# Preparing for Extreme Heterogeneity in High Performance Computing

**Jeffrey S. Vetter**

*With many contributions from FTG Group and Colleagues*

18th Annual Workshop on Charm++ and its Applications
October 20, 2020
University of Illinois at Urbana-Champaign (Virtual)

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

**U.S. DEPARTMENT OF ENERGY**

http://ft.ornl.gov    vetter@computer.org

# Highlights

## Recent trends in computing paint an ambiguous future for architectures

- 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

- **Entering an era of Extreme Heterogeneity**
- **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.

- **This is a crisis!**

## Programming systems must provide performance portability (beyond functional portability)!!

- Strive for 'Write once, run anywhere'
- Descriptive models of parallelism and data movement
- Introspective runtime systems
- Layered, modular, open source approaches required

- **Examples**
  - **ECP investments in LLVM**
    - **FORTRAN with GPU offloading**
  - **Programming FPGAs**
    - **Without Verilog**
  - **Memory systems are changing too**
    - **Language support for NVM**

**OAK RIDGE**
National Laboratory

Time for a short poll...

OAK RIDGE
National Laboratory

# History (circa 2010)

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

Yes

No

Revisionists ☺

**OAK RIDGE**
National Laboratory

# Future (circa 2030)

Q: Think forward 10 years. How many of you predict that our top 100 HPC systems will have the following architectural features?

Assume general purpose multicore CPU

- GPU
- FPGA/Reconfigurable processor
- Neuromorphic processor
- Deep learning processor
- Quantum processor
- RISC-V processor
- Some new unknown processor
- All/some of the above in one SoC

OAK RIDGE
National Laboratory

# Implications for Science Applications Teams

Q: Now, imagine you are building a new application with an expected ~3M LOC and 20 team members over the next 10 years. What on-node programming model/system do you use to future-proof your app?

Assume C and C++ (?)

Fortran XX

Metaprogramming, DSEL, etc (e.g., AMP, Kokkos, RAJA, SYCL)

CUDA, cu***, HIP, OpenCL

Directives: OpenMP, OpenACC

Python, Julia, Rust, R, Matlab, etc

Domain Specific Language (e.g., Claw, Hallide, PySL) or Domain Specific Framework (e.g., PetSc, AMReX)

Legion, Charm++, HPX, etc

Some new unknown programming approach

Some combination of the above

**OAK RIDGE**
National Laboratory

# Motivating Trends

# Foundries' Sales Show Hard Times Continuing

**Peter Clarke**

5/23/2016 09:33 P
2 comments

Like 6

LONDON--Taiwan
semiconductor se
market slow down

Both companies a
those they achieve
an annual basis in
TSMC and UMC a

eetasia.com

## GlobalFoundries Forfeit 7nm Manufacturing - EE Times Asia

6-7 min

SAN
the bl

Globa
than 5

subsic

---

**SEMICONDUCTOR** ENGINEERING

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

**MANUFACTURING, DESIGN & TEST**

### Uncertainty Grows For 5nm, 3nm

f 797  G+  74

*Nanosheets and nanowire FETs under development, but costs are skyrocketing. New packaging options could provide an alternative.*

DECEMBER 19TH, 2016 - BY: MARK LAPEDUS

**A**s several chipmakers ramp up their 10nm finFET processes, with 7nm just around the corner, R&D has begun for 5nm and beyond. In fact, some are already moving full speed ahead in the arena.

TSMC recently announced plans to build a new fab in Taiwan

---

## Samsung to Invest $115 Billion in Foundry & Chip Businesses by 2030

---

# Intel's 10nm Is Broken, Delayed Until 2019

**37 COMMENTS**

by Paul Alcorn April 26, 2018 at 6:30 PM

---

**DESIGNLINES** | WIRELESS AND NETWORKING DESIGNLINE

# GlobalFoundries Selling ASIC Business to Marvell

By Dylan McGrath, 05.20.19  1

### Another Step Toward the End of Moore's Law

Samsung and TSMC move to

---

## US Government's Aurora Supercomputer Delayed Due to Intel's 7nm Setback

By Anton Shilov 18 days ago

Intel-powered exascale supercomputer will come online later than expected.

f  Comments (7)

---

13 Dec 2019 | 20:20 GMT

## TSMC's 5-Nanometer Process on Track for First Half of 2020

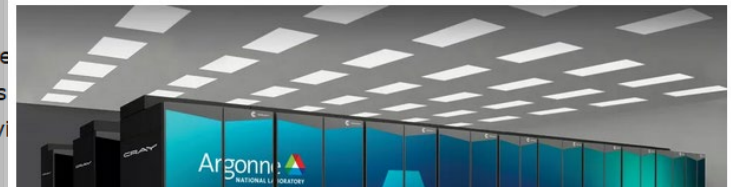Devices are 15 percent faster, 30 percent more energy efficient

By Samuel K. Moore

Photo: Taiwan Semiconductor Manufacturing Co.

The performance enhancement achieved by TSMC's new 5-nanometer process is partly due to the inclusion of a "high-mobility channel." How is it created? TSMC wouldn't reveal.

---

**Number of Foundries with a Cutting Edge Logic Fab**

| | | | | | | | |
|---|---|---|---|---|---|---|---|
| SilTerra | | | | | | | |
| X-FAB | | | | | | | |
| Dongbu HiTek | | | | | | | |
| ADI | ADI | | | | | | |
| Atmel | Atmel | | | | | | |
| Rohm | Rohm | | | | | | |
| Sanyo | Sanyo | | | | | | |
| Mitsubishi | Mitsubishi | | | | | | |
| ON | ON | | | | | | |
| Hitachi | Hitachi | | | | | | |
| Cypress | Cypress | Cypress | | | | | |
| Sony | Sony | Sony | | | | | |
| Infineon | Infineon | Infineon | | | | | |
| Sharp | Sharp | Sharp | | | | | |
| Freescale | Freescale | Freescale | | | | | |
| Renesas (NEC) | Renesas | Renesas | Renesas | Renesas | | | |
| SMIC | SMIC | SMIC | SMIC | SMIC | | | |
| Toshiba | Toshiba | Toshiba | Toshiba | Toshiba | | | |
| Fujitsu | Fujitsu | Fujitsu | Fujitsu | Fujitsu | | | |
| | | | TI | TI | | | |
| | | Panasonic | Panasonic | Panasonic | | | |
| | | STM | STM | STM | | | |
| | | UMC | UMC | UMC | | | |
| | | IBM | IBM | IBM | IBM | | |
| Foundries | GF | GF | GF | GF | GF | | |
| Samsung | Samsung | Samsung | Samsung | Samsung | Samsung | Samsung | |
| TSMC | TSMC | TSMC | TSMC | TSMC | TSMC | TSMC | |
| Intel | Intel | Intel | Intel | Intel | Intel | Intel | Future |
| nm | 45 nm/40 nm | 32 nm/28 nm | 22 nm/20 nm | 16 nm/14 nm | 10 nm | 7 nm | 5 nm |

