# Programming for the Future : Are We There Yet?
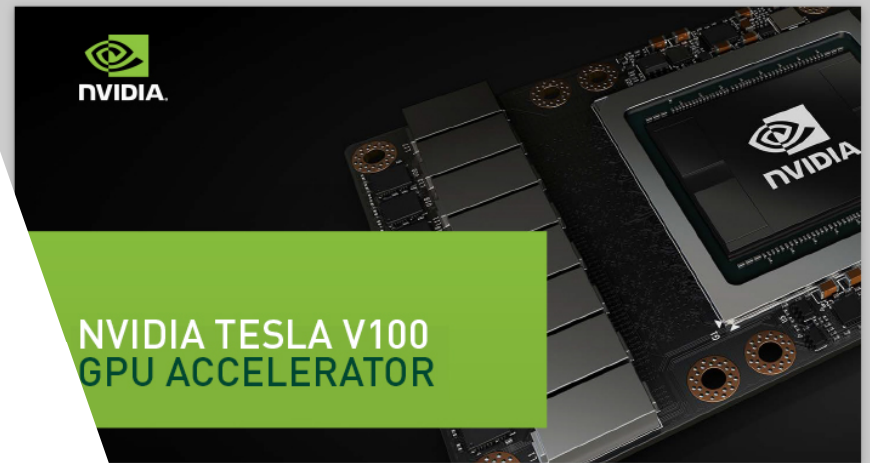
Barbara Chapman
Stony Brook University
Brookhaven National Lab

SPPEXA Workshop, Versailles, March 21, 2019

**iACS** INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# We Live in Interesting Times

Development of large-scale applications is challenging. It is becoming much more so.

- ❑ Increasing complexity of computer hardware
    - ❑ Diversity in large-scale platforms
- ❑ Growing diversity of applications and users
    - ❑ Traditional scientific computations, AI, combination, workflows
- ❑ Increasing need for dynamic program adaptation
    - ❑ To handle changes in computation or resources
- ❑ Changing expectations on part of application developers
    - ❑ Python, TensorFlow, PyTorch, …
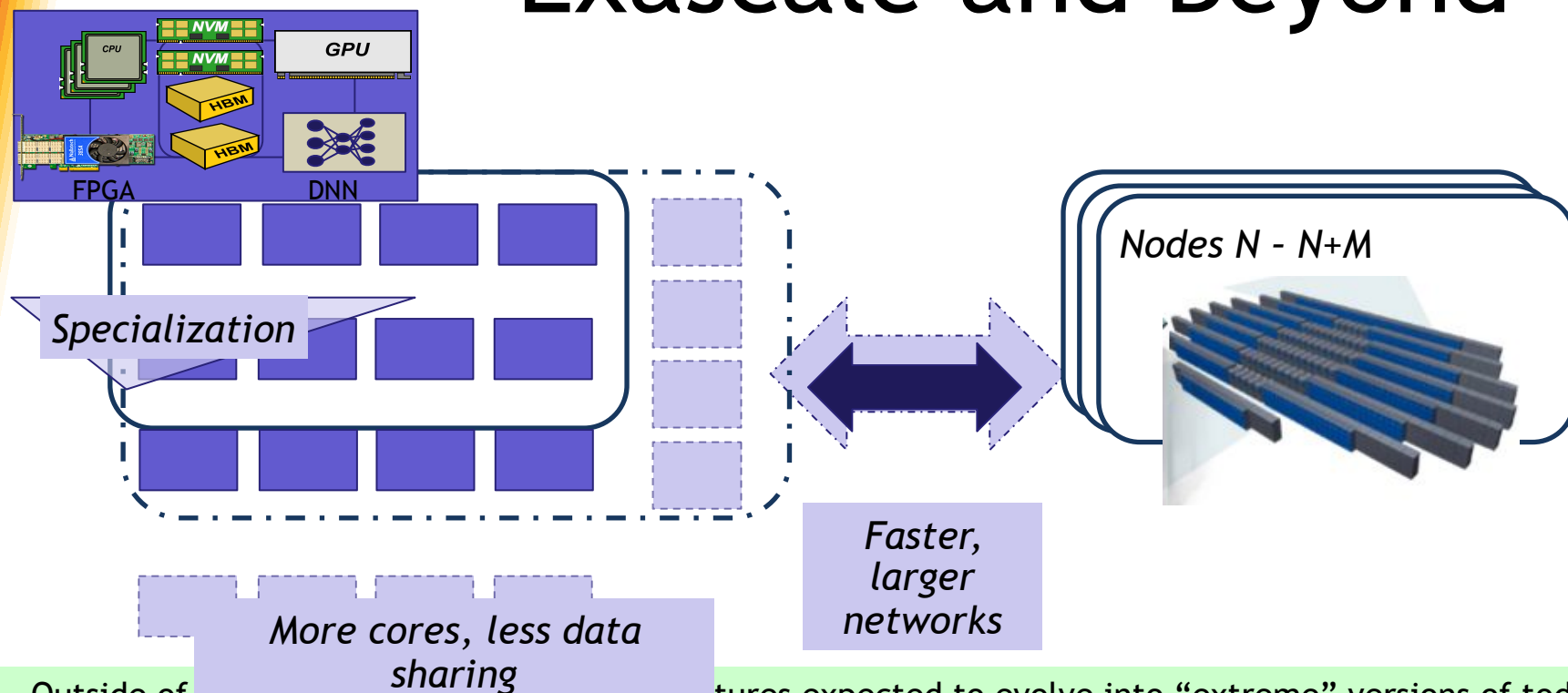- ❑ Scalable performance, performance portability, productivity; power saving



NVIDIA TESLA V100
GPU ACCELERATOR

intel
CORE i9
X-series

Xilinx®
FPGA

iACS INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# Evolution of DOE Leadership Class Systems

Accelerated node

| Name | Titan | Mira | Cori | Theta | Summit | Sierra | Perlmutter |
|---|---|---|---|---|---|---|---|
| System peak (PF) | 27 | 10 | Haswell: 2.81 KNL: 29.5 | 11.69 | 200 | 125 | |
| Peak Power (MW) | 9 | 4.8 | 4.2 | 1.7 | 13.3 | | 6 |
| Total system memory | 710TB | 768 TB | Haswell: 298.5 TB DDR4 KNL: 1.06 PB DDR4 + High Bandwidth Memory | 1475 TB: 843 DDR4 + 70 MCDRAM + 562 SSD | 2.8 PB: DDR4, HBM2, PB persistent, memory | 1.4 PB DDR4, HBM2, PB persistent, memory | |
| Node performance (TF) | 1.452 | 0.204 | Haswell: 1.178 KNL: 3.046 | 2.66 | >40 | | |
| Node Processors | AMD Opteron NVIDIA K20x | 64-bit PowerPC A2 | Intel Haswell Intel KNL | Intel KNL | 2 POWER9 6 NVIDIA Volta GPUs | 2 POWER9 4 NVIDIA Volta GPUs | AMD EPYC (Milan) NVIDIA GPU |
| System Size (nodes) | 18,688 nodes | 49,152 | Haswell; 2,388 nodes KNL: 9,688 nodes | 4,392 nodes | ~4600 nodes | 4320 | > 4000 node CPU-only partition |
| System Interconnect | Gemini | 5D Torus | Aries | Aries | Dual Rail EDR-IB | Dual Rail EDR-IB | Cray Slingshot |
| File System | 32 OB 1 TB/s Lustre | 26 PB 300 GB/s GPFS | 28 PB >700 GB/s Lustre | 10 PB 744 GB/s Lustre | 120 PB 1 TB/s GPFS | | 30 PB 4 TB/s Lustre |

iACS INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# Exascale and Beyond

