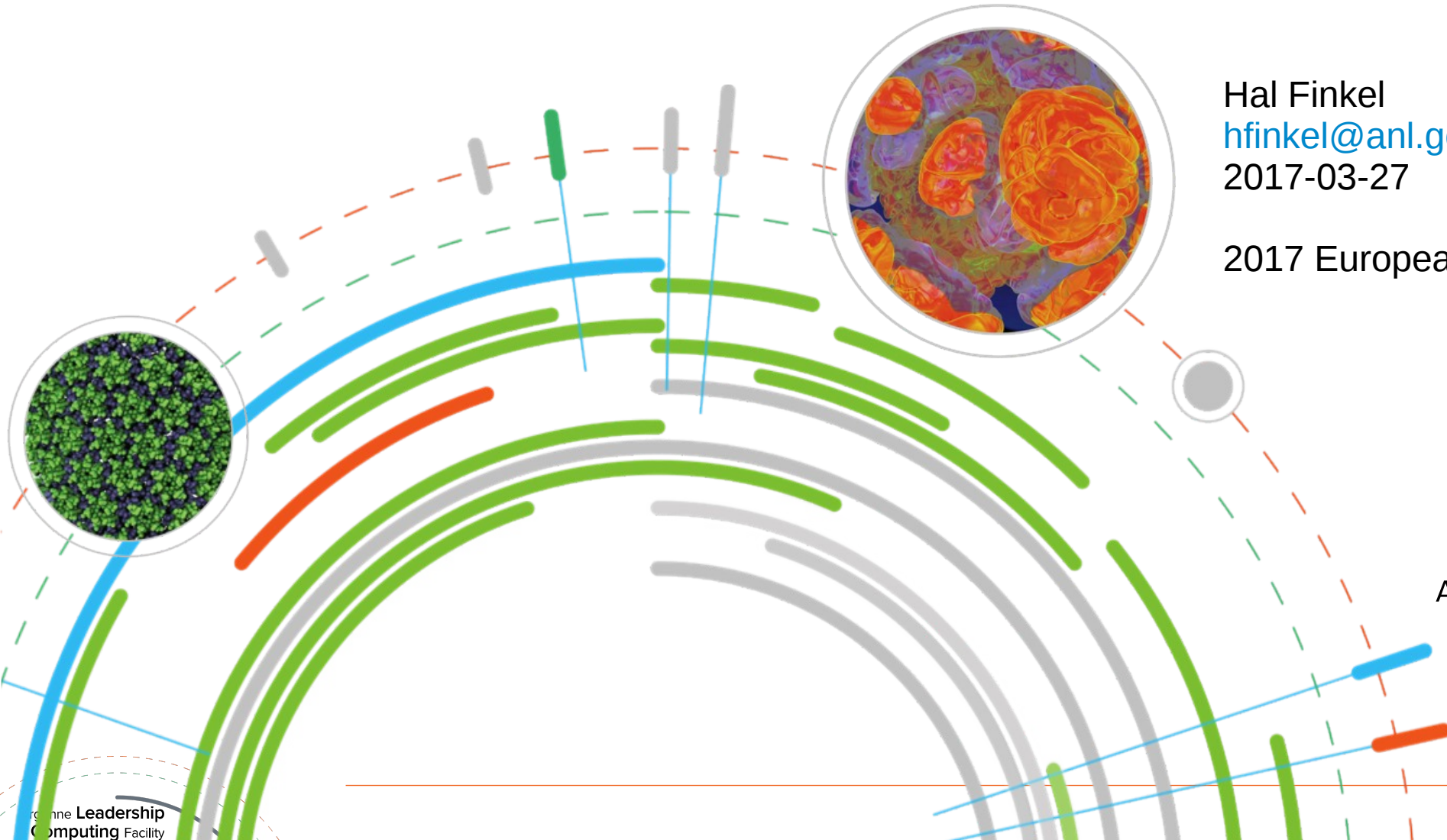


LLVM for the future of Supercomputing

Hal Finkel
hfinkel@anl.gov
2017-03-27

2017 European LLVM Developers' Meeting

Argonne **Leadership**
Computing Facility



What is Supercomputing?

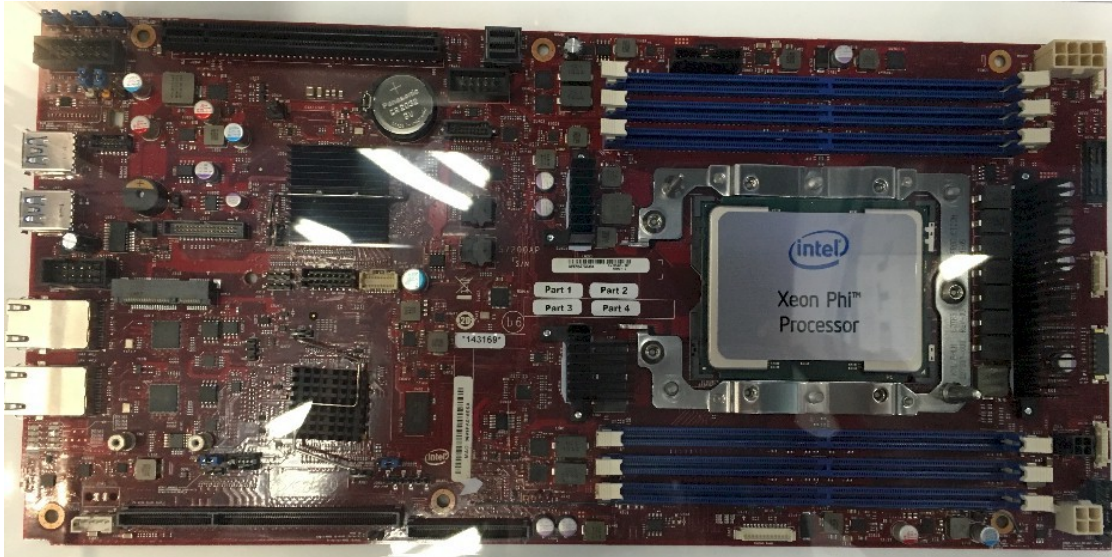
Computing for large, tightly-coupled problems.

Lots of computational
cap
lots of high

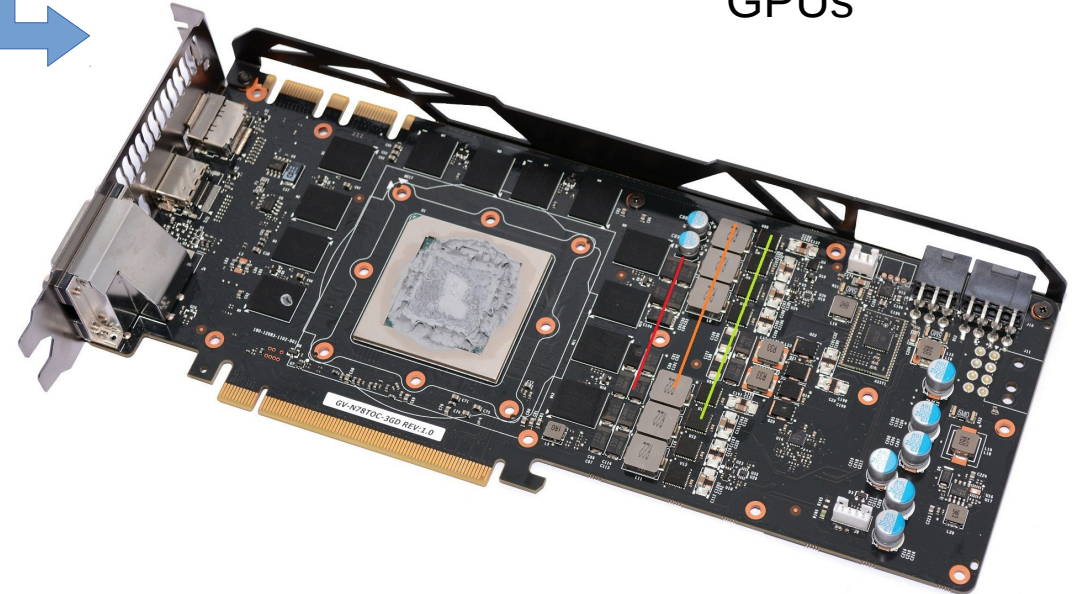
High computational density paired with a
high-throughput low-latency network.

Supercomputing “Swim Lanes”

“Many Core” CPUs



GPUs



<https://forum.beyond3d.com/threads/nvidia-pascal-speculation-thread.55552/page-4>

<http://www.nextplatform.com/2015/11/30/inside-future-knights-landing-xeon-phi-systems/>