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

# Sixth Wave of Computing



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

OAK RIDGE
National Laboratory

# 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 effectively for specific workloads
- Integrate components to boost performance and eliminate inefficiencies
- Workload specific memory+storage system design

## Emerging Technologies

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

OAK RIDGE
National Laboratory

# 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 effectively for specific workloads
- Integrate components to boost performance and eliminate inefficiencies
- Workload specific memory+storage system design
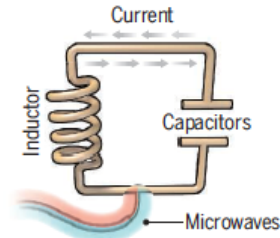
## Emerging Technologies

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

OAK RIDGE
National Laboratory

# Predictions for Transition Period

| 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 effectively for specific workloads<br>• Integrate components to boost performance and eliminate inefficiencies<br>• Workload specific memory+storage system design | • Investigate new computational paradigms<br>  • Quantum<br>  • Neuromorphic<br>  • Advanced Digital<br>  • Emerging Memory Devices |

OAK RIDGE
National Laboratory

# Quantum computing: Qubit design and fabrication have made recent progress but still face challenges

*Science 354, 1091 (2016) – 2 December*



## A bit of the action

In the race to build a quantum computer, companies are pursuing many types of quantum bits, or qubits, each with its own strengths and weaknesses.

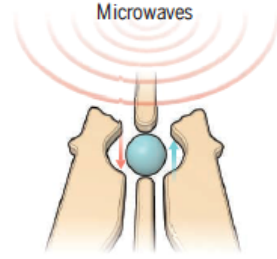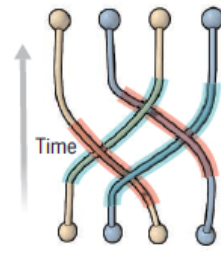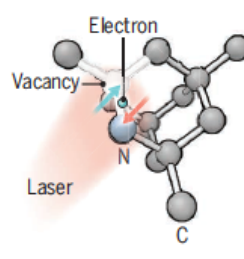| | **Superconducting loops** | **Trapped ions** | **Silicon quantum dots** | **Topological qubits** | **Diamond vacancies** |
|---|---|---|---|---|---|
| | A resistance-free current oscillates back and forth around a circuit loop. An injected microwave signal excites the current into superposition states. | Electrically charged atoms, or ions, have quantum energies that depend on the location of electrons. Tuned lasers cool and trap the ions, and put them in superposition states. | These "artificial atoms" are made by adding an electron to a small piece of pure silicon. Microwaves control the electron's quantum state. | Quasiparticles can be seen in the behavior of electrons channeled through semiconductor structures. Their braided paths can encode quantum information. | A nitrogen atom and a vacancy add an electron to a diamond lattice. Its quantum spin state, along with those of nearby carbon nuclei, can be controlled with light. |
| **Longevity (seconds)** | 0.00005 | >1000 | 0.03 | N/A | 10 |
| **Logic success rate** | 99.4% | 99.9% | ~99% | N/A | 99.2% |
| **Number entangled** | 9 | 14 | 2 | N/A | 6 |
| **Company support** | Google, IBM, Quantum Circuits | ionQ | Intel | Microsoft, Bell Labs | Quantum Diamond Technologies |
| ⊕ **Pros** | Fast working. Build on existing semiconductor industry. | Very stable. Highest achieved gate fidelities. | Stable. Build on existing semiconductor industry. | Greatly reduce errors. | Can operate at room temperature. |
| ⊖ **Cons** | Collapse easily and must be kept cold. | Slow operation. Many lasers are needed. | Only a few entangled. Must be kept cold. | Existence not yet confirmed. | Difficult to entangle. |

**Note:** Longevity is the record coherence time for a single qubit superposition state, logic success rate is the highest reported gate fidelity for logic operations on two qubits, and number entangled is the maximum number of qubits entangled and capable of performing two-qubit operations.

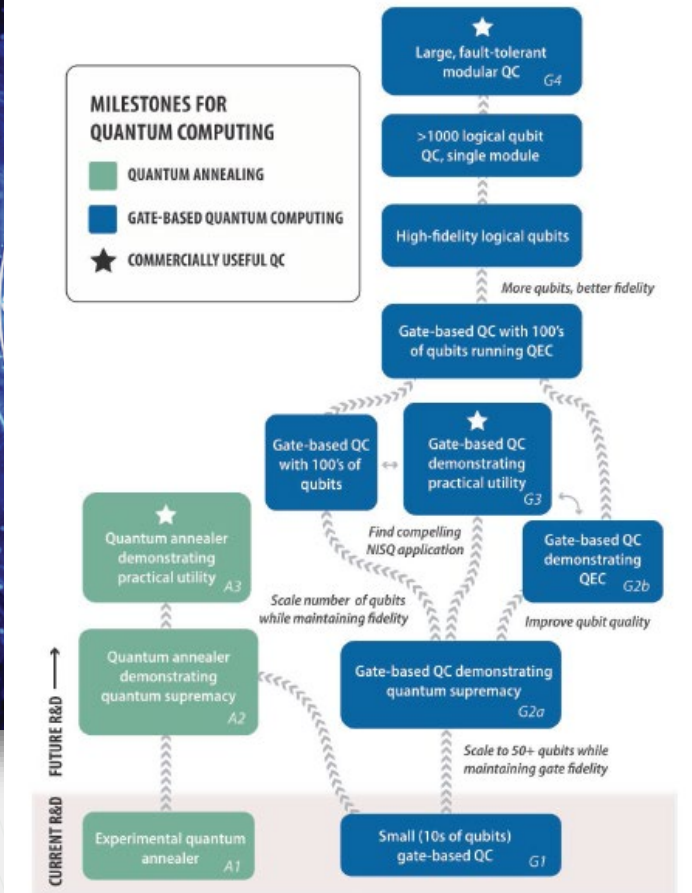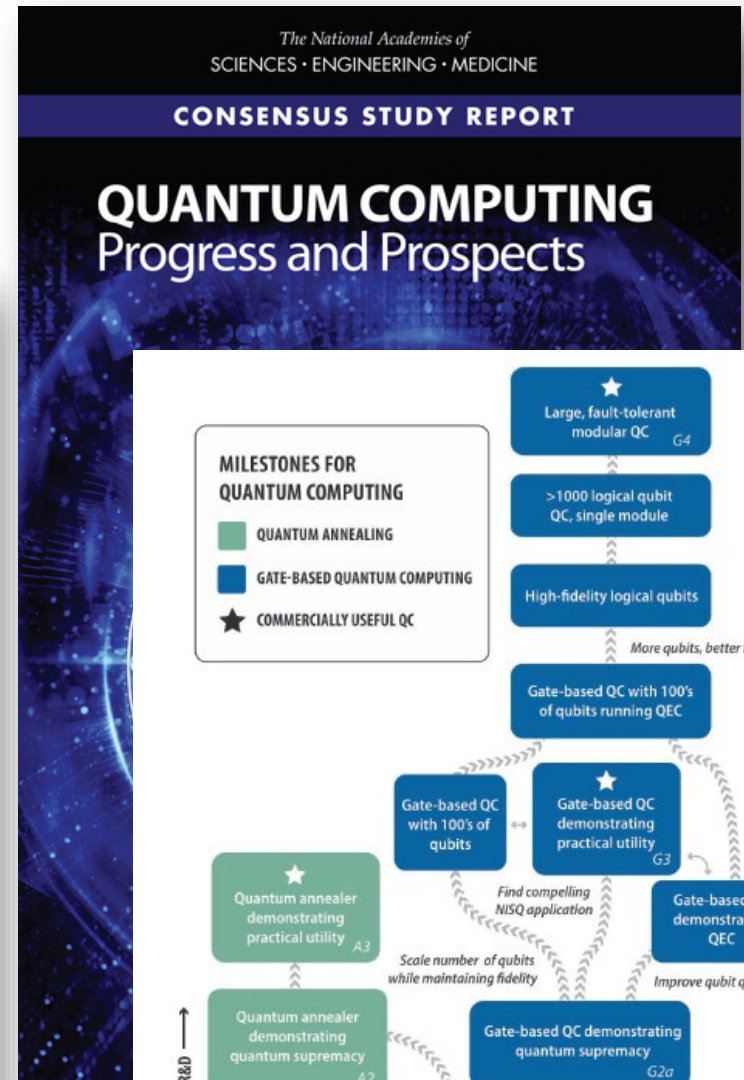## QUANTUM COMPUTING
### Progress and Prospects