**FPGA** · CPU · NVM · NVM · HBM · HBM · GPU · **DNN**

*Specialization*

**Nodes N – N+M**

*Faster, larger networks*
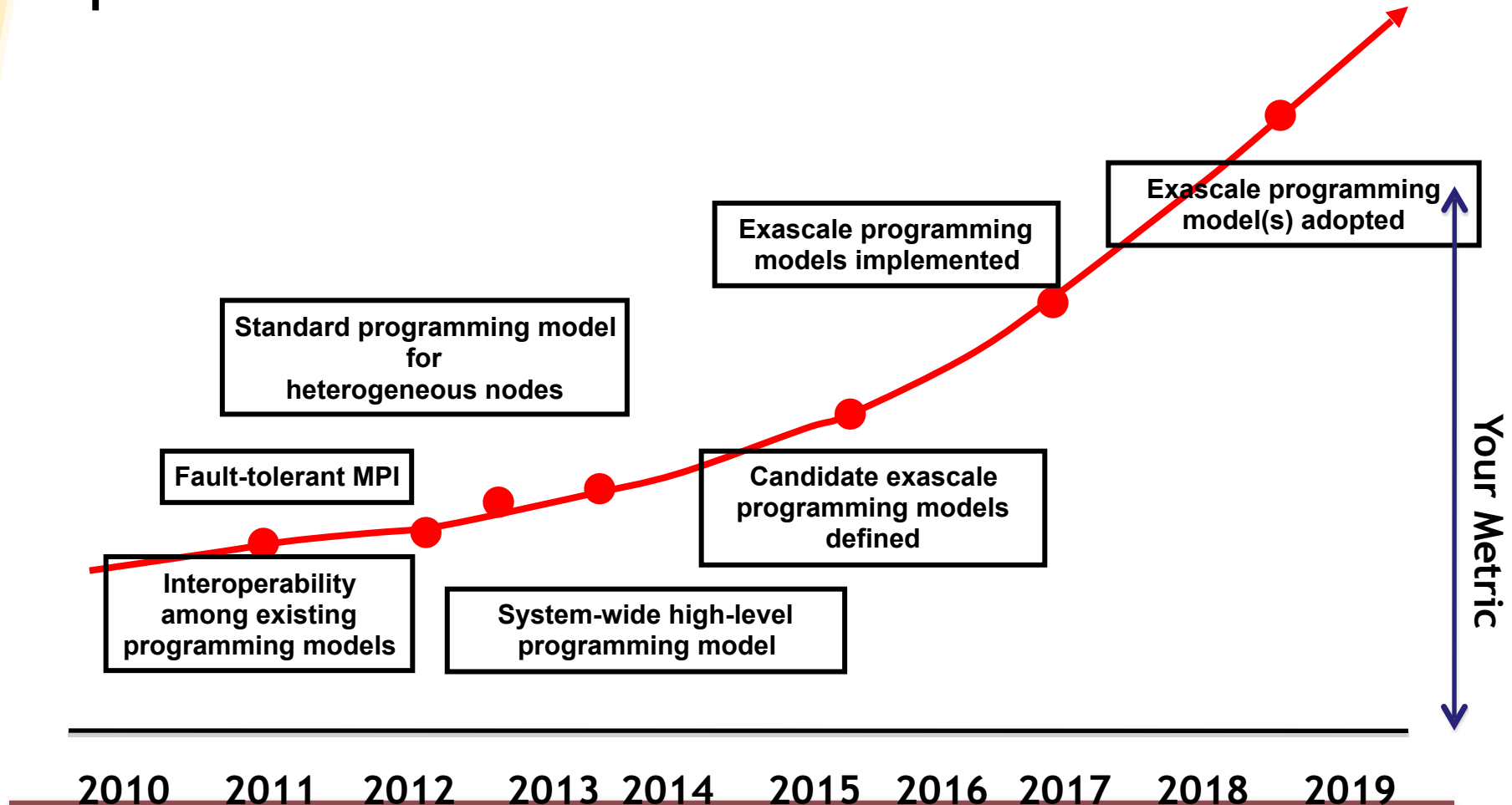
*More cores, less data sharing*

- Outside of quantum, neuromorphic, architectures expected to evolve into "extreme" versions of today's systems
  - 3D stacked processors, less cache, more on-die memory, more specialization, optical interconnects
- We need programming languages that meet tomorrow's needs as well as today's application goals
  - **Address needs of systems with diverse, extremely complex memory hierarchies**
  - **Able to handle more (and more kinds of) devices and high core counts**
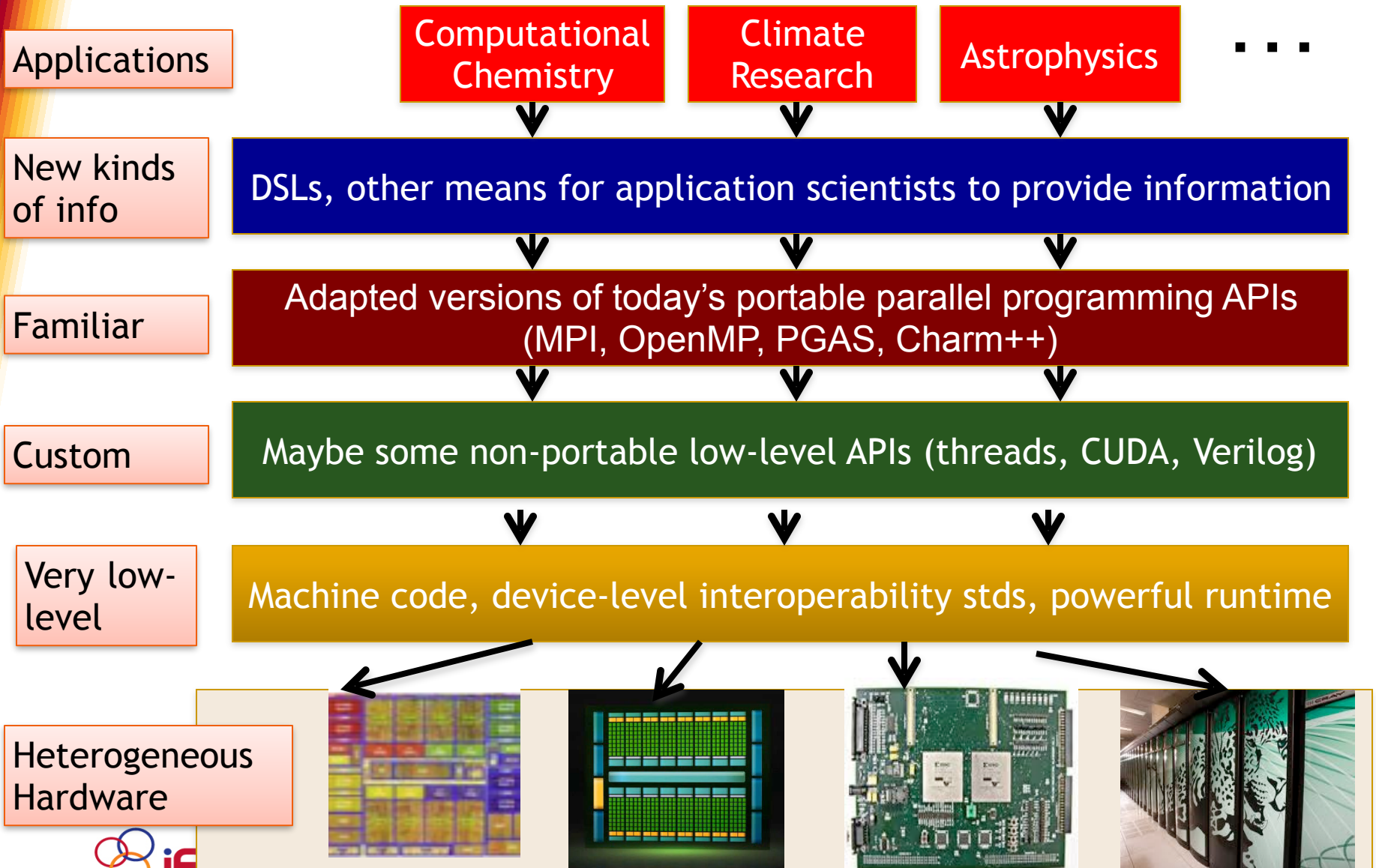  - **Facilitate interoperability, especially with internode approaches**