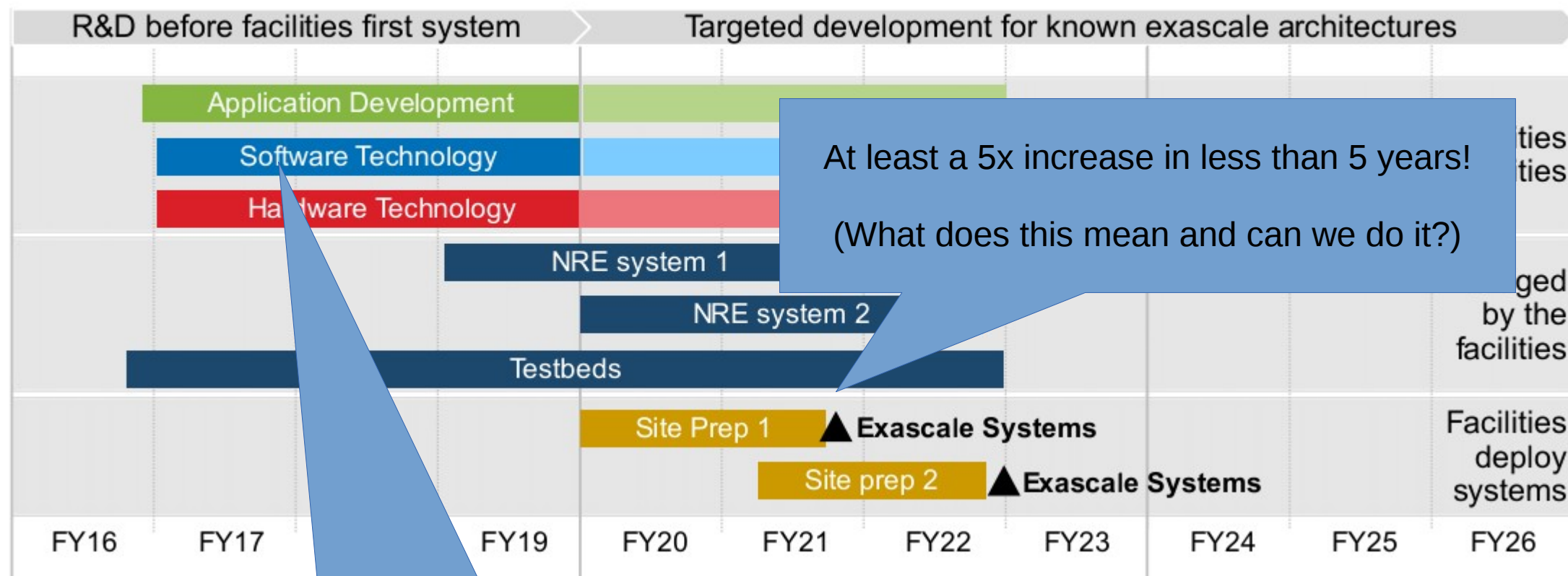
System attributes	NERSC Now	OLCF Now	ALCF Now	NERSC Upgrade	OLCF Upgrade	ALCF Upgrades	
Name Planned Installation	Edison	TITAN	MIRA	Cori 2016	Summit 2017-2018	Theta 2016	Aurora 2018-2019
System peak (PF)	2.6	27	10	> 30	200	>8.5	180
Peak Power (MW)	2	9	4.8	Only a 2.7x increase in power!			
Total system memory	357 TB	710TB	768TB	~1 PB DDR4 High Bandwidth Memory	~1 PB DDR4 High Bandwidth Memory	~1 PB DDR4 + High Bandwidth Memory	> 7 PB High Bandwidth Memory
Node performance (TF)	0.460	1	1	Our next system will have: 180 PF An 18x increase!			
Node processors	Intel Ivy Bridge	AMD Opteron Nvidia Kepler	PowerPC A2	core CPUs Intel Haswell CPU in data partition	& multiple Nvidia Volta GPUs	many core CPUs	Knights Hill Xeon Phi many core CPUs
	5,600	18,688	49,152	9,300 nodes	9,300 nodes	9,300 nodes	>50,000 nodes
			5D Torus	Aries	Dual Rail EDR-IB	Aries	2 nd Generation Intel Omni-Path Architecture
File System	168 GB/s, Lustre®	1 TB/s, Lustre®	26 PB 300 GB/s GPFS™	28 PB 744 GB/s Lustre®	120 PB 1 TB/s GPFS™	10PB, 210 GB/s Lustre initial	150 PB 1 TB/s Lustre®

The heterogeneous system with GPUs has 10x fewer nodes!

Still ~50,000 nodes.

Our next system is:
180 PF
(180 OPS)

High-level ECP technical project schedule



At least a 5x increase in less than 5 years!
(What does this mean and can we do it?)

We need to start preparing applications and tools now. exascaleproject.org for more information.

What Exascale Means To Us...

Exascale System	Goal
Delivery Date	2019-2020
Performance	1000 PF LINPACK and 300 PF on to-be-specified applications
Power Consumption*	20 MW
MTBAI**	6 days
Memory including NVRAM	128 PB
Node Memory Bandwidth	4 TB/s
Node Interconnect Bandwidth	400 GB/s

*Power consumption includes only power to the compute nodes, not to the storage or cooling systems.

**The mean time to application failure requiring a restart must be greater than 24 hours, and the asymptotic scaling must be greater than 0.5 over time. The system overhead to handle automatic restarts must not reduce application efficiency by more than half.

It means 5x the compute and 20x the memory on 1.5x the power!

<http://estrfi.cels.anl.gov/files/2011/07/RF>

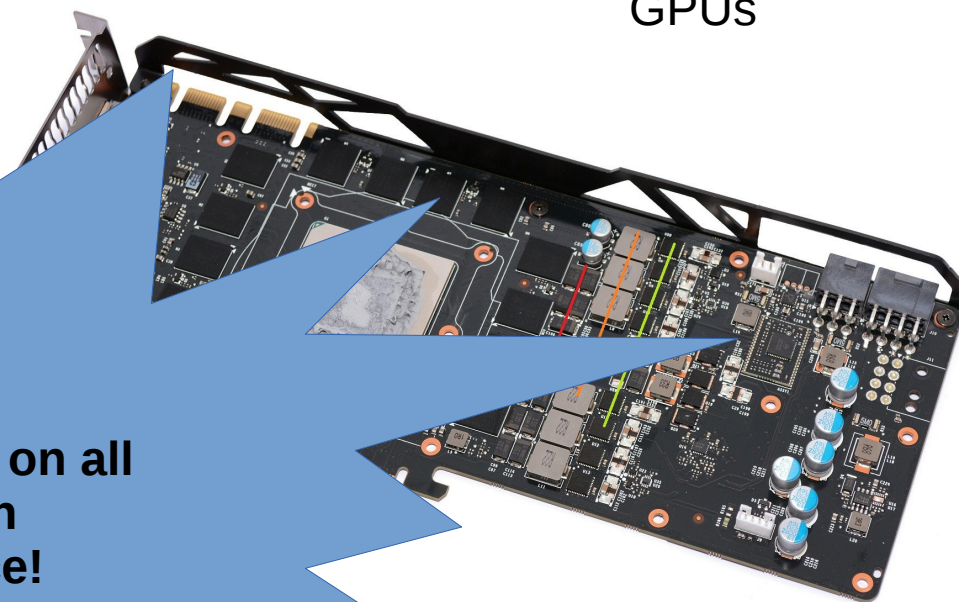
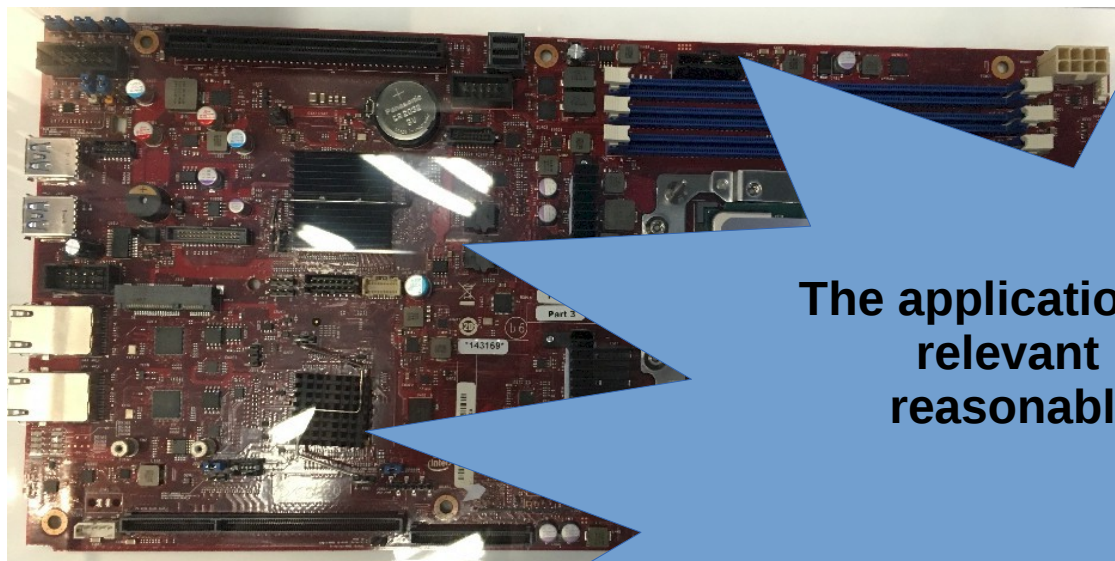
What do we want?

We Want Performance Portability!