FIGURE 7.4 An illustration of potential milestones of progress in quantum computing. The arrangement of milestones corresponds to the order in which the committee thinks they are likely to be achieved; however, it is possible that some will not be achieved, or that they will not be achieved in the order indicated.

http://nap.edu/25196

# Fun Question: when was the field effect transistor patented?

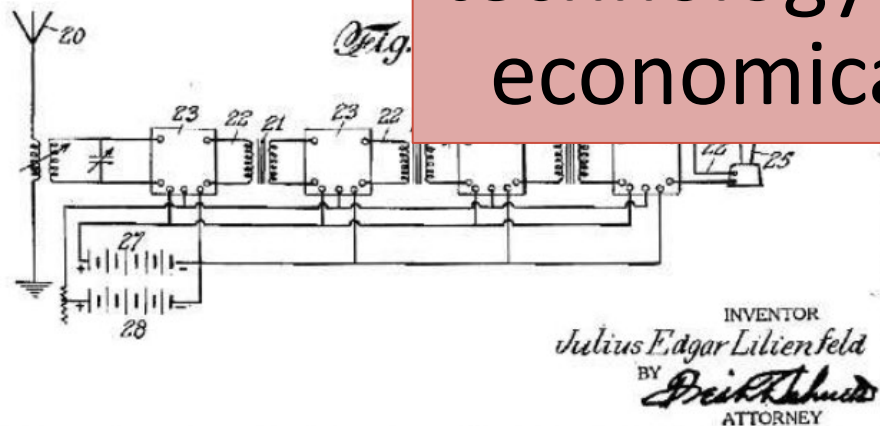## Lilienfeld patents field effect transistor, October 8, 1926

Jessica MacNeil - October 08, 2018

**6 Comments**

On this day in tech history, JE Lilienfeld filed a patent for a three-electrode structure using copper-sulfide semiconductor material, known today as a field-effect transistor.

Lilienfeld's patent for a **"method and apparatu[s] electric currents"** was granted on January 28, 1[...]

According to the patent, his invention was for [...] flow of electric current between two terminals [...] conducting solid by establishing a third potent[...] amplification of oscillating currents like those [...]

INVENTOR
*Julius Edgar Lilienfeld*
BY
ATTORNEY

**Google Patents**    lilienfeld controlling electric curre    🔍    1 of 25 ‹ ›

← Back to results    ✎ controlling; electric; currents; Assignee: lilienfeld;

### Method and apparatus for controlling electric currents

Images (1)

**US1745175A**
United States

📄 Download PDF    🔍 Find Prior Art
Σ Similar

**Inventor:** Lilienfeld Julius Edgar

**Worldwide applications**

1925 CA  1926 US

**Application US140363A events** ⓘ

| Date | Event |
|---|---|
| 1925-10-22 | Priority to CA272437T |
| 1926-10-08 | Application filed by Lilienfeld Julius Edgar |
| 1930-01-28 | Application granted |
| 1930-01-28 | Publication of US1745175A |
| 1947-01-28 | Anticipated expiration |
| 2020-02-16 | Application status is Expired - Lifetime |

junction depletion layer or carrier concentration layer; Details of semiconductor bodies or of electrodes thereof; Multistep manufacturing processes therefor

🔳 H01L29/78681 Thin film transistors, i.e. transistors with a channel being at least partly a thin film having a semiconductor body comprising AIIIBV or AIIBVI or AIVBVI semiconductor materials, or Se or Te

---

**Moral of this story**
It may take decades for a new technology to be manufacturable, economical, and usable, if ever.

51

OAK RIDGE
National Laboratory

# Predictions for Transition Period

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

OAK RIDGE
National Laboratory

# Various Markets already Experiencing these Architectural Trends

# DOE HPC Roadmap to Exascale Systems

| FY 2012 | FY 2016 | FY 2018 | FY 2021 | FY 2022 | FY 2023 |
|---------|---------|---------|---------|---------|---------|