# IESP Programming Models
## International Exascale Software Project, 2010-11

**Proposed timeline**

# A Layered Programming Approach

Applications

Computational Chemistry | Climate Research | Astrophysics | ...

New kinds of info

DSLs, other means for application scientists to provide information

Familiar

Adapted versions of today's portable parallel programming APIs (MPI, OpenMP, PGAS, Charm++)

Custom

Maybe some non-portable low-level APIs (threads, CUDA, Verilog)

Very low-level

Machine code, device-level interoperability stds, powerful runtime

Heterogeneous Hardware

# Summit

2 Power9s, 6 GPUs per node

27,648 NVIDIA Tesla V100s, each with:

- 5120 CUDA cores
- 640 Tensor cores
- 300 GB/s BW (NVLink 2.0)
- 20MB registers, 16MB cache, 16GB HBM2 @900 GB/s
- 7.5 DP TFLOPS; 15 SP TFLOPS, 120 FP16 TFLOPS

Tensor cores do mixed precision multiply add of 4x4 matrices



| Type | Size | Range | $u = 2^{-t}$ |
|------|------|-------|--------------|
| half | 16 bits | $10^{\pm 5}$ | $2^{-11} \approx 4.9 \times 10^{-4}$ |
| single | 32 bits | $10^{\pm 38}$ | $2^{-24} \approx 6.0 \times 10^{-8}$ |
| double | 64 bits | $10^{\pm 308}$ | $2^{-53} \approx 1.1 \times 10^{-16}$ |
| quadruple | 128 bits | $10^{\pm 4932}$ | $2^{-113} \approx 9.6 \times 10^{-35}$ |



$$D = AB + C$$

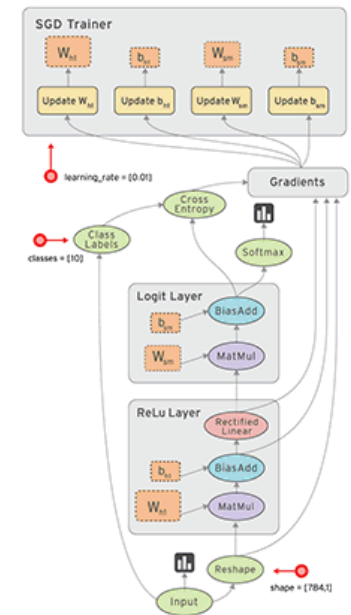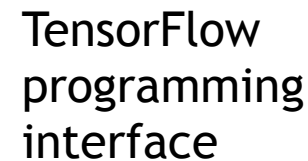INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

The Modeling & Simulation community can benefit from utilizing mixed / reduced precision

- Eg: Possible to achieve 4x FP64 peak for 64bit LU on V100 with iterative mixed precision (Dongarra et al.)

# Programming: HPC vs. DL

Nodes in a cluster

SMP Multicore Architecture

**Processor** | **Processor** | **Processor**

M | M | M

**Bus or Switch Network**

**Shared Memory**

**Network for Data Exchange**

Programmed using MPI

Programmed using OpenMP

TensorFlow programming interface

INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# OpenMP

Portable parallel programming since 1997

- Compiler directives
- Data, task, SIMD parallelism
- Multicores, GPUs
- User specifies the strategy, not the details

Maintained by industry consortium
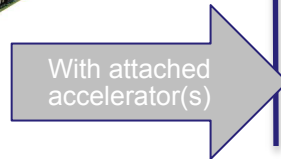
- It is now easy for academics to join

The mission of the OpenMP ARB (Architecture Review Board) is to standardize directive-based multi-language high-level parallelism that is performant, productive and portable.

- OpenMP 4+ supports heterogeneous systems (accelerators/devices)

- Accelerator model
  - Host device and attached
  - One or more ta      devices

Single device attached

Multiple devices attached

With attached accelerator(s)

Host Device
(CPU Multicore)

Xeon Phi(s) –
(Accelerator and self-hosted)

GPU(s)

iACS INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

- BerkeleyGW is a C++ application which computes the excited state properties of materials

- GPP contains the self-energy computation: large matrix reductions over complex arrays in a single loop nest of 4 loops

The OpenMP implementation with XL compiler achieves approximately the same run time as a tuned CUDA implementation



*Lower is better*

Results from: Rahulkumar Gayatri, "A Case Study for Performance Portability Using OpenMP 4.5", WACCPD-18

Chris Daley, NERSC

INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# Emerging OpenMP Features