Application
(One Maintainable Code Base)

“Many Core” CPUs

GPUs



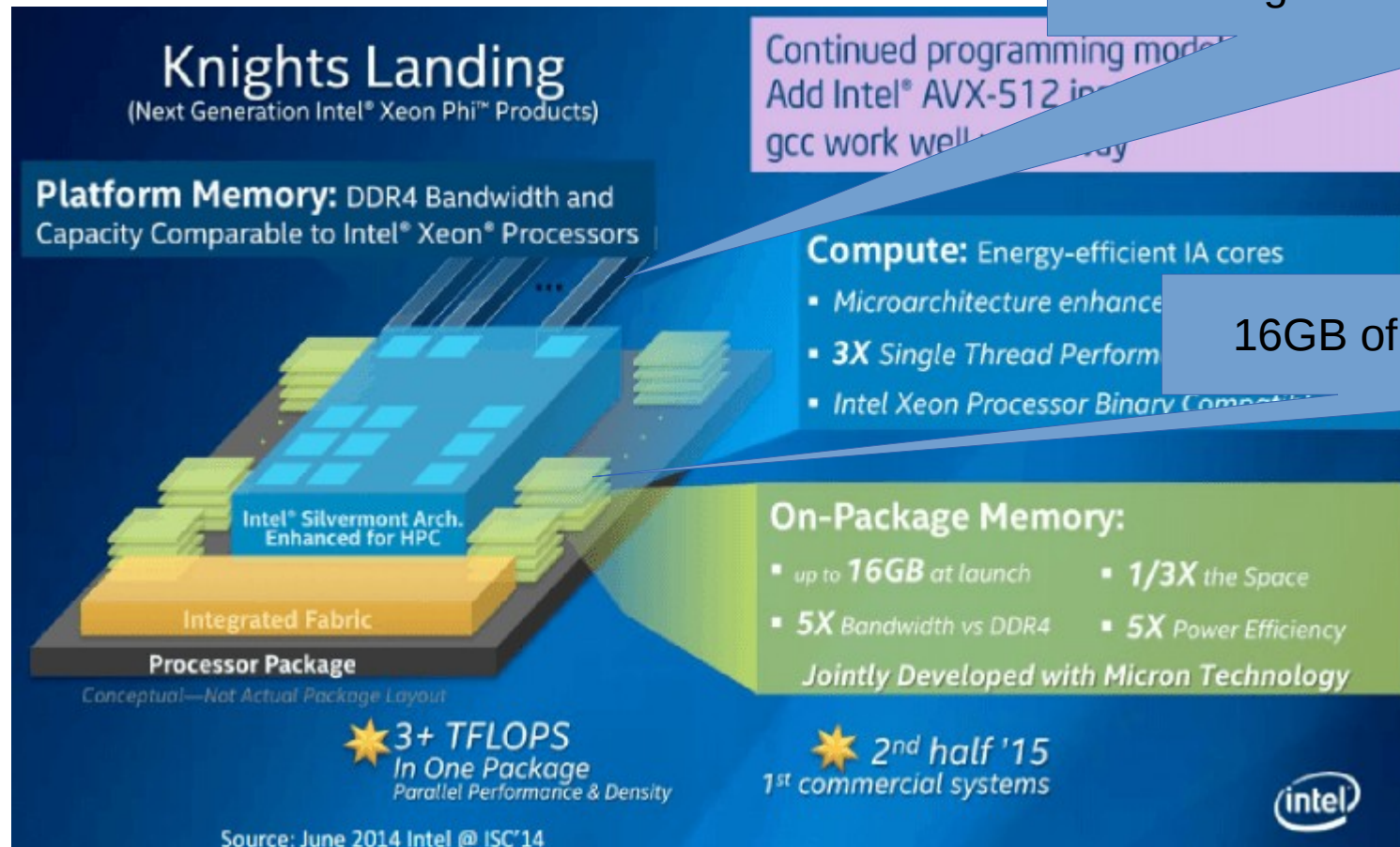
**The application should run on all
relevant hardware with
reasonable performance!**

<http://www.nextplatform.com/2015/11/11/next-platform-computing-facility/>

[m/threads/multi-processor-speculation-thread.55552/page-4](#)

Let's Talk About Memory...

Intel Xeon Phi HBM



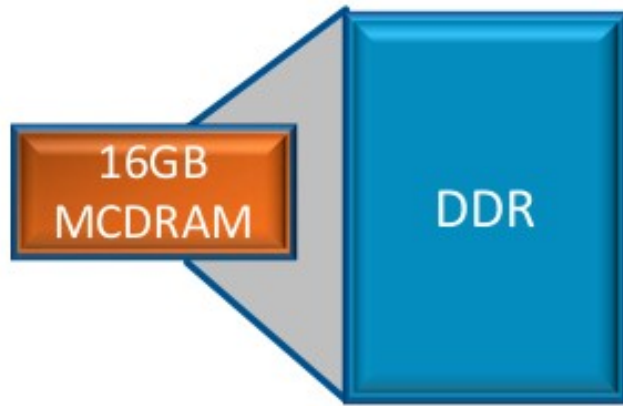
Large amounts of regular DRAM far away.

16GB of high-bandwidth on-package memory!

<http://www.techenablement.com/preparing-knights-landing-stay-hbm-memory/>

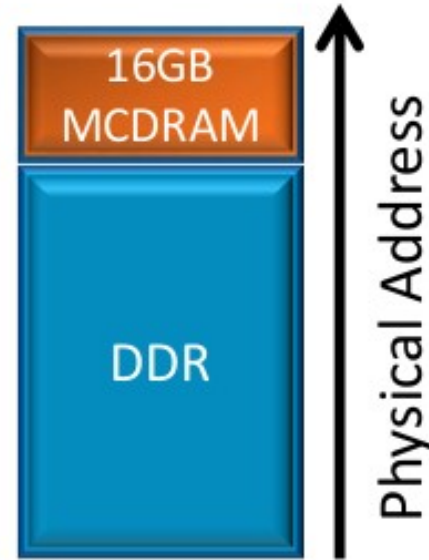
Three Modes. Selected at boot

Cache Mode



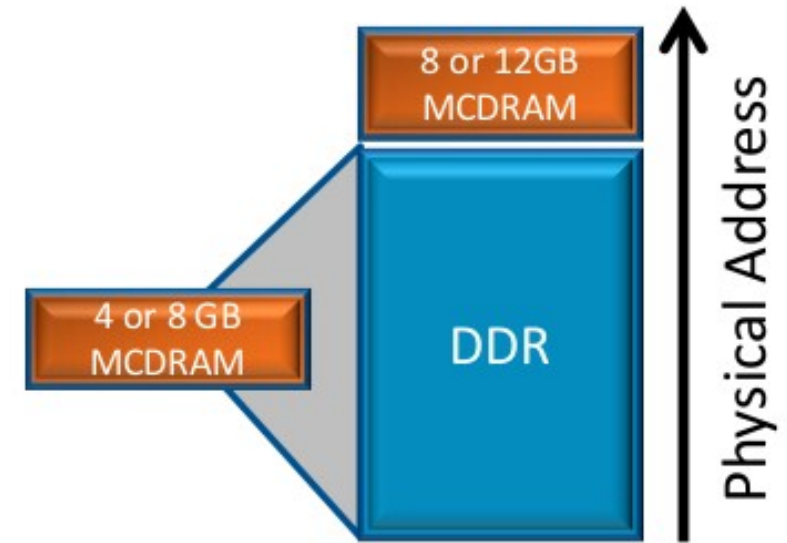
- SW-Transparent, Mem-side cache
- Direct mapped. 64B lines.
- Tags part of line
- Covers whole DDR range

Flat Mode



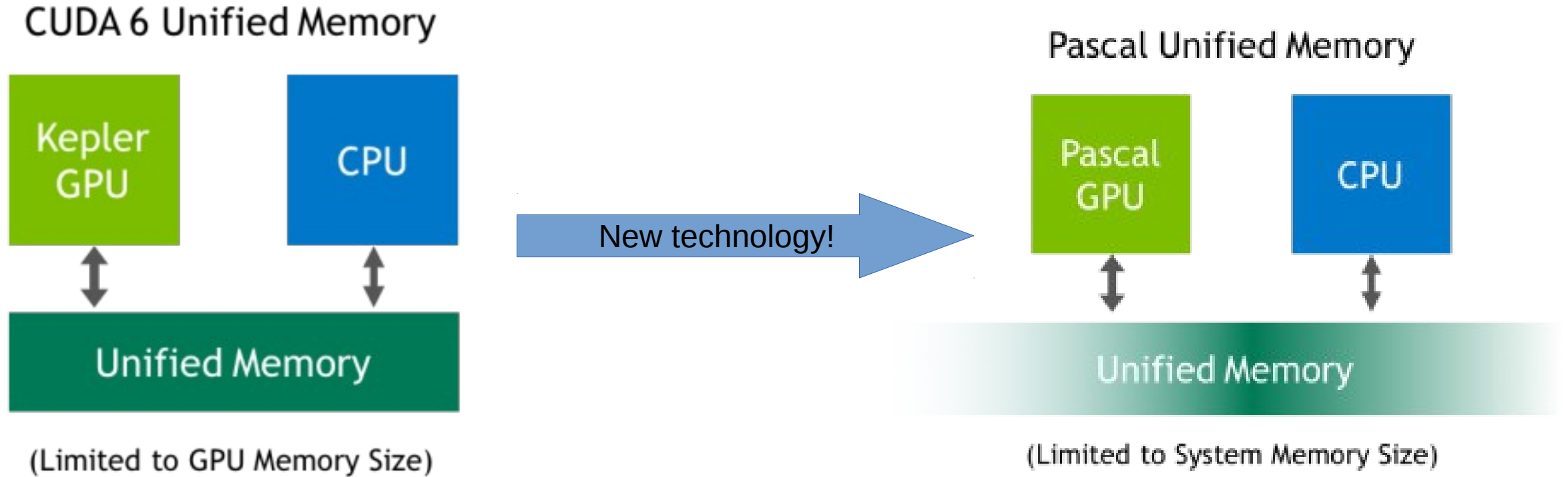
- MCDRAM as regular memory
- SW-Managed
- Same address space

Hybrid Mode



- Part cache, Part memory
- 25% or 50% cache
- Benefits of both

CUDA Unified Memory



Unified memory enables “lazy” transfer on demand – will mitigate/eliminate the “deep copy” problem!

CUDA UM (The Old Way)

CPU Code

```
void sortfile(FILE *fp, int N) {  
    char *data;  
    data = (char *)malloc(N);  
  
    fread(data, 1, N, fp);  
  
    qsort(data, N, 1, compare);  
  
    use_data(data);  
  
    free(data);  
}
```

CUDA 6 Code with Unified Memory

```
void sortfile(FILE *fp, int N) {  
    char *data;  
    cudaMallocManaged(&data, N);  
  
    fread(data, 1, N, fp);  
  
    qsort<<<...>>>(data, N, 1, compare);  
    cudaDeviceSynchronize();  
  
    use_data(data);  
  
    cudaFree(data);  
}
```


CUDA UM (The New Way)

CPU Code

```
void sortfile(FILE *fp, int N) {  
    char *data;  
    data = (char *)malloc(N);  
  
    fread(data, 1, N, fp);  
  
    qsort(data, N, 1, compare);  
  
    use_data(data);  
  
    free(data);  
}
```

Pascal Unified Memory*

```
void sortfile(FILE *fp, int N) {  
    char *data;  
    data = (char *)malloc(N);  
  
    fread(data, 1, N, fp);  
  
    qsort<<<...>>>(data, N, 1, compare);  
    cudaDeviceSynchronize();  
  
    use_data(data);  
  
    free(data);  
}  
*with operating system support
```

Pointers are “the same” everywhere!