**Titan**
**ORNL**
Cray/AMD/NVIDIA

**Mira**
**ANL**
IBM BG/Q

decommissioned

**Theta**
**ANL**
Cray/Intel KNL

**Cori**
**LBNL**
Cray/Intel Xeon/KNL

**Summit**
**ORNL**
IBM/NVIDIA

**FRONTIER**
**ORNL**
HPE/AMD

**Aurora**
**ANL**
Intel/HPE

**Perlmutter**
**LBNL**
HPE/AMD/NVIDIA

**Exascale Systems**

**Sequoia**
**LLNL**
IBM BG/Q

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

**Sierra**
**LLNL**
IBM/NVIDIA

**CROSSROADS**
**LANL/SNL**
HPE/Intel

**EL CAPITAN**
**LLNL**
HPE/AMD

ECP EXASCALE COMPUTING PROJECT

Version 1.6
September 30, 2020

65

# Future -> Open Source Hardware Enables a Rapid Design of Specialized Chips and Effectively Mass Customization

## RISC-V Ecosystem

**Software**

| Open-source software: | Commercial software: |
|---|---|
| Gcc, binutils, glibc, Linux, BSD, LLVM, QEMU, FreeRTOS, ZephyrOS, LiteOS, SylixOS, ... | Lauterbach, Segger, Micrium, ExpressLogic, ... |

**RISC-V Foundation** — ISA specification

**Hardware**

| Open-source cores: | Commercial cores: |
|---|---|
| Rocket, BOOM, RI5CY, Ariane, PicoRV32, Piccolo, SCR1, Hummingbird, ... | Andes, Bluespec, Codasip, Cortus, Nuclei, SiFive, Sy... |

### DARPA — 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
Add to Cart
or 1-Click Checkout
Buy now with 1-Click®
Source: Amazon
Large Box
Source: NVIDIA

Distribution Statement "A" (Approved for Public Release, Distribution Unlimited)     23

A. Olofsson, 2018

### Open-source computing

## A new blueprint for microprocessors challenges the industry's giants

*RISC-V is an alternative to proprietary designs*
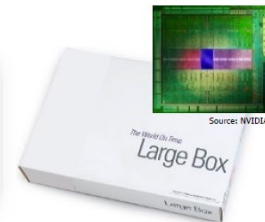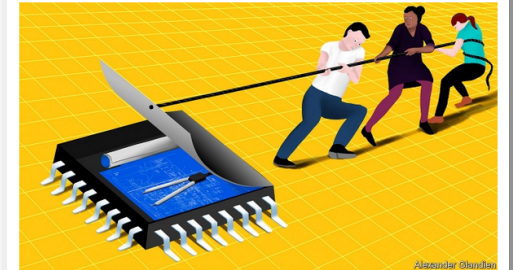
Print edition | Science and technology ›
Oct 3rd 2019

Most microprocessors—the chips that do the grunt work in computers—are built around designs, known as instruction-set architectures (ISAs), which are owned either by Intel, an American giant, or by Arm, a Japanese one. Intel's ISAs power desktop computers, servers and laptops. Arm's power phones, watches and other mobile devices. Together, these two firms dominate the market. Almost every one of the 5.1bn mobile phones on the planet, for example, relies on an Arm-designed ISA. The past year, however, has seen a boomlet in chips made using an ISA called RISC-V. If boomlet becomes boom, it may change the chip industry dramatically, to the detriment of Arm and Intel, because unlike the ISAs from those two firms, which are proprietary, RISC-V is available to anyone, anywhere, and is free.

An ISA is a standardised description of how a chip works at the most basic level, and instructions for writing software to run on it. To draw an analogy, a house might have two floors or three, five bedrooms or six, one bathroom or two. That is up to the architect. An ISA, however, is the equivalent of insisting that the same sorts of electrical sockets and water inlets and outlets be put in the same places in every appropriate room, so that an electrician or a plumber can find them instantly and carry the correct kit to connect to them.

National Laboratory

# Take away message →
# During this Sixth Wave transition, Complexity is our major challenge!

**Architecture**

- How do we design future systems so that they are better than current systems on important applications?
- Simulation and modeling are more difficult
- Entirely possible that the new system will be slower than the old system!
- Expect 'disaster' procurements

**Programmability**

- How do we design applications with some level of performance portability?
- Software lasts much longer than transient hardware platforms
- Proper abstractions for flexibility and efficiency
- Adapt or die

# Implication for Applications and Software → Explosion of (incomplete) Programming Systems →

## Programming Models

- OpenMP
- OpenMP Offload
- OpenACC
- SYCL
- DPC++
- HIP
- CUDA
- OpenCL
- Kokkos
- Raja
- *Many others…*

## Implementations (SYCL Example)



Fig. 1: Current SYCL implementations and their corresponding API backends.

Inconsistent: feature support, performance, tools ecosystem

# Software Strategies for Extreme Heterogeneity

# Strategies for Programming Systems in this Era of Rapidly Designed, Diverse Architectures

## Goals

- Strive for 'Write one, run anywhere'

- Descriptive models of parallelism and data movement that enable effective code generation

- Introspective runtime systems

- Layered, modular, open source approaches required
  - One organization can't do it all

## Examples

- Contributing to LLVM
  - FORTRAN with GPU offloading

- Programming FPGAs
  - Without Verilog

- Memory systems are changing too
  - Language support for NVM

**OAK RIDGE**
National Laboratory

# Contributing to LLVM

# 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) | Software Technology (ST) | Hardware and Integration (HI) |
|---|---|---|
| Develop and enhance the predictive capability of applications critical to the DOE | Produce expanded and vertically integrated software stack to achieve full potential of exascale computing | 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 |
|---|---|---|

EXASCALE COMPUTING PROJECT

# ECP is Improving the LLVM Compiler Ecosystem

## LLVM
- Very popular open source compiler infrastructure
- Easily extensible
- Widely used and contributed to in industry
- Permissive license
- Used for heterogeneous computing

## +SOLLVE
- Enhancing the implementation of OpenMP in LLVM
- Unified memory
- OMP Optimizations
- Prototype OMP features for LLVM
- OMP test suite
- Tracking OMP implementation quality

## +PROTEAS-TUNE
- Core optimization improvements to LLVM
- OpenACC capability for LLVM
  - Clacc
  - Flacc
- Autotuning for OpenACC and OpenMP in LLVM
- Integration with Tau performance tools

## +FLANG
- Developing an open-source, production Fortran frontend
- Upstream to LLVM public release
- Support for OpenMP and OpenACC
- Recently approved by LLVM

## +HPCToolkit
- Improvements to OpenMP profiling interface OMPT
- OMPT specification improvements
- Refine HPCT for OMPT improvements

## +ATDM
- Enhancing LLVM to optimize template expansion for FlexCSI, Kokkos, RAJA, etc.
- Flang testing and evaluation

## Vendors
- Increasing dependence on LLVM
- Collaborations with many vendors using LLVM
  - AMD
  - ARM
  - Cray
  - HPE
  - IBM
  - Intel
  - NVIDIA

Active involvement with broad LLVM community: LLVM Dev, EuroLLVM

https://github.com/llvm-doe-org/llvm-project

# Leveraging LLVM Ecosystem to Meet a Critical ECP (community) need : FORTRAN

- Fortran support continues to be an ongoing requirement

- Flang project started in NNSA funding NVIDIA/PGI to open source compiler front-end into LLVM ecosystem