- New features in OpenMP 5.0
  - Memory model: deal with memory heterogeneity

    ```
    #pragma omp allocate(A) allocator(omp_high_bw_mem_alloc)
    ```
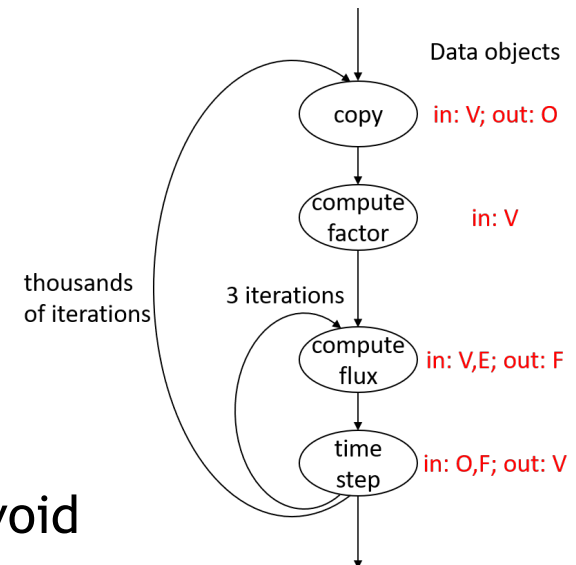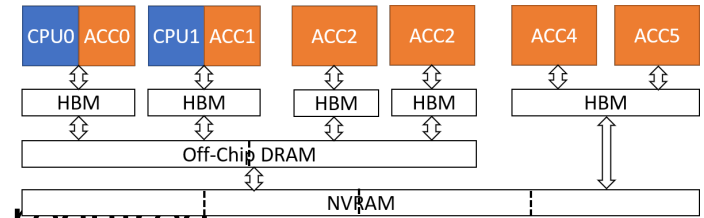
  - "concurrent" directive: descriptive parallel method

    ```
    #pragma omp for order(concurrent)  for (int i = 0; i < N; i++) {…}
    ```

- Also support for unified memory, deep copy of data, metadirectives, task affinity, and more

- Significant implementation effort to support these well

iACS INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# Extreme Heterogeneity: Memory Management

- Expect significantly more complex memory systems in the future

- Programming environment innovations are required to cope with such complexity

- Programming model
  - Provide high level data abstractions
  - Improve memory systems' programmability
  - Integrate performance modeling

- Compiler and runtime
  - Unified memory optimizations (prefetching, pinning) to avoid thrashing
  - Coarse grained data optimization
  - Better usage of hardware capabilities e.g. avoid caching of non-temporal accesses

| CPU0 | ACC0 | CPU1 | ACC1 | ACC2 | ACC2 | | ACC4 | ACC5 |

HBM HBM HBM HBM HBM

Off-Chip DRAM

NVRAM

Data objects

copy — in: V; out: O

compute factor — in: V

thousands of iterations

3 iterations

compute flux — in: V,E; out: F

time step — in: O,F; out: V

INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# OpenMP Loop Feature

- Work-sharing directive: split loop iterations among threads

- Significant challenges: large reduction basic blocks, irregular accesses (e.g. a[b[i]]), deeply nested conditionals, deep copy / allocation at target (where is my data?)

- Strategies vary per architecture

```
#pragma omp for order(concurrent
for (i = 1; i < N; i++) {
  f(A,i);
}
```

i=[1..N/4]          i=[N/2+1..3N/4]

thread1    thread2    thread3    thread4

i=[N/4+1..N/2]          i=[3N/4+1..N]

```
void f(double A[N+1], int i) {
  if (A[i] > 0.5)
    A[i] += i;
}
```
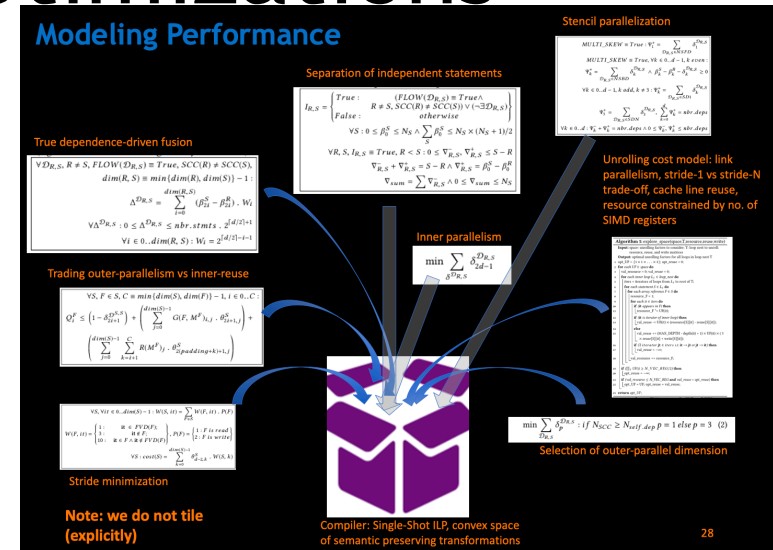
Where is A?
Where should
it be?

# Model-driven, Composable and Multi-Target Compiler Optimizations

**Motivation**: reduce performance gap between general purpose and domain specific compiler frameworks

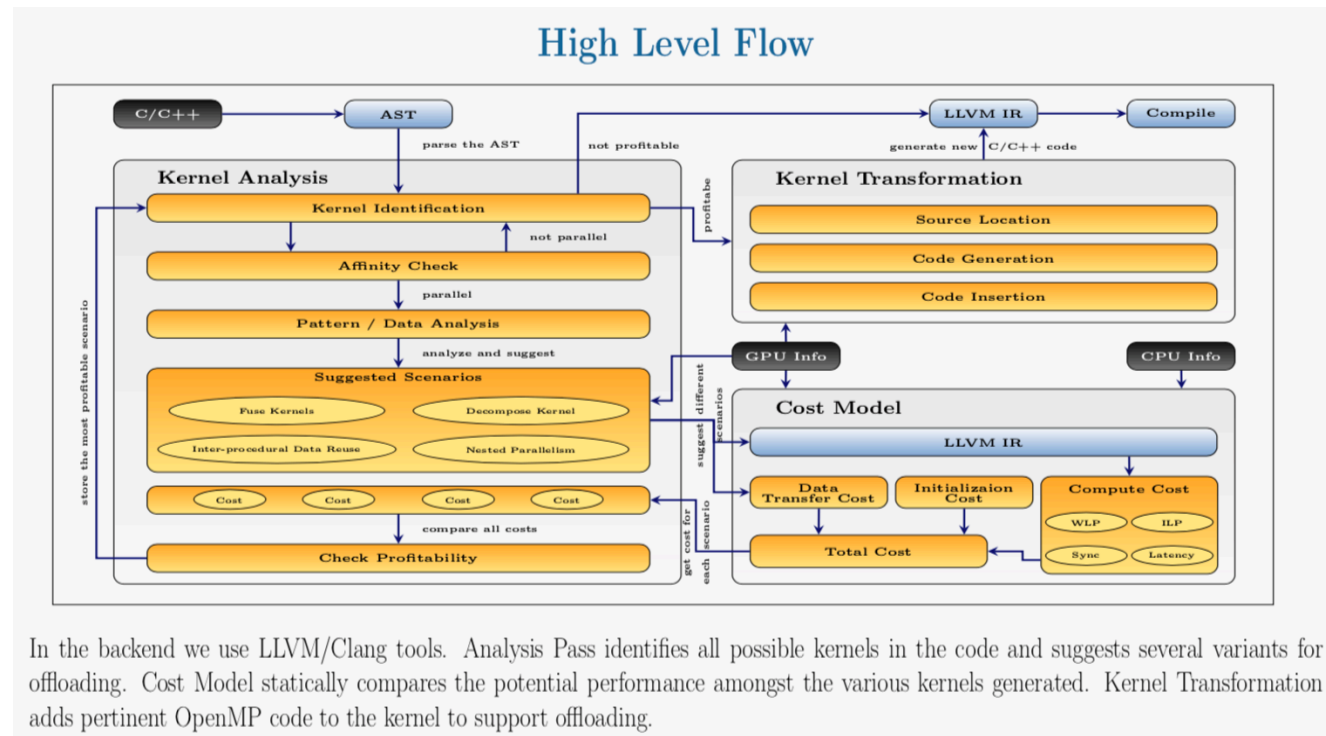**Impact**: substantial performance improvements, reduced application tuning time