How Do We Get Performance Portability?

How Do We Get Performance Portability? Shared Responsibility!

Applications and
Solver Libraries

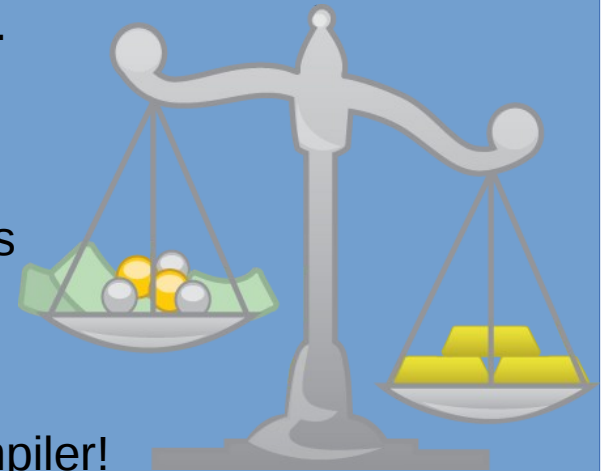
Applications and solver libraries must be flexible and parameterized! Why?

Trade-offs between...

- basis functions
- resolution
- Lagrangian vs. Eulerian representations
- renormalization and regularization schemes
- solver techniques
- evolved vs computed degrees of freedom
- and more...

cannot be made by a compiler!

Autotuning can help.

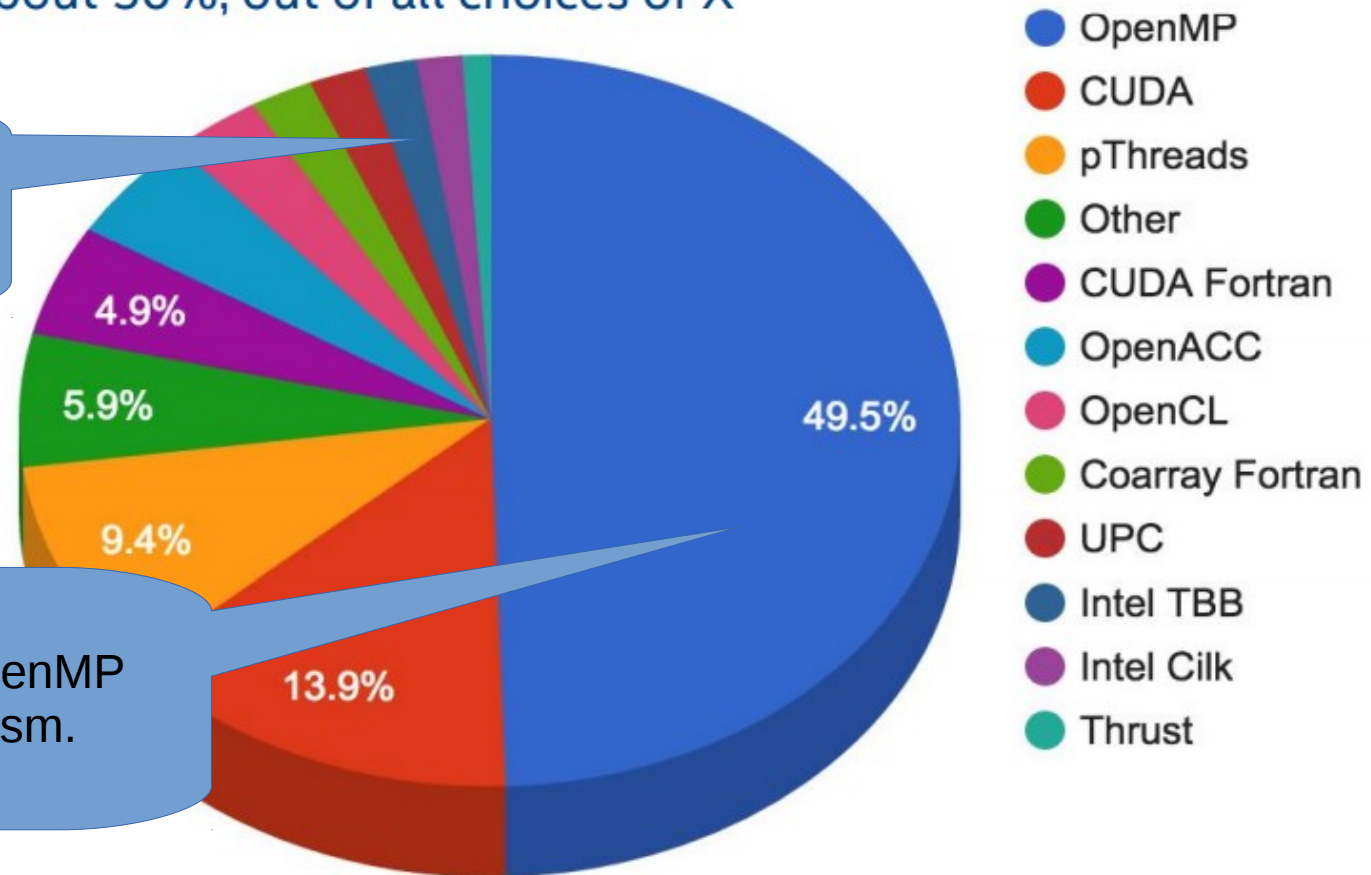


How do we express parallelism - MPI+X?

✓ OpenMP is about 50%, out of all choices of X

A minority of applications use abstraction libraries (TBB and Thrust on this chart)

In 2015, many codes use OpenMP directly to express parallelism.



Courtesy of Yun (Helen) He, Alice Koniges, et. al., (NERSC) at OpenMPCon'2015

<http://llvm-hpc2-workshop.github.io/slides/Tian.pdf>

How do we express parallelism - MPI+X?

But this is changing...

- We're seeing even greater adoption of OpenMP, but...
- Many applications are not using OpenMP directly. Abstraction libraries are gaining in popularity.

Often uses OpenMP and/or other compiler directives under the hood.

BB and Thrust.

Use of C++ Lambdas.

- RAJA (<https://github.com/LLNL/RAJA>)

```
RAJA::ReduceSum<reduce_policy, double> piSum(0.0);

RAJA::forall<execute_policy>(begin, numBins, [=](int i) {
    double x = (double(i) + 0.5) / numBins;
    piSum += 4.0 / (1.0 + x * x);
});
```

- Kokkos (<https://github.com/kokkos>)

```
int sum = 0;
// The KOKKOS_LAMBDA macro replaces
// the capture-by-value clause [=].
// It also handles any other syntax
// needed for CUDA.
Kokkos::parallel_reduce (n, KOKKOS_LAMBDA (const int i,
int& lsum) {

    lsum += i*i;
}, sum);
```


How do we express parallelism - MPI+X?

And starting with C++17, the standard library has parallel algorithms too...

Table 2 — Table of parallel algorithms

adjacent_difference	adjacent_find	all_of	any_of
copy	copy_if	copy_n	count
count_if	equal	exclusive_scan	fill
fill_n	find	find_end	find_first_of
find_if	find_if_not	for_each	for_each_n
generate	generate_n	includes	inclusive_scan
inner_product	inplace_merge	is_heap	is_heap_until
is_partitioned	is_sorted	is_sorted_until	lexicographical_compare
max_element	merge	min_element	minmax_element
mismatch	move	none_of	nth_element
partial_sort	partial_sort_copy	partition	partition_copy
reduce	remove	remove_copy	remove_copy_if
remove_if	replace	replace_copy	replace_copy_if
replace_if	reverse	reverse_copy	rotate
rotate_copy	search	search_n	set_difference
set_intersection	set_symmetric_difference	set_union	sort
stable_partition	stable_sort	swap_ranges	transform
transform_exclusive_scan	transform_inclusive_scan	transform_reduce	uninitialized_copy
uninitialized_copy_n	uninitialized_fill	uninitialized_fill_n	unique
unique_copy			

[Note: Not all algorithms in the Standard Library have counterparts in [Table 2](#). — end note]

// For example:

```
std::sort(std::execution::par_unseq, vec.begin(), vec.end()); // parallel and vectorized
```