- SOLLVE is improving OpenMP dialect, implementation, and core optimizations

- PROTEAS-TUNE is creating OpenACC dialect and improving MLIR

- ECP projects are contributing many changes upstream to LLVM core, MLIR, etc

- Many others are contributing: backends for processors, optimizations in toolchain, ...
  - Google contributed MLIR



EXASCALE COMPUTING PROJECT

160

# PROTEAS-TUNE: Clacc – OpenACC in Clang/LLVM

- Develop production-quality, standard-conforming traditional OpenACC compiler and runtime support by extending Clang and LLVM

  - Build on existing OpenMP infrastructure

- Enable research and development of source-level OpenACC tools

  - Design compiler to leverage Clang/LLVM ecosystem extensibility

  - E.g., Pretty printers, analyzers, lint tools, and debugger and editor extensions

- Actively contribute improvements to the OpenACC specification

- Actively contribute upstream all Clang and LLVM improvements that are mutually beneficial

  - Many contributions are already in LLVM

- Open-source with multiple collaborators (vendors, universities)



Clacc: Translating OpenACC to OpenMP in Clang, Joel E. Denny, Seyong Lee, and Jeffrey S. Vetter, 2018 IEEE/ACM 5th Workshop on the LLVM Compiler Infrastructure in HPC (LLVM-HPC), Dallas, TX, USA, (2018).

# Programming FPGAs with OpenACC

OAK RIDGE
National Laboratory

# 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.
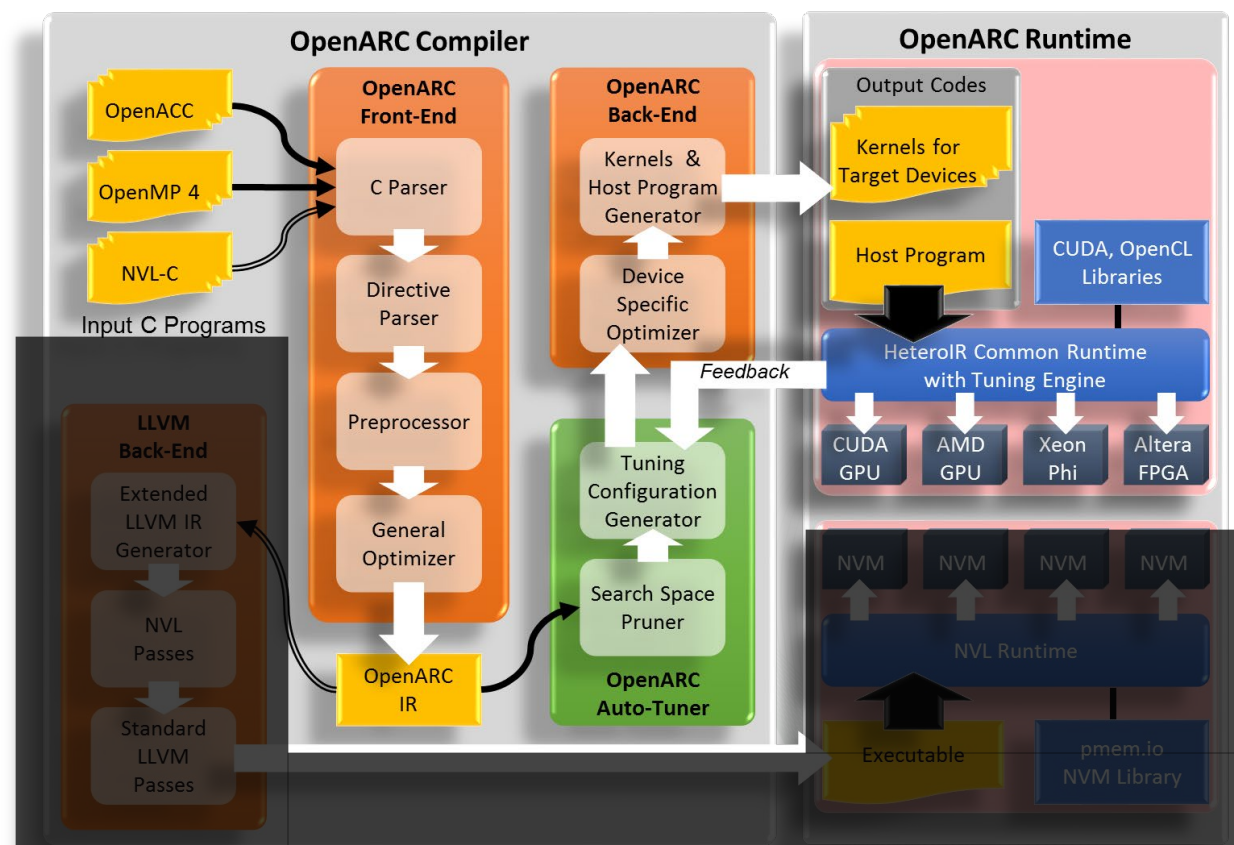    - None of these use a portable, open standard.

# 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.

OAK RIDGE
National Laboratory

# Directive-based Strategy with OpenARC: Open Accelerator Research Compiler

- Open-Sourced, High-Level Intermediate Representation (HIR)-Based, Extensible Compiler Framework.

  - Perform source-to-source translation from OpenACC C to target accelerator models.

    - Support full features of OpenACC V1.0 ( + array reductions and function calls)

    - Support both CUDA and OpenCL as target accelerator models

  - Provide common runtime APIs for various back-ends

  - Can be used as a research framework for various study on directive-based accelerator computing.

    - Built on top of Cetus compiler framework, equipped with various advanced analysis/transformation passes and built-in tuning tools.