Martin Kong: To appear in ACM Programming Language Design and Implementation, PLDI'19

# Kernel Fusion/Decomposition for Automatic GPU-offloading

**Motivation**: automatic GPU-offloading capabilities in LLVM to maximize application performance and user productivity

**Impact**: immediate benefit for medium to large scale mathematical libraries (e.g. Grid++ Lattice QCD parallel library)
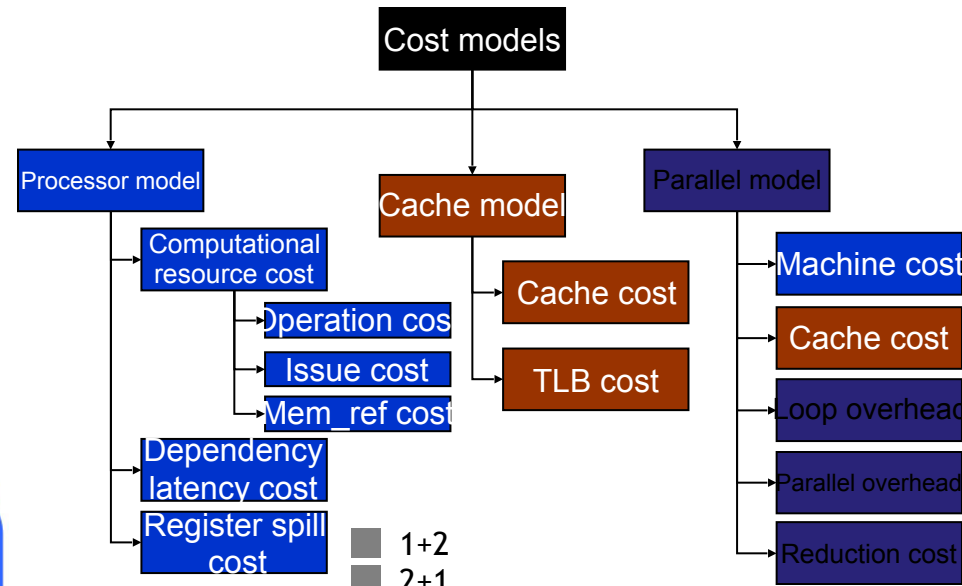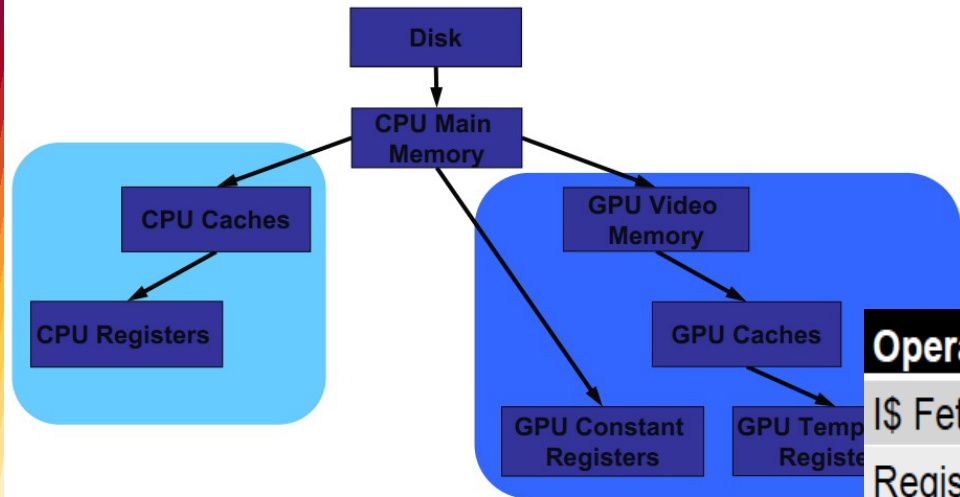


In the backend we use LLVM/Clang tools. Analysis Pass identifies all possible kernels in the code and suggests several variants for offloading. Cost Model statically compares the potential performance amongst the various kernels generated. Kernel Transformation adds pertinent OpenMP code to the kernel to support offloading.