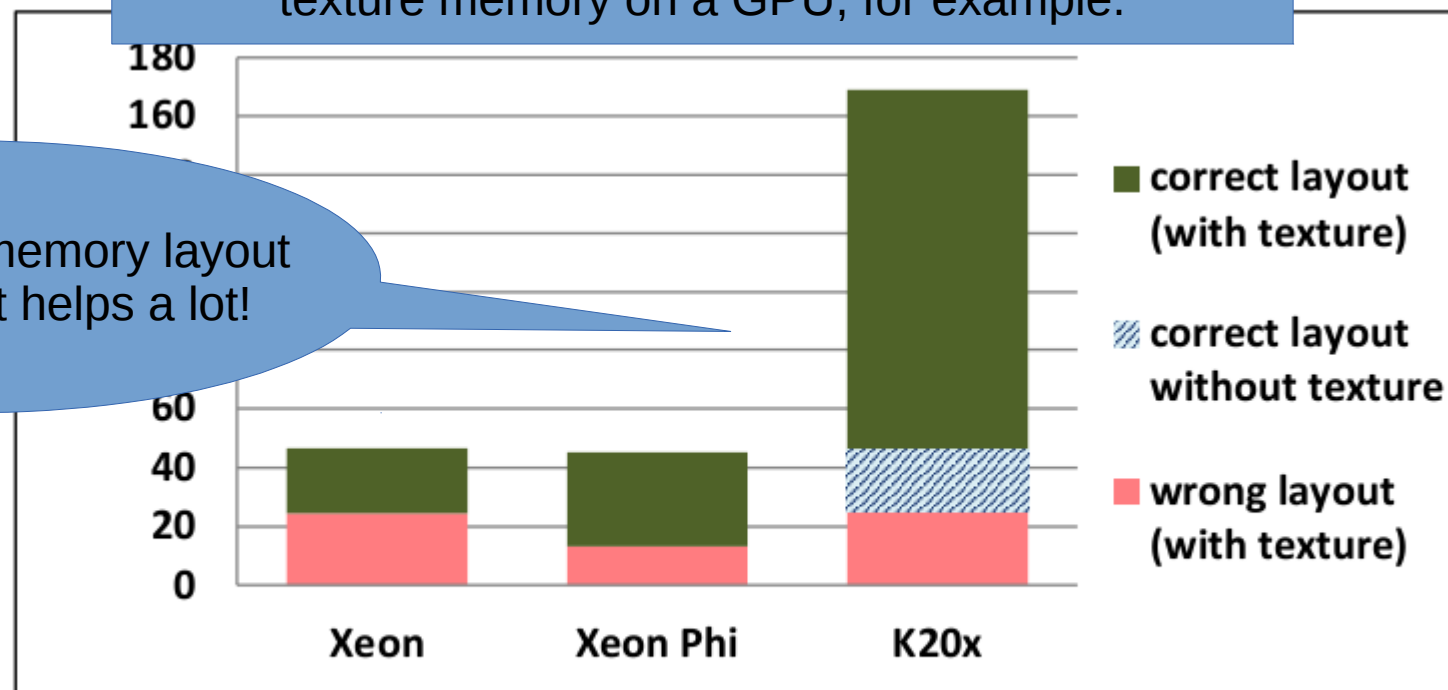
What About Memory?

It is really hard for compilers to change memory layouts and generally determine what memory is needed where. The Kokkos C++ library has memory placement and layout policies:

```
View<const double ***, Layout, Space, MemoryTraits<RandomAccess>> name (...);
```

Constant random-access data might be put into texture memory on a GPU, for example.

Using the right memory layout and placement helps a lot!



- Large loss in performance with wrong layout

The Exascale Computing Project – Improvements at All Levels

Applications
and Solver Libraries

Over 30 Application and Library Teams

Libraries Abstracting
Memory and Parallelism

Kokkos, RAJA, etc.

Compilers and Tools

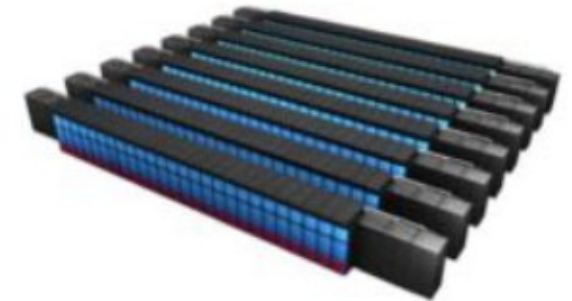
SOLLVE, PROTEAS,
Y-Tune, ROSE,
Flang, etc.

Hardware
Technology

Hardware technology
elements

Exascale
Systems

Integrated exascale
supercomputers

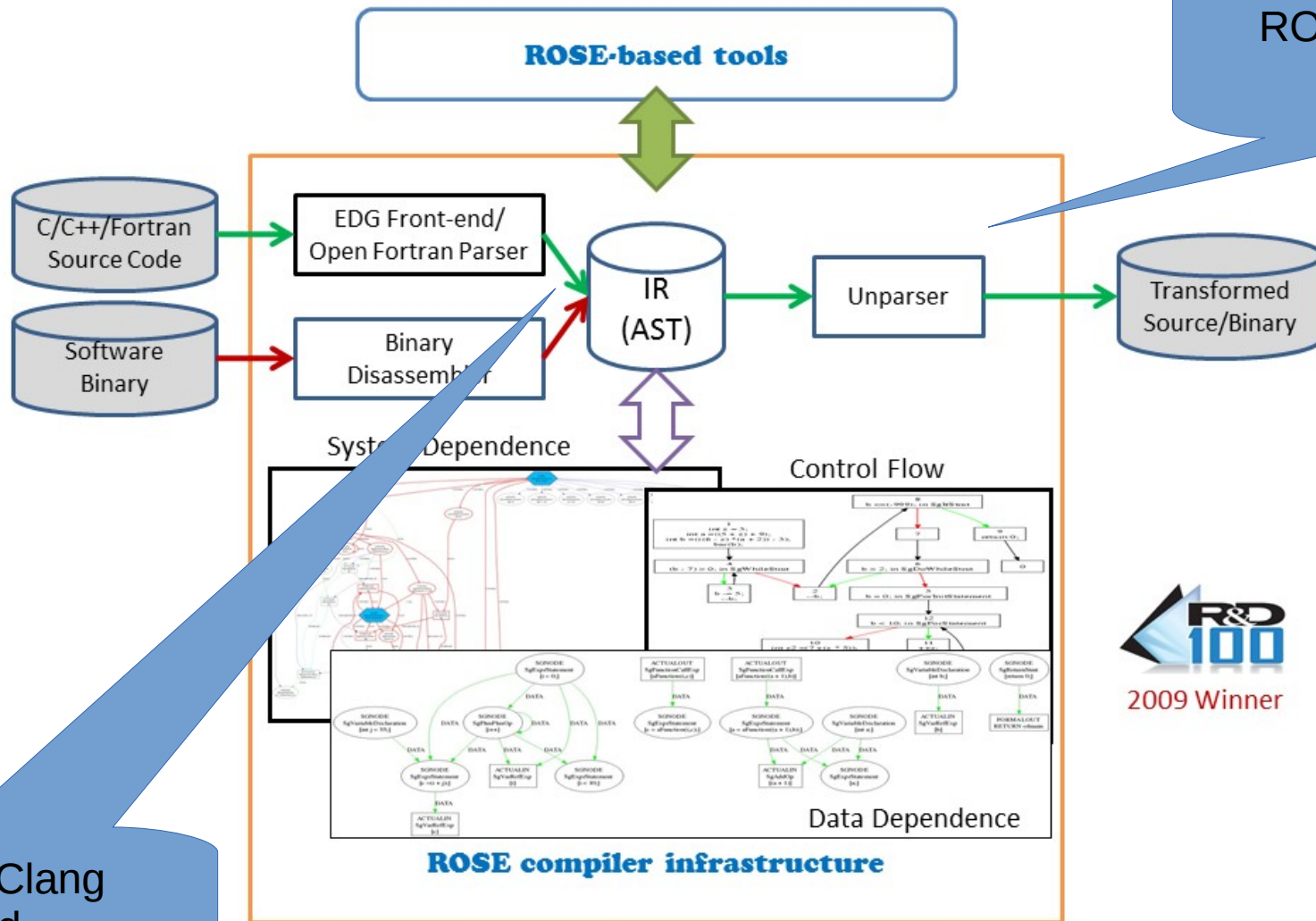


Now Let's Talk About LLVM...

LLVM Development in ECP



ROSE - Advanced Source-to-Source Rewriting



ROSE can generate LLVM IR.

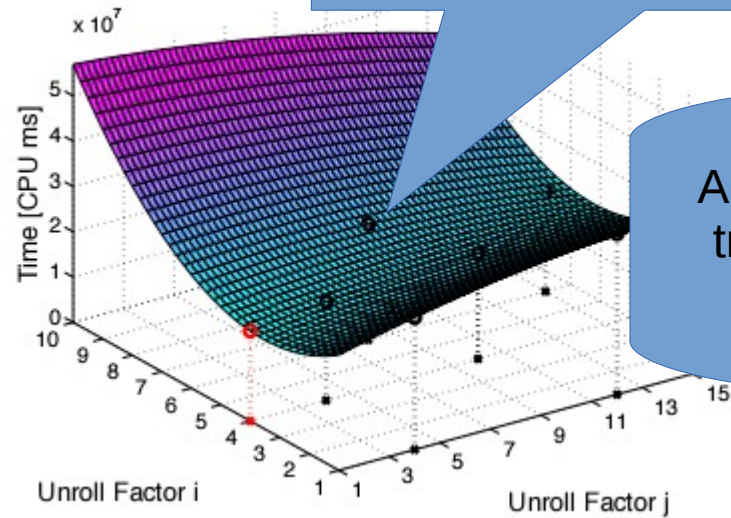
ROSE can use Clang as a frontend.



<http://rosecompiler.org/>

Y-Tune

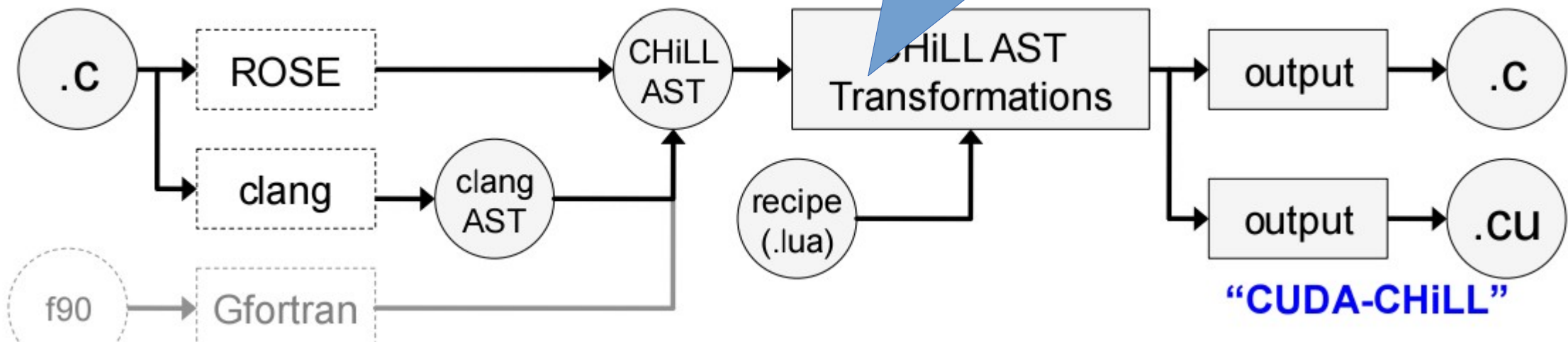
Machine-learning assisted search and optimization.