    - OpenARC's IR provides an AST-like syntactic view of the source program, easy to understand, access, and transform the input program.
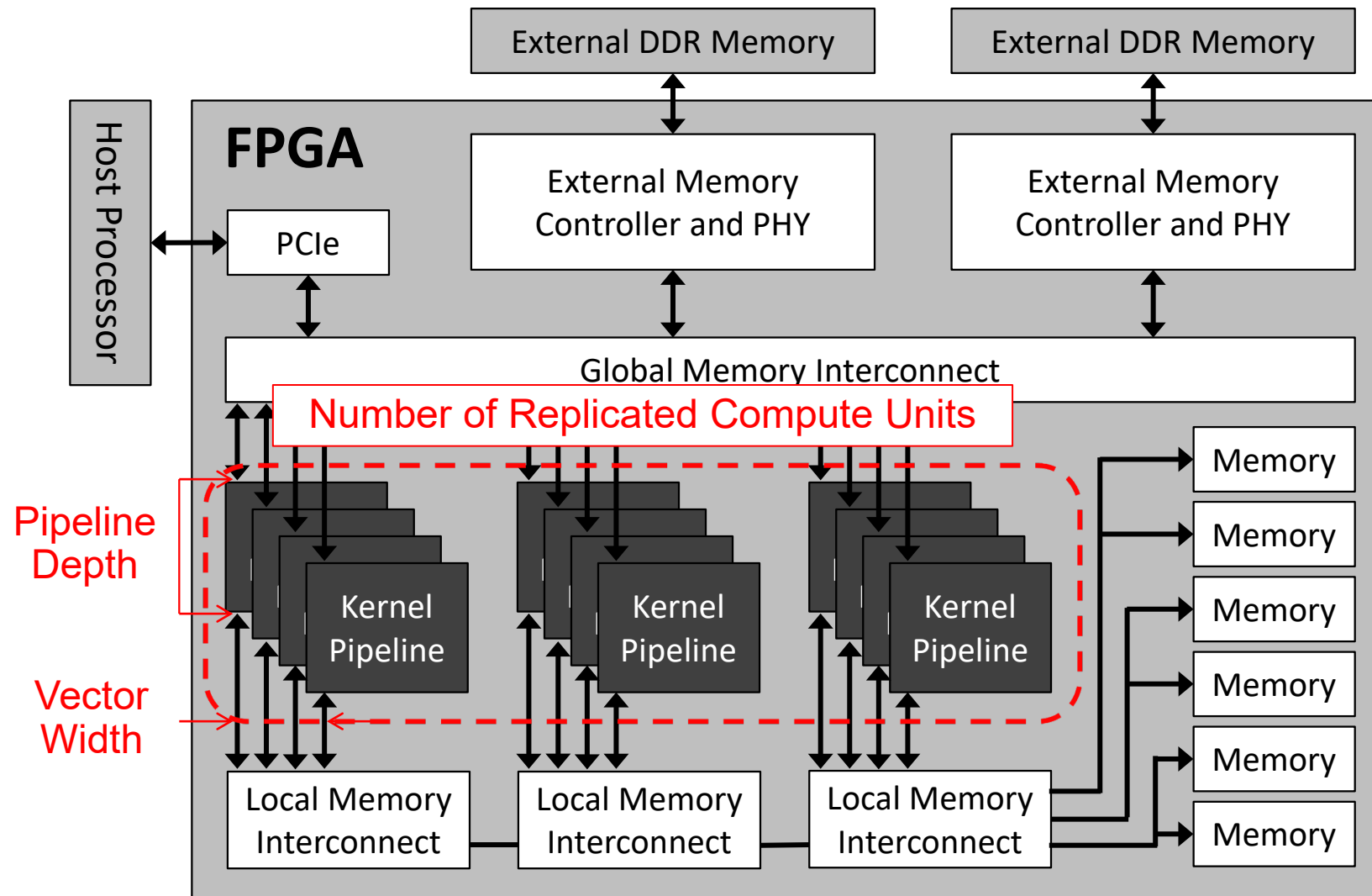
S. Lee and J.S. Vetter, "OpenARC: Open Accelerator Research Compiler for Directive-Based, Efficient Heterogeneous Computing," in *ACM Symposium on High-Performance Parallel and Distributed Computing (HPDC). Vancouver: ACM, 2014*

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).

S. Lee, J. Kim, and J.S. Vetter, "OpenACC to FPGA: A Framework for Directive-based High-Performance Reconfigurable Computing," Proc. IEEE International Parallel & Distributed Processing Symposium (IPDPS), 2016, 10.1109/IPDPS.2016.28.

OAK RIDGE
National Laboratory

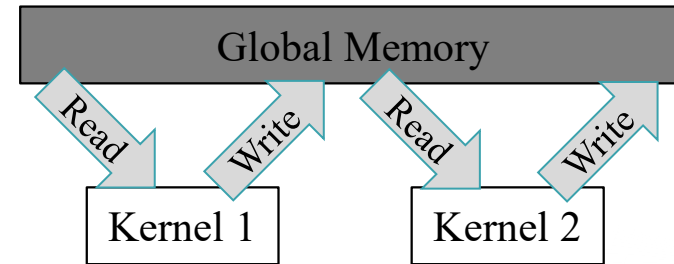# Baseline Translation of OpenACC-to-FPGA

- 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.
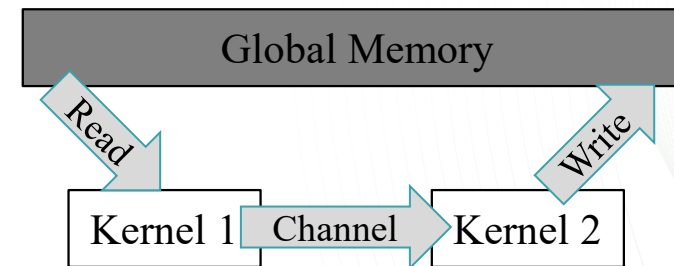  - Reuse most of compiler passes for kernel generation.

# 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

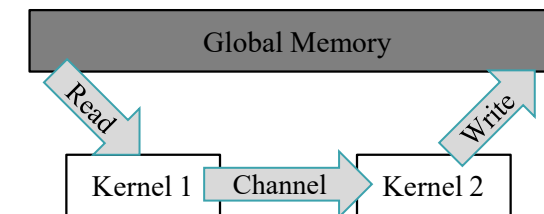Kernel communications through global memory in OpenACC
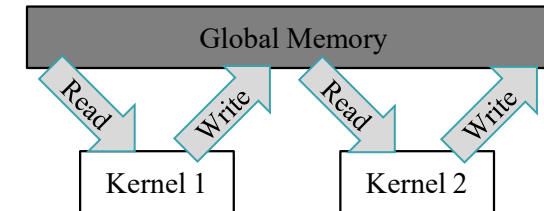
Kernel communications with Altera channels

# 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]; }
}
```



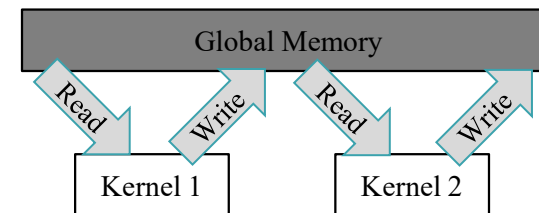(b) Altera OpenCL code with channels

