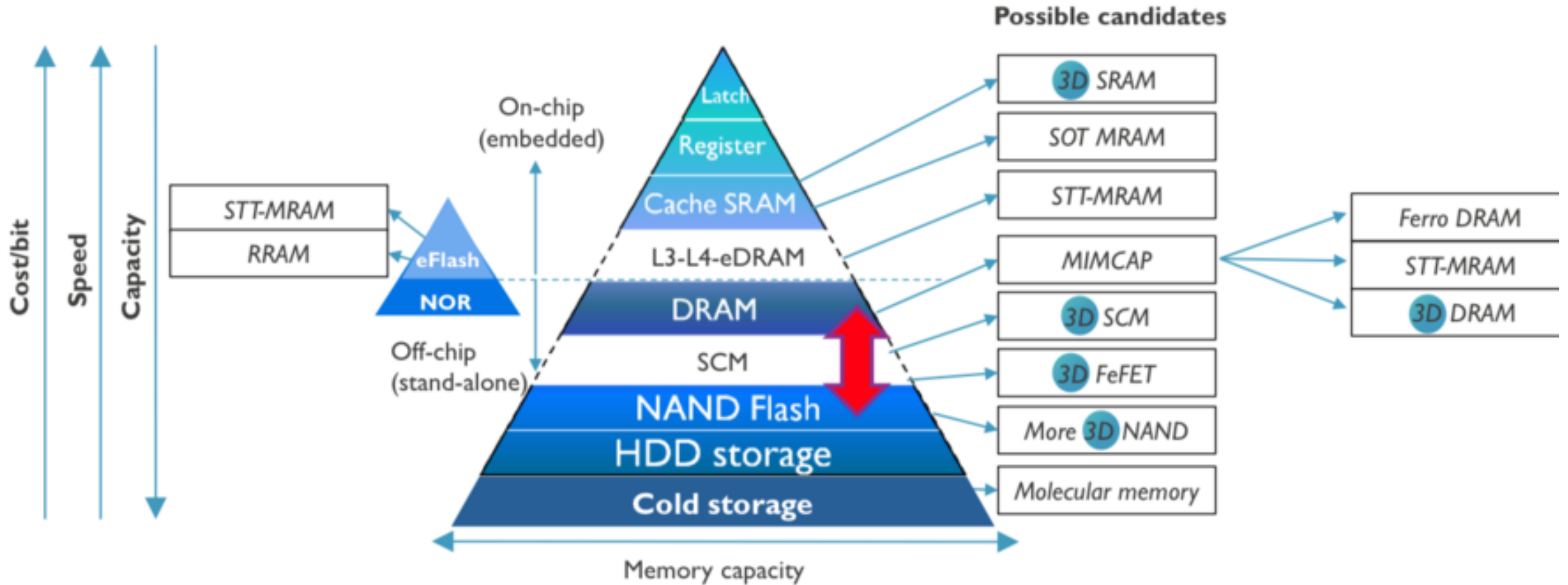Fig. 4. (a) A typical 1T1R structure of RRAM with HfO$_x$; (b) HR-TEM image of the TiN/Ti/HfO$_x$/TiN stacked layer; the thickness of the HfO$_2$ is 20 nm.

H.S.P. Wong, H.Y. Lee, S. Yu et al., "Metal-oxide RRAM," Proceedings of the IEEE, 100(6):1951-70, 2012.

# Complexity in the Expanding and Diversifying Memory Hierarchy



Image Source: IMEC

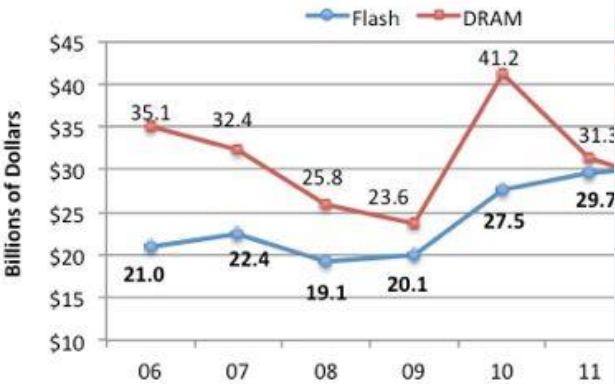# NVRAM Technology Continues to Improve – Driven by Broad Market Forces

# 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 μs |
| Z-NAND | ✗ | 800 | 3.2 | 0.5x | 750K IOPS | 12-20 μs |
| NAND Flash | ✗ | >1,000 | <3 | 0.1x | 50K IOPS | 25-125 μs |