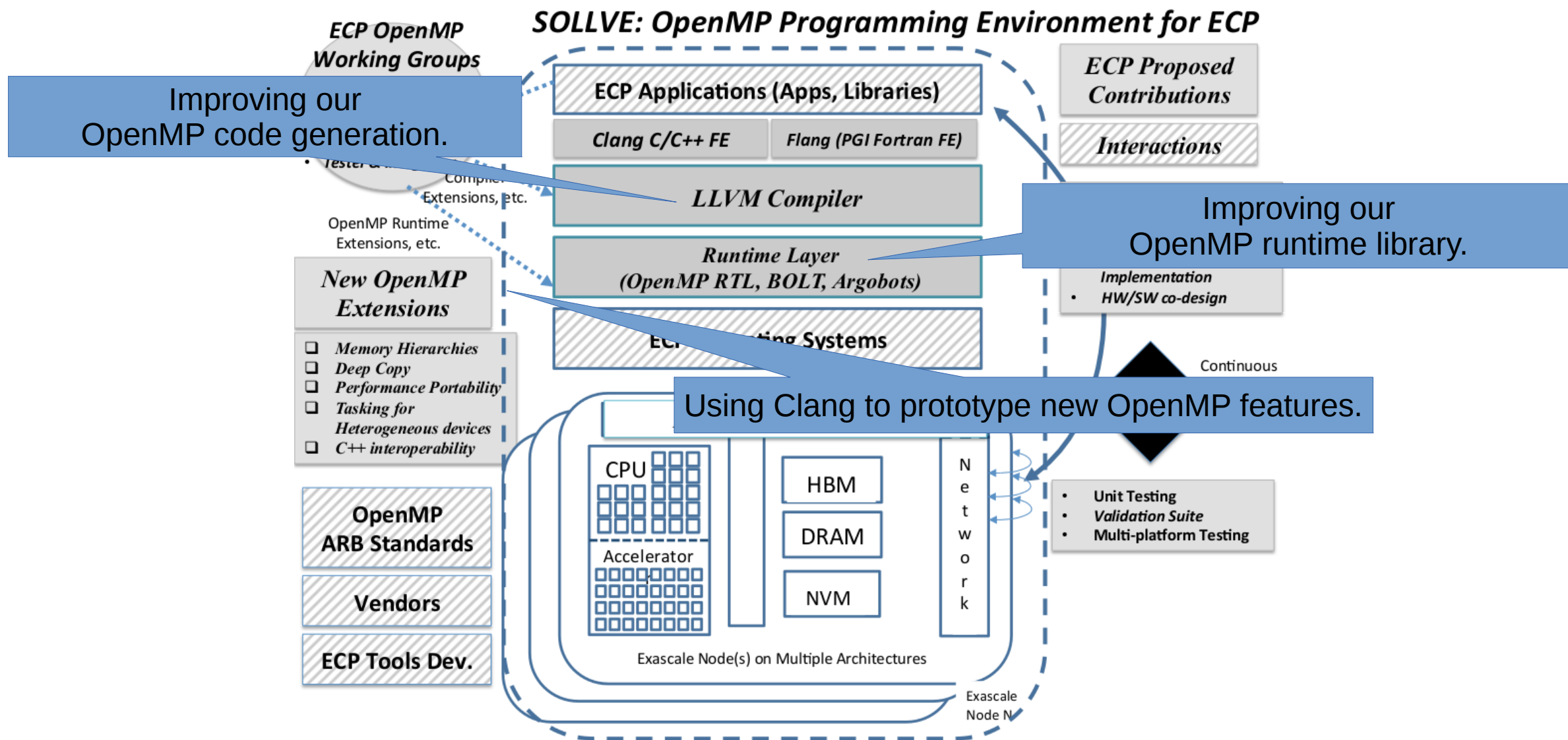
Y-Tune's scope includes improving LLVM for:

- Better optimizer feedback to guide search (via pragmas)

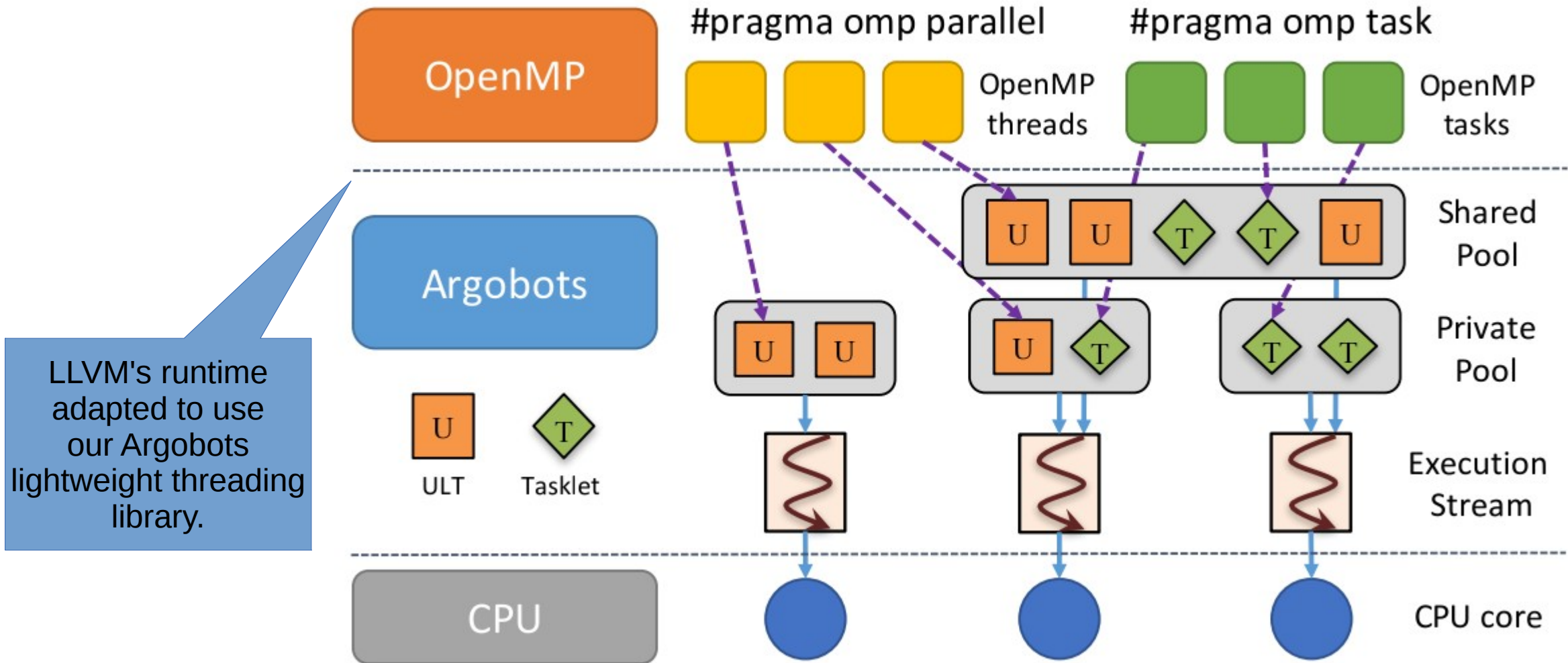
Advanced polyhedral and **application-specific** operator transformations. We can deal with the combined space of compiler-assisted and algorithm tuning!



SOLLVE – “Scaling Openmp with LLVM for Exascale performance and portability”



BOLT - “**B**OLT is **O**penMP over **L**ightweight **T**hreads” (Now Part of SOLLVE)



<http://www.bolt-omp.org/>

BOLT - “BOLT is OpenMP over Lightweight Threads” (Now Part of SOLLVE)