```
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]; }
}
```
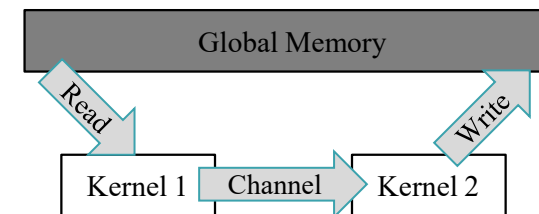


Global Memory

Read   Write   Read   Write

Kernel 1        Kernel 2

Kernel-pipelining transformation
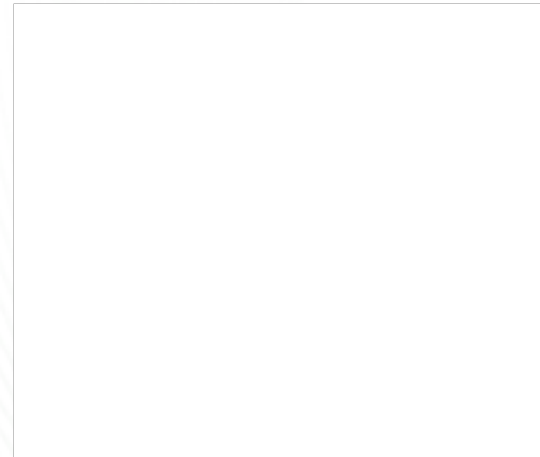
Valid under specific conditions

(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];}
}
```
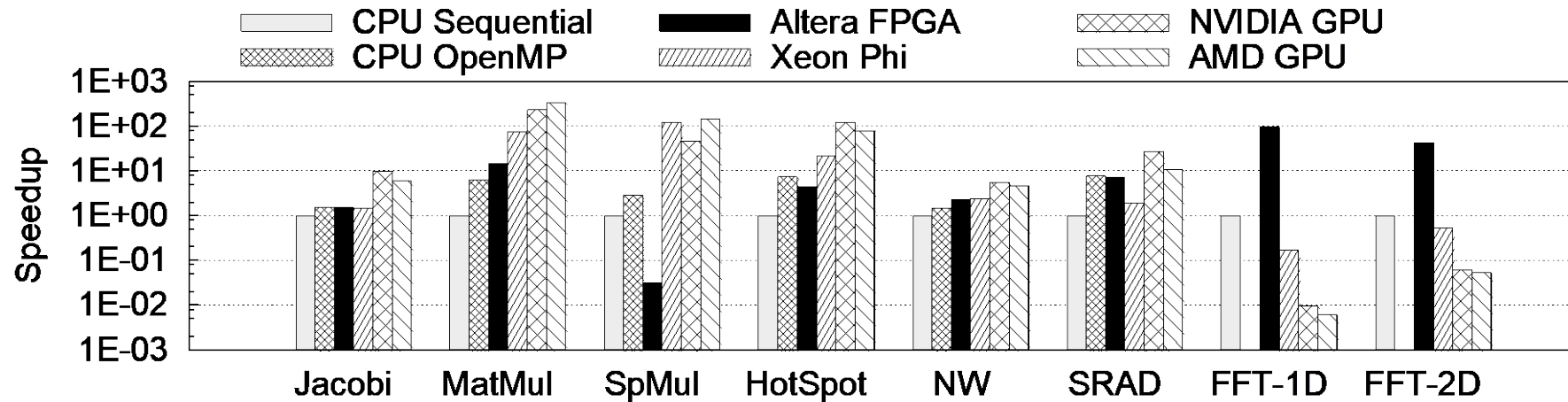


Global Memory

Read                        Write

Kernel 1   Channel   Kernel 2

OAK RIDGE
National Laboratory

# FPGA-specific Optimizations

- Single work-item

- Collapse

- <u>Reduction</u>

- Sliding window

- (Branch-variant code motion)

- (Custom unrolling)

Directive-Based Programming for High-Performance FPGA Computing

OAK RIDGE
National Laboratory

# Overall Performance of OpenARC FPGA Evaluation



Legend: CPU Sequential, CPU OpenMP, Altera FPGA, Xeon Phi, NVIDIA GPU, AMD GPU

Y-axis: Speedup (1E+03 to 1E-03)

X-axis categories: Jacobi, MatMul, SpMul, HotSpot, NW, SRAD, FFT-1D, FFT-2D

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 Hierarchy is Specializing too



Image Source: IMEC

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

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

# 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

# Design Goals: Modular implementation



- Core is common compiler middle-end
- Multiple complier front ends for multiple high-level languages:
  - For now, just OpenARC for NVL-C
- Multiple runtime implementations:
  - For now, just Intel's pmem (pmemobj)

# NVL-C: Programming Model

- 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);
}
```

Denny, J.E., Lee, S., and Vetter, J.S.: 'NVL-C: Static Analysis Techniques for Efficient, Correct Programming of Non-Volatile Main Memory Systems'. Proc. Proceedings of the 25th ACM International Symposium on High-Performance Parallel and Distributed Computing, Kyoto, Japan2016 pp. Pages

**OAK RIDGE**
National Laboratory

# Programming Model: NVM Pointers

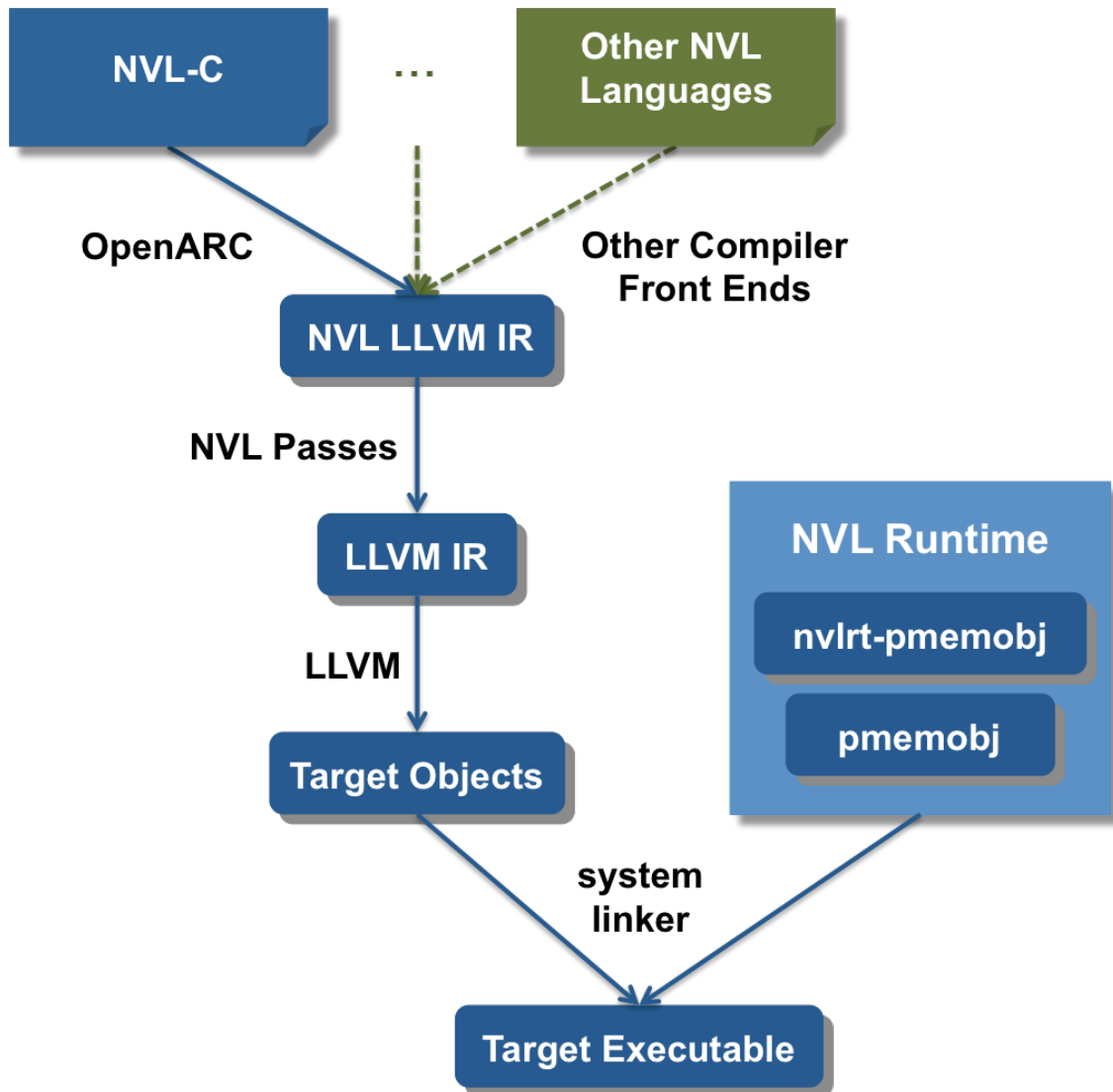```
#include <nvl.h>
struct list {
  int value;
  nvl struct list *next;
};
void add(int k, nvl struct list *after) {
  struct list *node
    = malloc(sizeof(struct list));
  node->value = k;
  node->next  = after->next;
  after->next = node;
}
```