Third Place, ACM Student Research Competition at the International Conference of Code Generation and Optimization (CGO'19, Washington DC)
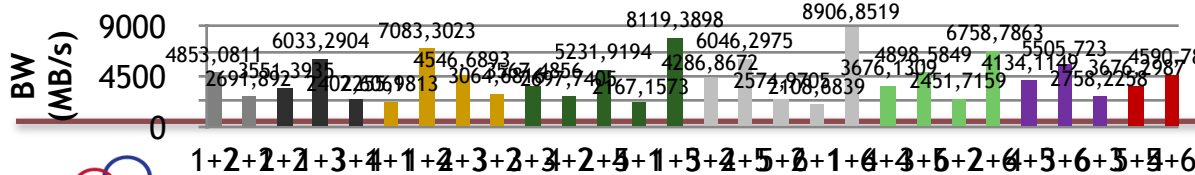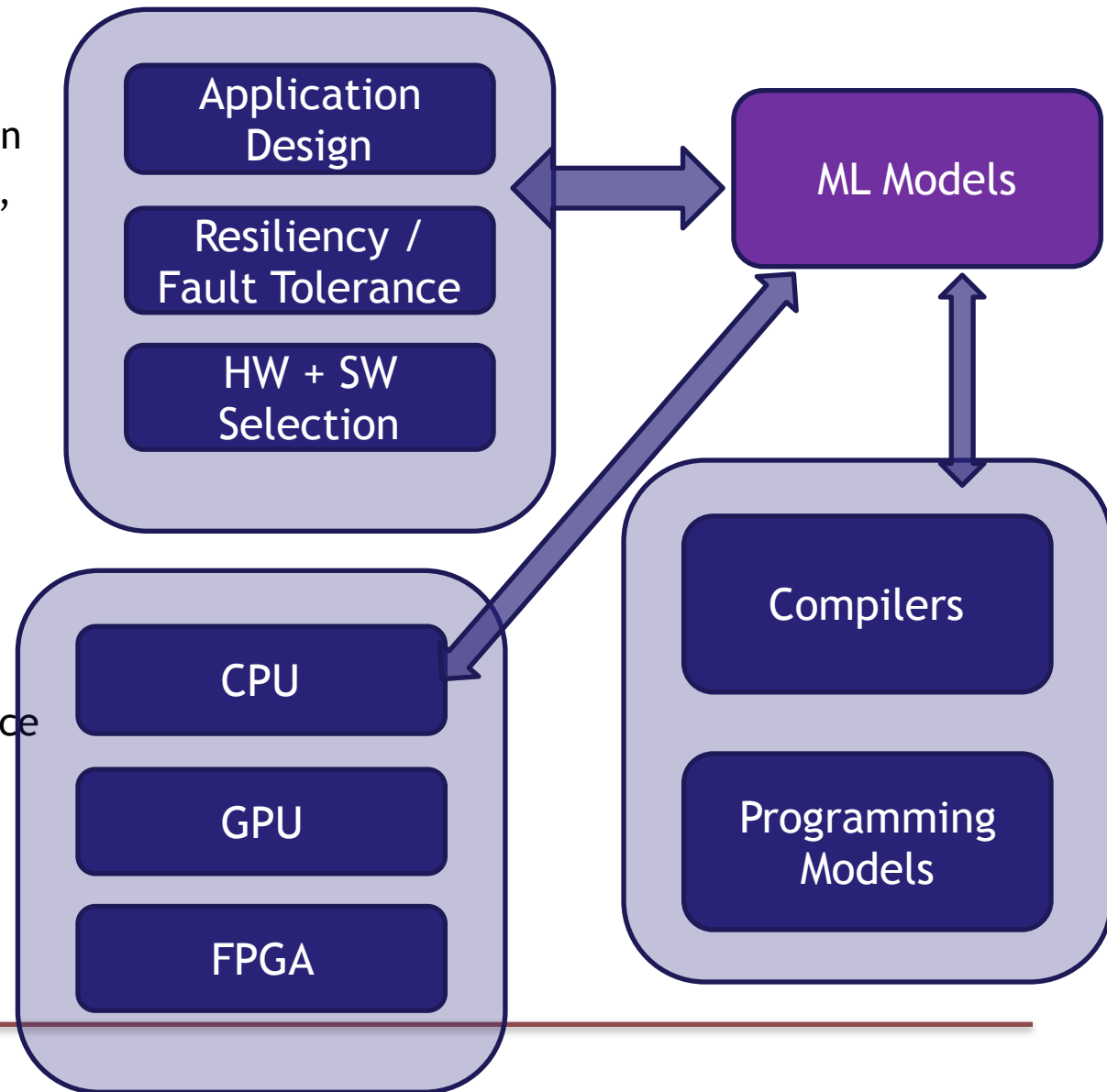Alok Mishra, Martin Kong, Barbara Chapman

**iACS** INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# Model-Driven Translation



**Cost models**

Processor model → Computational resource cost → Operation cost, Issue cost, Mem_ref cost; Dependency latency cost; Register spill cost

Cache model → Cache cost, TLB cost

Parallel model → Machine cost, Cache cost, Loop overhead, Parallel overhead, Reduction cost

Disk → CPU Main Memory → CPU Caches → CPU Registers; GPU Video Memory → GPU Caches → GPU Constant Registers, GPU Temp Registers

| Operation | Energy (pJ) | DP FLOPs | Insts* |
|---|---|---|---|
| I$ Fetch | 33 | 0.67 | 2.0 |
| Register Access (3W) | 10.5 | 0.2 | 0.6 |
| Access 3 D$ | 100 | 2 | 6 |
| Access 3 L2 D$ | 460 | 9 | 27 |
| Access 3 off chip | 762 | 15 | 45 |
| Access 3 from DRAM | 6000 | 120 | 360 |

**HT3 BW vs Threads**

BW (MB/s)

9000, 4500, 0

4853,0811  6033,2904  7083,3023  8119,3898  8906,8519  6046,2975  6758,7863  5505,723  4590,789
2691,892  3551,3935  2402,2606  4546,6893  5231,9194  4286,8672  3676,1300  4898,5849  4134,1149  3675,2987
2506,9813  3064,6867  2697,7405  2667,1573  2574,9705  2100,0839  2451,7159  2758,2258

**Thread Configuration (# of remote + # of local threads)**

Legend: 1+2, 2+1, 3+2, 3+1, 4+1, 1+4, 2+3, 3+2, 3+3, 4+2, 2+4, 5+1, 1+5, 3+4, 2+5, 5+2, 6+1, 1+6, 4+4, 3+5, 6+2, 2+6, 4+5

iACS INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# Exploiting ML in Software Stack

❑Design: determine best target machine given inherent application traits (memory or compute driven, accuracy?)

❑Compiler/runtime support for stringent power caps

❑Resiliency and fault tolerance: when and what to checkpoint?

❑Compiler options and runtime features embody a large and complex optimization space

❑Impact on application performance

❑Different application classes require different optimization strategies

Application Design

Resiliency / Fault Tolerance

HW + SW Selection

ML Models

CPU

GPU

FPGA

Compilers

Programming Models

iACS INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# Changing Workloads
## Data Analytic Computing (DAC) and Modeling and Simulation (M&S)

**DAC or AI**
- DAC community is demanding increased computational intensity, hence facing barriers to scalability along with new requirements for interoperability, robustness, and reliability of results
- Physics-Informed ML

**M&S**
- M&S community is demanding a more dynamic interaction between analysis and simulations by increasing use of large-scale data analytics

DAC and M&S have traditionally relied on different hardware and software stacks moving to a coherent platform for M&S and DAC benefits both while maximizing returns on R&D investments

**Challenge: Software layers of the HPC environment**
- A more agile and reusable HPC software portfolio that is equally capable in DAC and M&S will improve productivity, increase reliability and trustworthiness, and increase sustainability

**Challenge: Application design**
- Componentization, extensibility, scalability, reusability, interoperability
- Build on libraries and motif-based toolkits

**iACS** INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

Doug Kothe, Exascale Computing Project

# Large Scale Petroleum: Combining HPC and Data Analytics



Seismic Data Processing and Interpretation
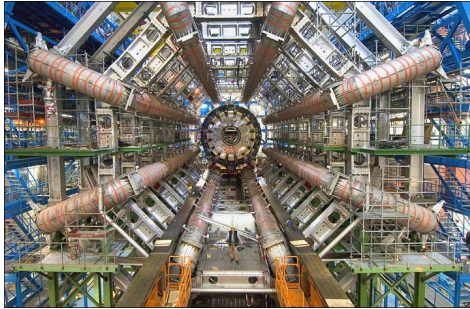
- A combination of Physics and Data Science
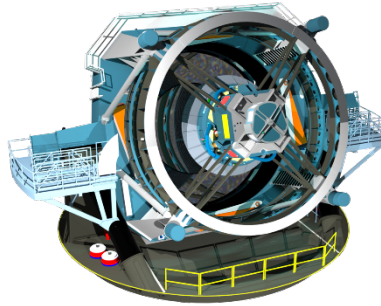
Reservoir Simulation

Oilfield Data Analytics

iACS INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# Scalable Learning Matters!