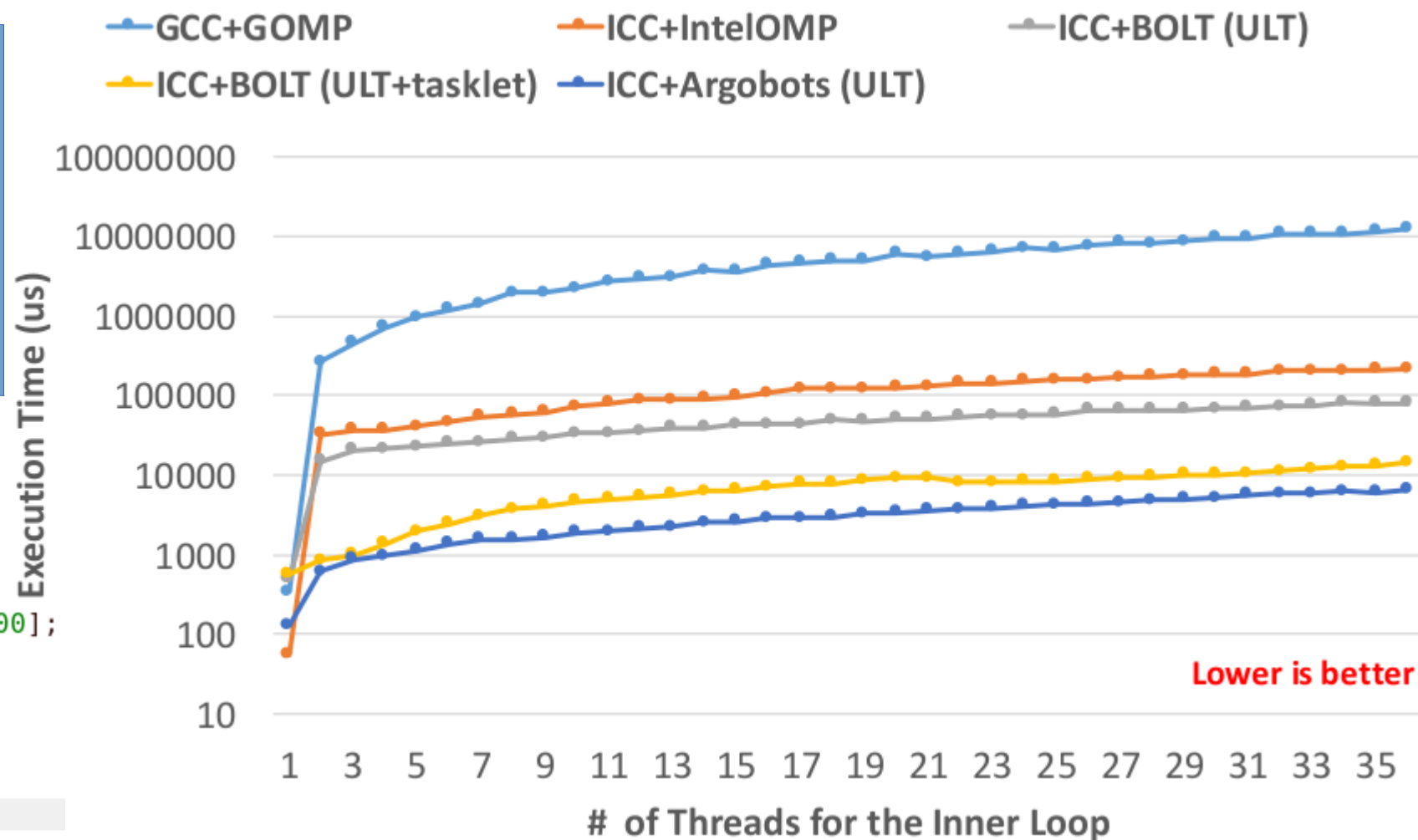
BOLT beats other runtimes by at least 10x on this nested parallelism benchmark.

Critical use case for composability!

```
int in[1000][1000], out[1000][1000];

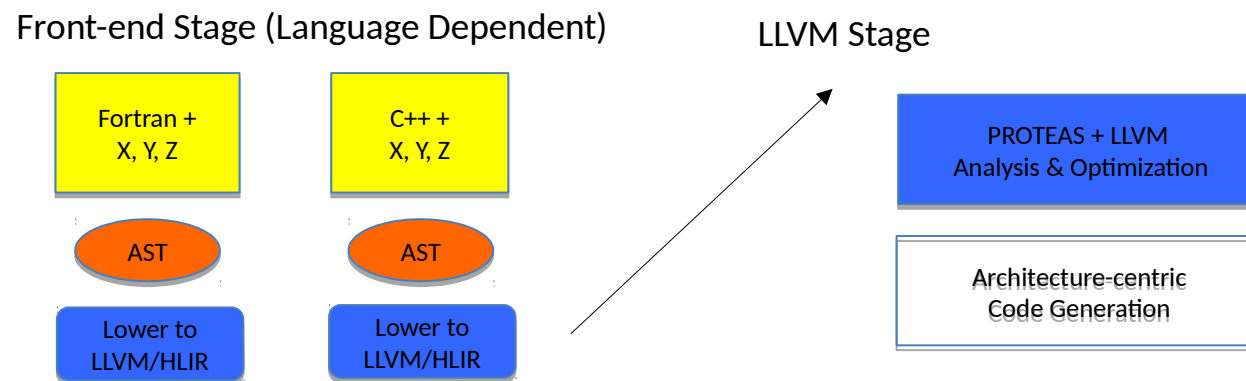
#pragma omp parallel for
  for (i = 0; i < 1000; i++)
    lib_compute(i);

void lib_compute(int x) {
  #pragma omp parallel for
  for (j = 0; j < 1000; j++)
    out[x][j] = compute(in[x][j]);
}
```



PROTEAS – “**PRO**gramming **T**oolchain for **E**merging **A**rchitectures and **S**ystems”

- Developing IR-level representations of parallelism constructs.
- Implementing optimizations on those representations to enable performance-portable programming.
- Exploring how to expose other aspects of modern memory hierarchies (such as NVM).



(Compiler) Optimizations for OpenMP Code

OpenMP is already an abstraction layer. Why can't programmers just write the code optimally?

- Because what is optimal is different on different architectures.
- Because programmers use abstraction layers and may not be able to write the optimal code directly:

```
in library1:  
void foo() {  
    std::for_each(std::execution::par_unseq, vec1.begin(), vec1.end(), ...);  
}
```

```
in library2:  
void bar() {  
    std::for_each(std::execution::par_unseq, vec2.begin(), vec2.end(), ...);  
}
```

```
foo(); bar();
```

(Compiler) Optimizations for OpenMP Code

```
void foo(double * restrict a, double * restrict b, etc.) {  
    #pragma omp parallel for  
    for (i = 0; i < n; ++i) {  
        a[i] = e[i]*(b[i]*c[i] + d[i]) + f[i];  
        m[i] = q[i]*(n[i]*o[i] + p[i]) + r[i];  
    }  
}
```



Split the loop

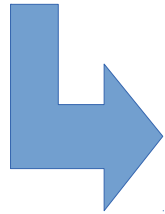


Or should we fuse instead?

```
void foo(double * restrict a, double * restrict b, etc.) {  
    #pragma omp parallel for  
    for (i = 0; i < n; ++i) {  
        a[i] = e[i]*(b[i]*c[i] + d[i]) + f[i];  
    }  
    #pragma omp parallel for  
    for (i = 0; i < n; ++i) {  
        m[i] = q[i]*(n[i]*o[i] + p[i]) + r[i];  
    }  
}
```

(Compiler) Optimizations for OpenMP Code

```
void foo(double * restrict a, double * restrict b, etc.) {  
    #pragma omp parallel for  
        for (i = 0; i < n; ++i) {  
            a[i] = e[i]*(b[i]*c[i] + d[i]) + f[i];  
        }  
    #pragma omp parallel for  
        for (i = 0; i < n; ++i) {  
            m[i] = q[i]*(n[i]*o[i] + p[i]) + r[i];  
        }  
}
```



(we might want to fuse
the parallel regions)

```
void foo(double * restrict a, double * restrict b, etc.) {  
    #pragma omp parallel  
    {  
        #pragma omp for  
            for (i = 0; i < n; ++i) {  
                a[i] = e[i]*(b[i]*c[i] + d[i]) + f[i];  
            }  
        #pragma omp for  
            for (i = 0; i < n; ++i) {  
                m[i] = q[i]*(n[i]*o[i] + p[i]) + r[i];  
            }  
    }  
}
```

(Compiler) Optimizations for OpenMP Code

In order to implement non-trivial parallelism optimizations, we need to move from “early outlining” to “late outlining.”

The optimizer misses:

- Point aliasing information from the parent function
- Loop bounds (and other loop information) from the parent function
- And more...

But perhaps most importantly, it forces us to decide early **how** to lower the parallelism constructs. With some analysis first, after inlining, we can do a much better job (especially when targeting accelerators).

Equivalent of:

```
parallel_for_body(...) {
```

```
parallel_loop(&parallel_for_body, ...);
```

OpenMP does not know
about the loop or the relationship
between the code in the outlined body
and the parent function.

(Compiler) Optimizations for OpenMP Code

An example of where we might generate very different code after analysis...

```
#pragma omp target
{
  // This is a "serial" region on the device.
  foo();
  // So it this.
}

void foo() {
  #pragma omp parallel for
  for (int I = 0; I < N; ++I) { ... }
}
```

On a GPU, you launch some number of SIMT threads: there is no "serial" device execution. To support this general model, we need to generate a complex state machine in each GPU thread. This:

- Wastes resources
- Adds extra synchronization
- Increases register pressure



With late lowering, we could do an analysis to determine that there is no serial code in the parallel region and:

- Generate the (efficient) code the user expects.
- Analyze memory accesses and potentially use local/shared/texture memory.

(Compiler) Optimizations for OpenMP Code

In order to implement non-trivial parallelism optimizations, we need to move from “early outlining” to “late outlining.”

These markers are currently being designed. They might be intrinsics, perhaps also using operand bundles, but also require several special properties:

- `alloca` instructions inside the region must stay inside the region.
- The region markers must appear to capture/access pointers used inside the region (regions might run more than once, or after function returns, etc.).
- For loops, prevent the introduction of new loop-carried dependencies (duplicate induction variables, etc.).
- UB if exception-triggered unwinding occurs.

LVM IR equivalent of:

```
@begin_parallel_for  
for (...) {  
    ...  
}
```

...
__run_parallel_for

...
body of:

```
body(...) {  
    ...  
}
```

```
void foo() {  
    __run_parallel_loop(&parallel_for_body, ...);  
}
```

If we don't also handle C++ lambdas using this kind of mechanism, we won't get the full benefit!