*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: Bare NVM Pointers

- NVM pointers are wide:
  - Facilitates: automatic reference counting, pointer constraints, transactions
  - NVM pointers must be decoded for target architecture's load/store
  - Bare NVM pointer = virtual address with all NVL-C metadata stripped away
- NVM pointer hoisting is important for performance:
  - Conversion to bare NVM pointer is many instructions longer than load/store
  - In tight loop, the performance penalty is severe
  - If conversion is loop-invariant, it can be hoisted
  - Currently, we implement per application with an informal NVL-C extension
  - Future work: eliminate extension and automate in compiler

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.)

# 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: Undo logs

```
#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;
    }
  }
}
```

- Before every NVM store, transaction creates undo log to back up old data

- Undo log contains metadata plus old data being overwritten

- Problem: large overhead because an undo log is created for every element of `a` (every iteration of `j` loop)

OAK RIDGE
National Laboratory

# Programming Model: Transactions: `clobber` clause

```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) \
            clobber(a[*i:1])
    {
      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;
    }
  }
}
```

- **`clobber`** clause suppresses undo logs

- Durability after transaction commit is still guaranteed

OAK RIDGE
National Laboratory

# Evaluation: LULESH



- ExM = use SSD as extended DRAM

- T1 = BSR + transactions

- T2 = T1 + `backup` clauses

- T3 = T1 + `clobber` clauses

- BlockNVM = `msync` included

- ByteNVM = `msync` suppressed

- `backup` **is important for performance**
- `clobber` **cannot be applied because old data is needed**

# 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

OAK RIDGE
National Laboratory

# Recap

**Recent trends in computing paint an ambiguous future for architectures**

- 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

- **Entering an era of Extreme Heterogeneity**
- **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.

- **This is a crisis!**

**Programming systems must provide performance portability (beyond functional portability)!!**

- Strive for 'Write once, run anywhere'
- Descriptive models of parallelism and data movement
- Introspective runtime systems
- Layered, modular, open source approaches required

- **Examples**
  - **ECP investments in LLVM**
    - **FORTRAN with GPU offloading**
  - **Programming FPGAs**
    - **Without Verilog**
  - **Memory systems are changing too**
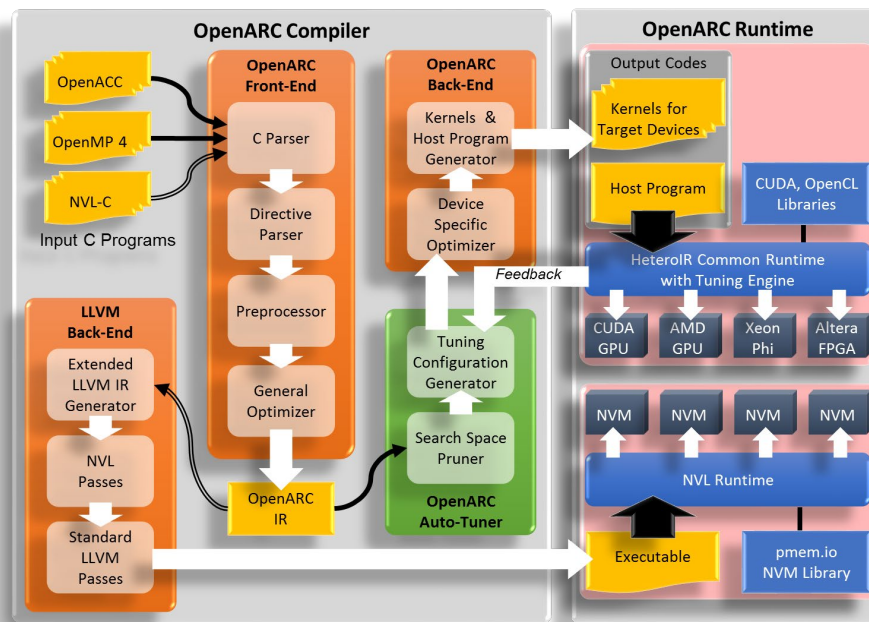    - **Language support for NVM**

- ~~Visit us~~
  - ~~We host interns and other visitors year round~~
    - ~~Faculty, grad, undergrad, high school, industry~~

- Jobs at ORNL
  - Postdoctoral Research Associate in Computer Science
  - Software Engineer
  - Computer Scientist
  - Visit https://jobs.ornl.gov

- Contact me vetter@ornl.gov



**OAK RIDGE**
National Laboratory

# 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