**Large Hadron Collider**

25GB/s, >200PB
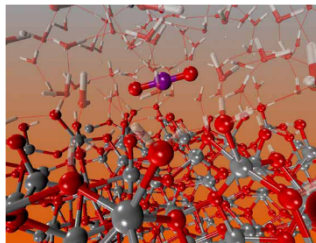
**LSST**

20TB/night, 73PB

**Transmission EM**

3GB/s

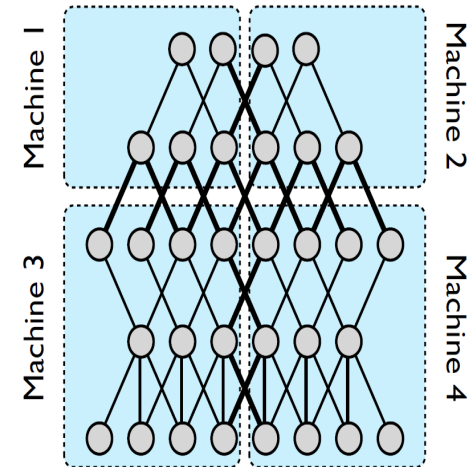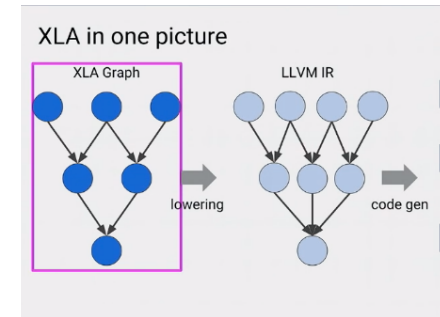**ExaFEL (ECP)**
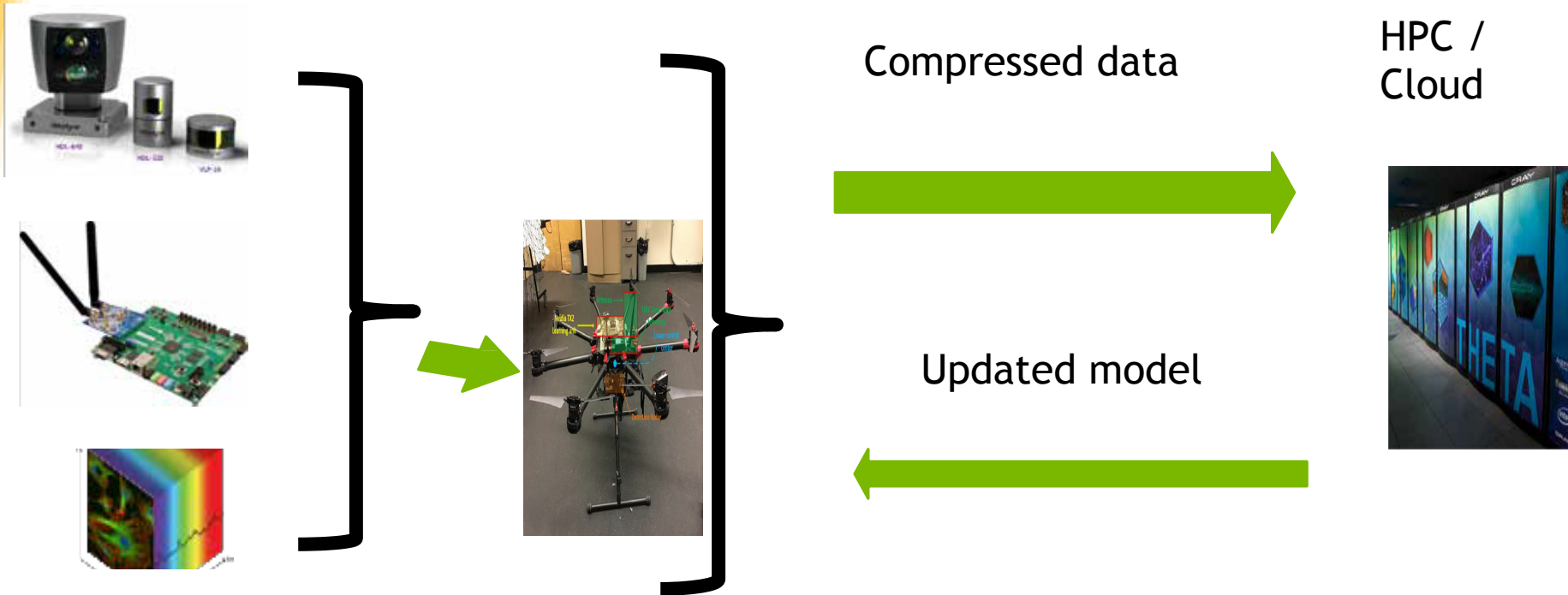
10PB

**MD Trajectories**

~32PB / simulation

**Meta-genomics**

~15GB / sample

iACS INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE

# Machine Learning for Scientific Applications

- Today's machine learning frameworks are not easy to scale
  - Expensive communication / synchronization in SGD
  - Significant effort optimizing frameworks for CPU, GPU
  - Use of OpenMP
  - Research making advances in distributed parallelization
- State of art
  - TensorFlow, PyTorch,…
  - Heavy optimization of operations
  - Otherwise, much lacking in compiler technology
- Language and compiler enhancements
  - Improved single node translation
  - Efficient translation of associated user code in Python
  - DSL features for enhanced ML?
  - Path toward integration with scientific application code?
  - Representation of data and computation flow for wholistic optimization?

# ML at the Edge



Compressed data

HPC / Cloud

Updated model

Inference at (potentially many) edge devices, major training in central computer.  Regular updates of model must be sent back to edge.