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

RIDGE Laboratory

# Programming NVM Systems Portably

# 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…

OAK RIDGE
National Laboratory

# 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: 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**

**OAK RIDGE** National Laboratory

# DRAGON Operations: Key Components



- **Three memory spaces:**
  - **GPU Mem (GM)** as 1$^{st}$ level cache
  - **Host Mem (HM)** as 2$^{nd}$ 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

https://github.com/pakmarkthub/dragon

P. Markthub, M.E. Belviranli *et al.*, "DRAGON: Breaking GPU Memory Capacity Limits with Direct NVM Access," in SC18, 2018

307

# 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.

```c
#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



NVL-C · · · Other NVL Languages
OpenARC / Other Compiler Front Ends
ARES HLIR
ARES LLVM Passes
LLVM IR + Metadata, Intrinsics, Run-time calls
NVL Runtime
libnvlrt-pmemobj
libpmemobj
LLVM
Target Objects
system linker
Target Executable

OAK RIDGE National Laboratory

# Design Goals: Familiar programming interface

```c
#include <nvl.h>
struct list {
  int value;
  nvl struct list *next;
};
void add(int k, nvl struct list *after) {
  nvl struct list *node
    = nvl_alloc_nv(heap, 1, struct list);
  node->value = k;
  node->next  = after->next;
  after->next = node;
}
```

- Small set of C language extensions:
  - Header file
  - Type qualifiers
  - Library API
  - Pragmas
- Existing memory interfaces remain:
  - NVL-C is a superset of C
  - Unqualified types as specified by C
  - Local/global variables stored in volatile memory (DRAM or registers)
  - Use existing C standard libraries for HDD

OAK RIDGE
National Laboratory

# 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

OAK RIDGE
National Laboratory

http://ft.ornl.gov/research/openarc

# 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

**OAK RIDGE**
National Laboratory

# Programming Model: NVM memory management

- Hybrid of traditional HDD and DRAM programming interfaces

- NVM storage organized into *NVM heaps* identified by file names

- NVM heaps can be managed using normal file system commands

- Within an NVM heap, memory always allocated dynamically

| NVM | HDD analogue |
|---|---|
| nvl_heap_t | FILE |
| nvl_open | fopen |
| nvl_close | fclose |
| mv, rm, ls, *etc.* | mv, rm, ls, *etc.* |

| NVM | DRAM analogue |
|---|---|
| nvl T* | T* |
| nvl_alloc_nv | malloc |
| *automatic* | free |

OAK RIDGE
National Laboratory

# Programming Model: Accessing NVM

Volatile Memory
(registers, stack, bss,
heap)

heap → nvl_heap_t

root

```
nvl_heap_t *heap =
  nvl_open("A.nvl");
```

NVM Heap A
("A.nvl")

How do we access allocations
within an NVM heap?

```
nvl T *root =
  nvl_get_root(heap, T);
```

Checksum error if `T` is
incorrect type.

Set root with `nvl_set_root`.

Before first `nvl_set_root`,
`nvl_get_root` returns null.

OAK RIDGE
National Laboratory

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



**Volatile Memory (registers, stack, bss, heap)**

V-to-NV

**NVM Heap A ("A.nvl")**

run-time error

inter-heap NV-to-NV

**NVM Heap B ("B.nvl")**

NV-to-V

compile-time error

intra-heap NV-to-NV

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

OAK RIDGE
National Laboratory

# Programming Model: Transactions: MATMUL Example

```c
#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

OAK RIDGE
National Laboratory

# Programming Model: Transactions: MATMUL Example

```c
#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

**OAK RIDGE** National Laboratory

# 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

OAK RIDGE
National Laboratory

# 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

# NVM Implications

# Implications

1. Device and architecture trends will have major impacts on HPC in coming decade
   1. NVM in HPC systems is real!
   2. Entirely possible to have an Exabyte of NVM in upcoming systems!
2. Performance trends of system components will create new opportunities and challenges
   1. Winners and losers
3. Sea of NVM allows/requires applications to operate differently
   1. Sea of NVM will permit applications to run for weeks without doing I/O to external storage system
   2. Applications will simply access local/remote NVM
   3. Longer term productive I/O will be 'occasionally' written to Lustre, GPFS
   4. Checkpointing (as we know it) will disappear
4. Requirements for system design will change
   1. Increase in byte-addressable memory-like message sizes and frequencies
   2. Reduced traditional IO demands
   3. KV traffic could have considerable impact – need more applications evidence
   4. Need changes to the operational mode of the system

OAK RIDGE
National Laboratory

# 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

**OAK RIDGE**
National Laboratory

# Acknowledgements

- Contributors and Sponsors
  - Future Technologies Group: http://ft.ornl.gov
  - US Department of Energy Office of Science
    - Exascale Computing Project
    - DOE Vancouver Project: https://ft.ornl.gov/trac/vancouver
    - DOE Blackcomb Project: https://ft.ornl.gov/trac/blackcomb
    - SciDAC RAPIDS Project
  - US DARPA

OAK RIDGE
National Laboratory

# Bonus Material