LLVM Optimizer

ARES/HeteroIR – Predecessors to PROTEAS

- Developed a high-level IR targeted by OpenACC (and other models).
- <http://ft.ornl.gov/research/openarc>
- <https://github.com/lanl/ares>

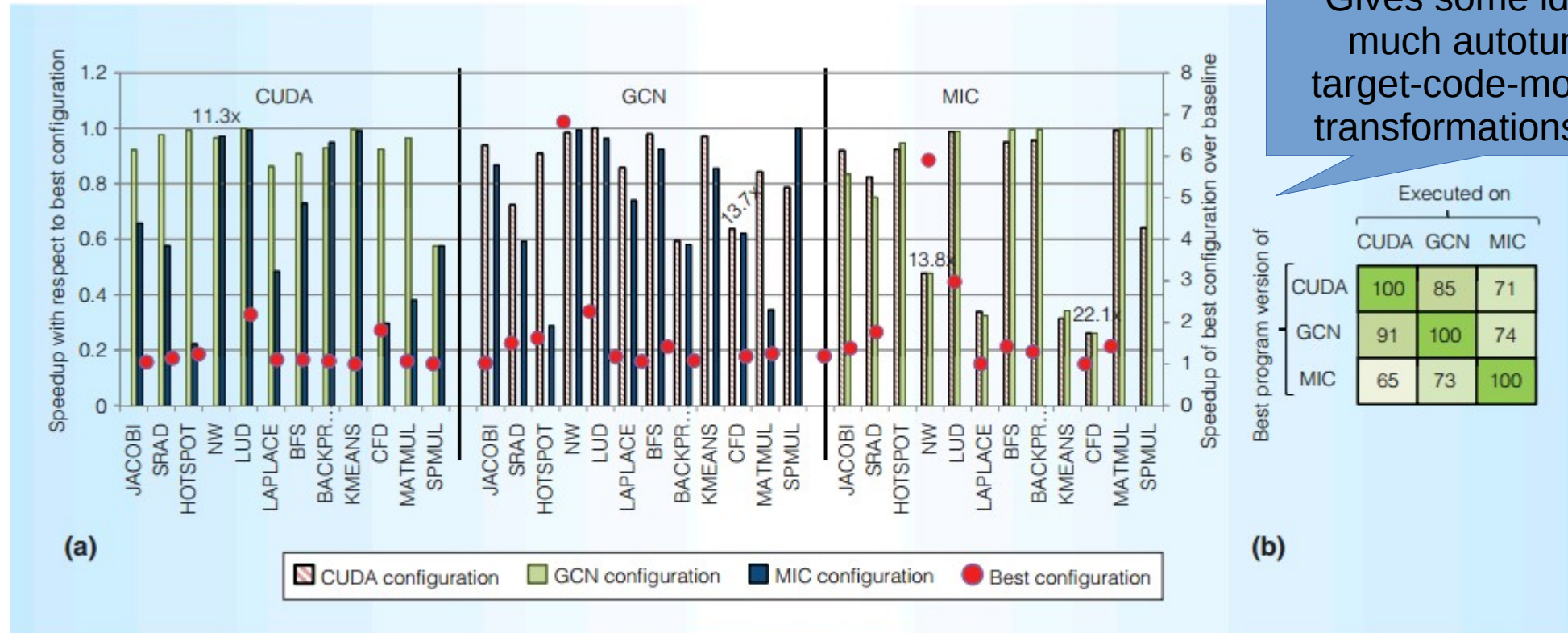
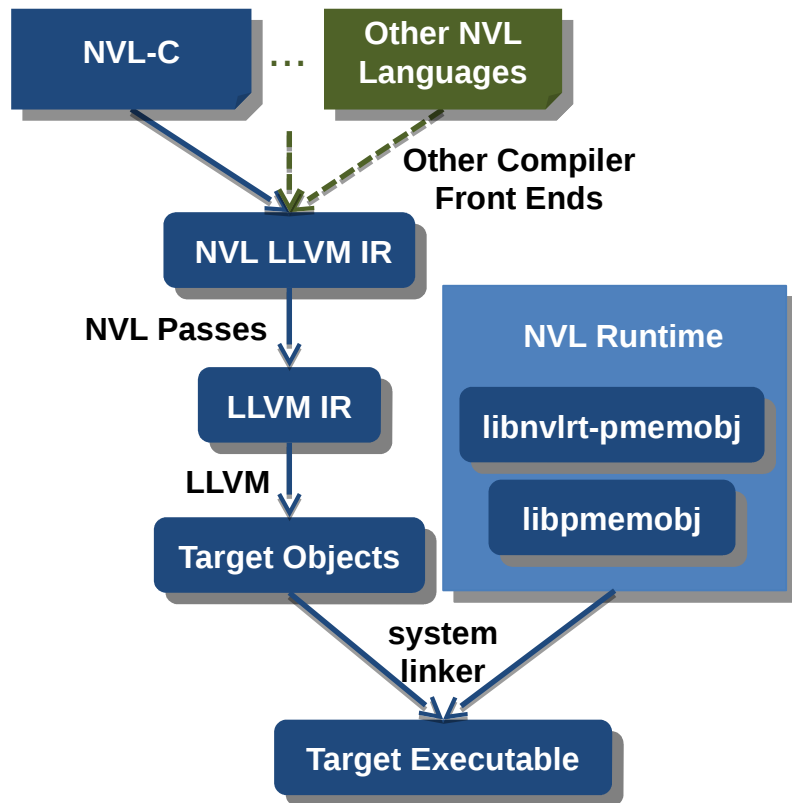


Figure 4. Performance portability of the OpenACC programming model as obtained by OpenARC. (a) Performance portability evaluation. (b) Performance portability achieved across benchmarks. Better performance portability is achieved among the GPU architectures.

NVL-C – Predecessors to PROTEAS

- Experimenting with how to use NVM
- <http://ft.ornl.gov/research/nvl-c>



Extensions to C
with transactions.

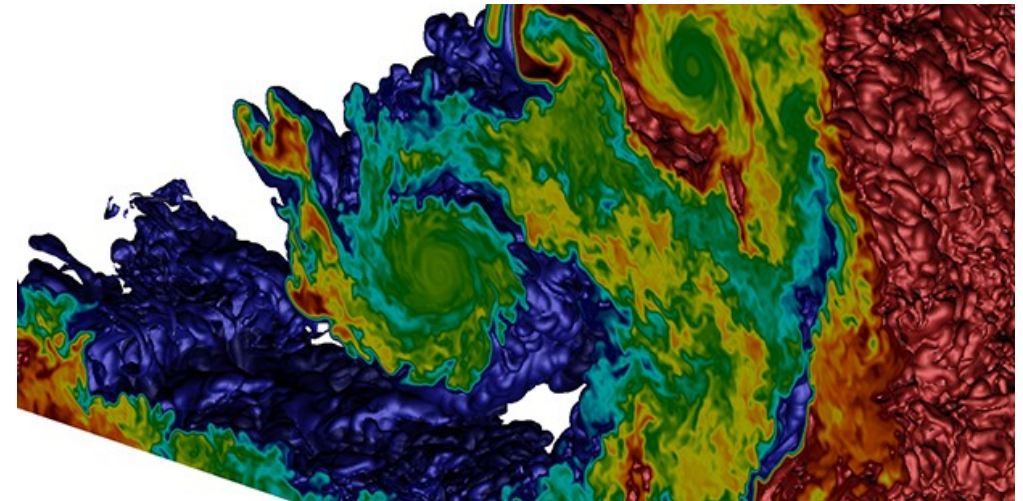
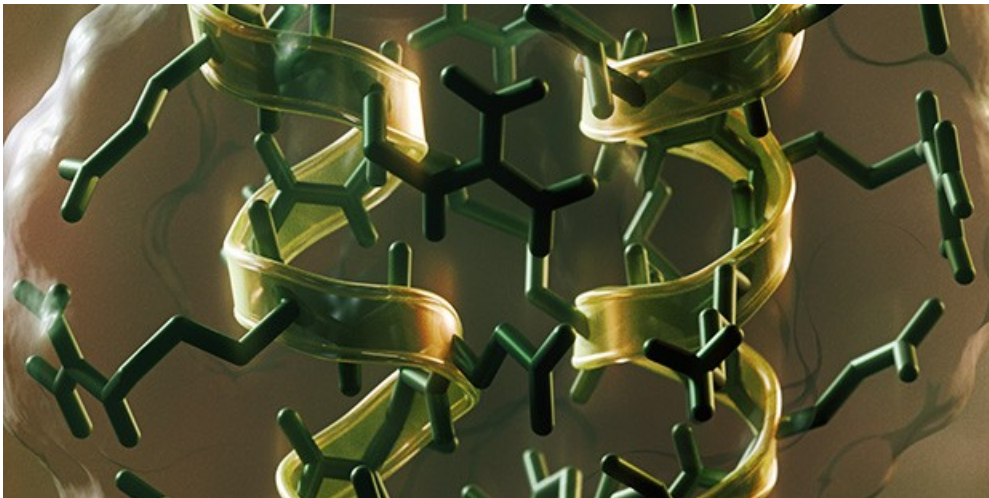
```
#include <nvl.h>
struct list {
    int value;
    nvl struct list *next;
};
void remove(int k) {
    nvl_heap_t *heap
        = nvl_open("foo.nvl");
    nvl struct list *a
        = nvl_get_root(heap, struct list);
    #pragma nvl atomic
    while (a->next != NULL) {
        if (a->next->value ==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

Will high-performance NVM fundamentally change the way that people write software?

(Work by Joel E. Denny, Seyong Lee, and Jeffrey S. Vetter)

In Conclusion...

- Future HPC hardware will be diverse.
- Work is needed on applications, abstraction libraries, and compilers (and related tools).
- Enhancing LLVM to understand parallelism provides an enabling underlying technology for performance portability!



Acknowledgments

- The LLVM community (including our many contributing vendors)
- ALCF, ANL, and DOE
- ALCF is supported by DOE/SC under contract DE-AC02-06CH11357