# One Size Fits All?  Tasks and Data Flow

- Inspired by the data flow execution model (Dennis, 1970s)
  - Cilk, TBB, OpenMP,...
  - Legion, HPX, Parsec,...
  - Google Cloud Dataflow, Tensorflow, ....
- Not a "natural" approach for many applications
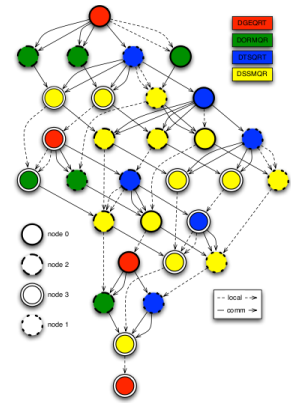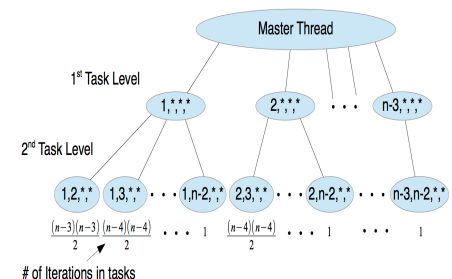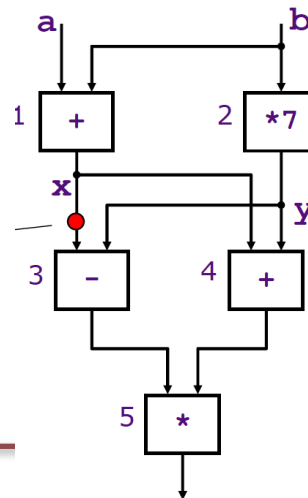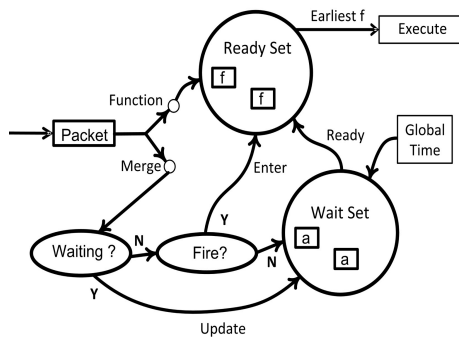  - Implementations often forget data locality
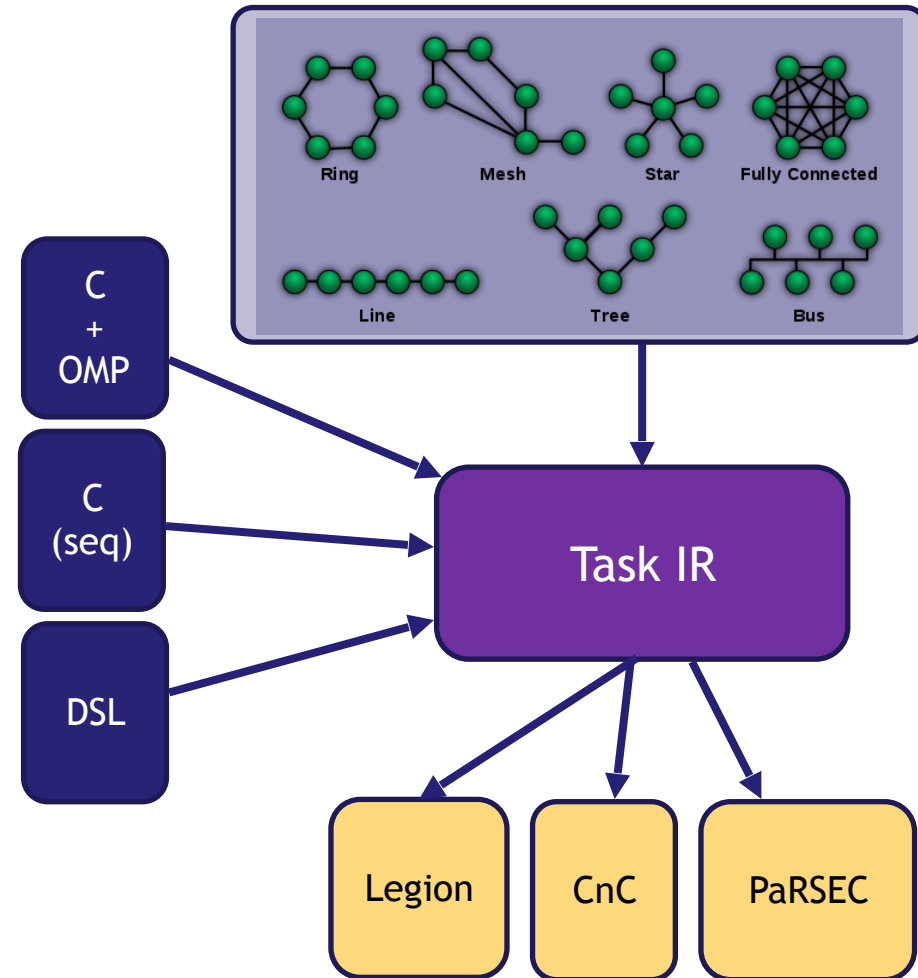
Fig. 5.   DAG of QR for a 4x4 tile matrix.

# High-Performance and High-Productivity Programming With Task-Based Execution

The performance promise of task / dataflow runtimes is hampered by the lack of high-productivity systems that encourage their use

Can we design a unified system for distributed scientific applications targeting tasking / dataflow runtimes?:
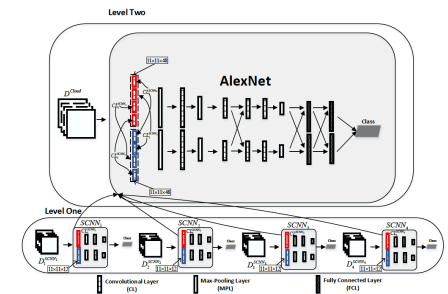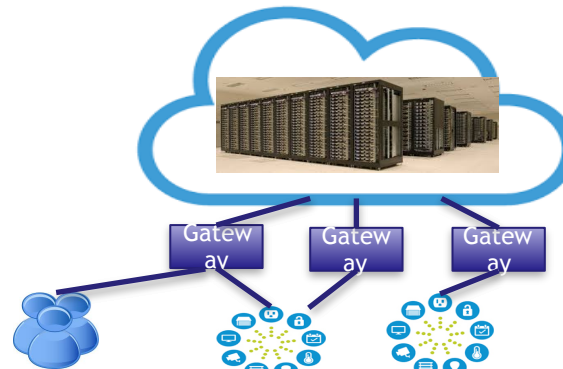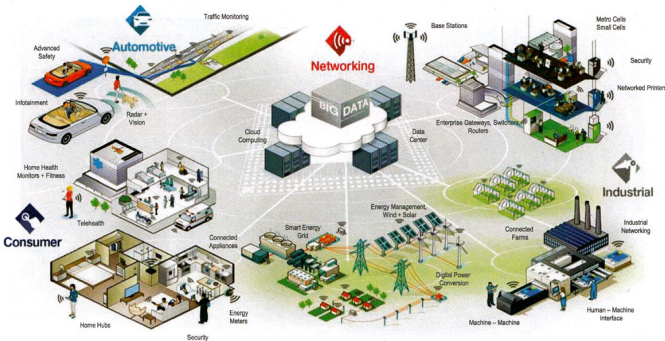
- ❑ Implement compiler transformations that optimize programs for a variety of target runtimes (Legion, Concurrent Collections, ParSEC) and platforms
- ❑ Design a common intermediate representation to allow for several input specifications: including sequential, explicitly parallel (e.g OpenMP) and domain specific languages
- ❑ Embed topological information and exploit it to find optimal task and data mappings that take locality fully into account
- ❑ Develop intra-node work and data partitioning strategies to exploit multi-GPU execution
- ❑ Can we extend this to include ML frameworks?



Ring  Mesh  Star  Fully Connected

Line  Tree  Bus

C + OMP

C (seq)

DSL

Task IR

Legion   CnC   PaRSEC

# Are We There Yet?

- High Performance Computing, Cloud Computing, Edge Computing, Fog Computing
  - Computing anywhere, anytime, any devices
  - Extreme Heterogeneity, Deep Memory Hierarchy
  - Performance, Resilience, Elasticity, Productivity, Power
- Data-driven, scientific and AI code in (peaceful?) co-existence



**How do we program these complex systems?**


iACS INSTITUTE FOR ADVANCED COMPUTATIONAL SCIENCE