

# Programming your GPU with OpenMP

OpenMP®

Tim Mattson

The Bonn2025 repository is a temporary repository. I will delete it in a month or two.

For long term access to my materials, go to my "archival storage" repositories.

For my parallel computing courses (hands-on):

```
git clone https://github.com/tgmattso/ParProgForPhys.git
```

For more general computer Science lectures for people interested in the sciences:

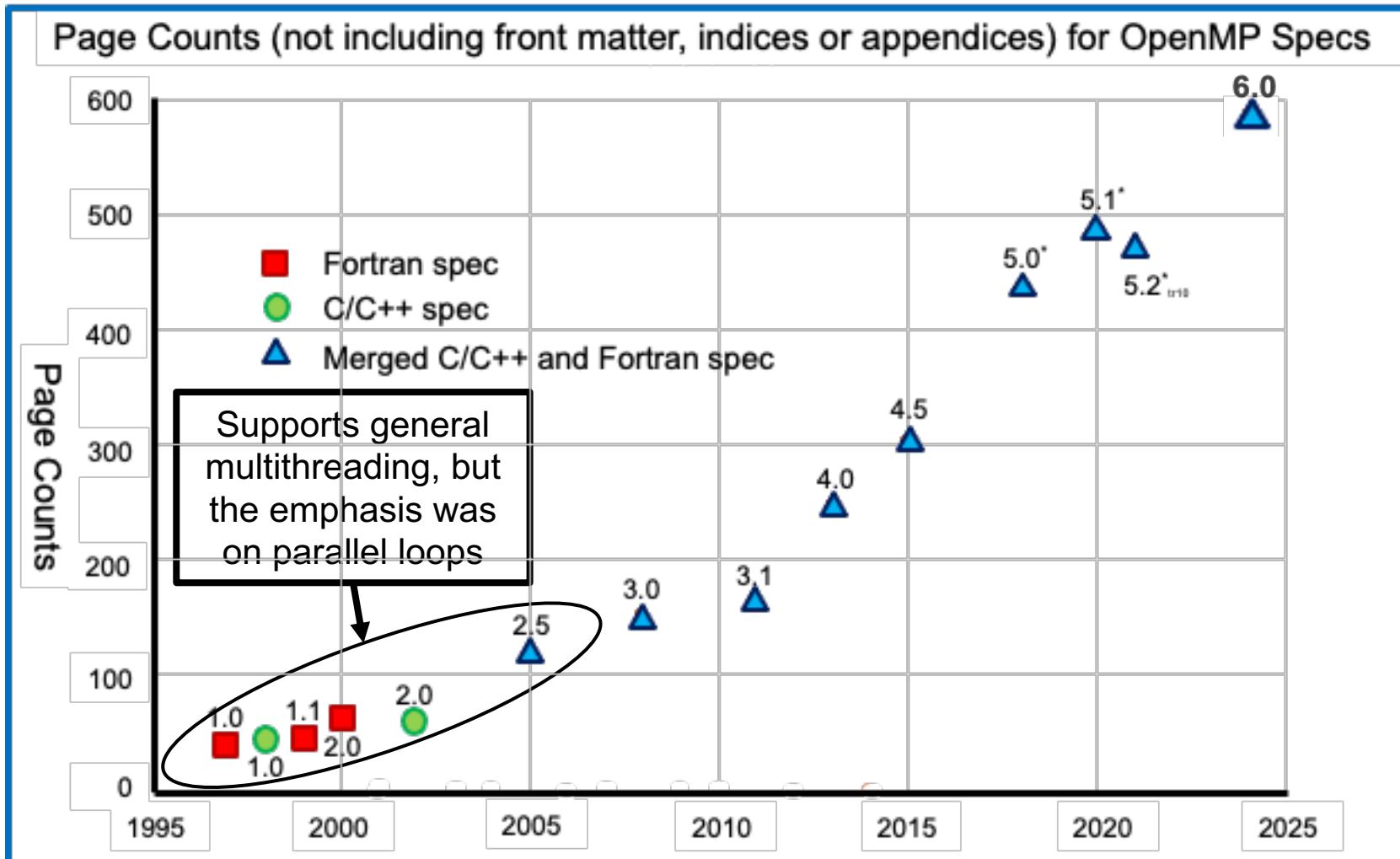
```
git clone https://github.com/tgmattso/CompSciForPhys.git
```

```
git clone https://github.com/tgmattso/Bonn2025.git
```

Much of this content was created with Tom Deakin and Simon McIntosh-Smith of the University of Bristol

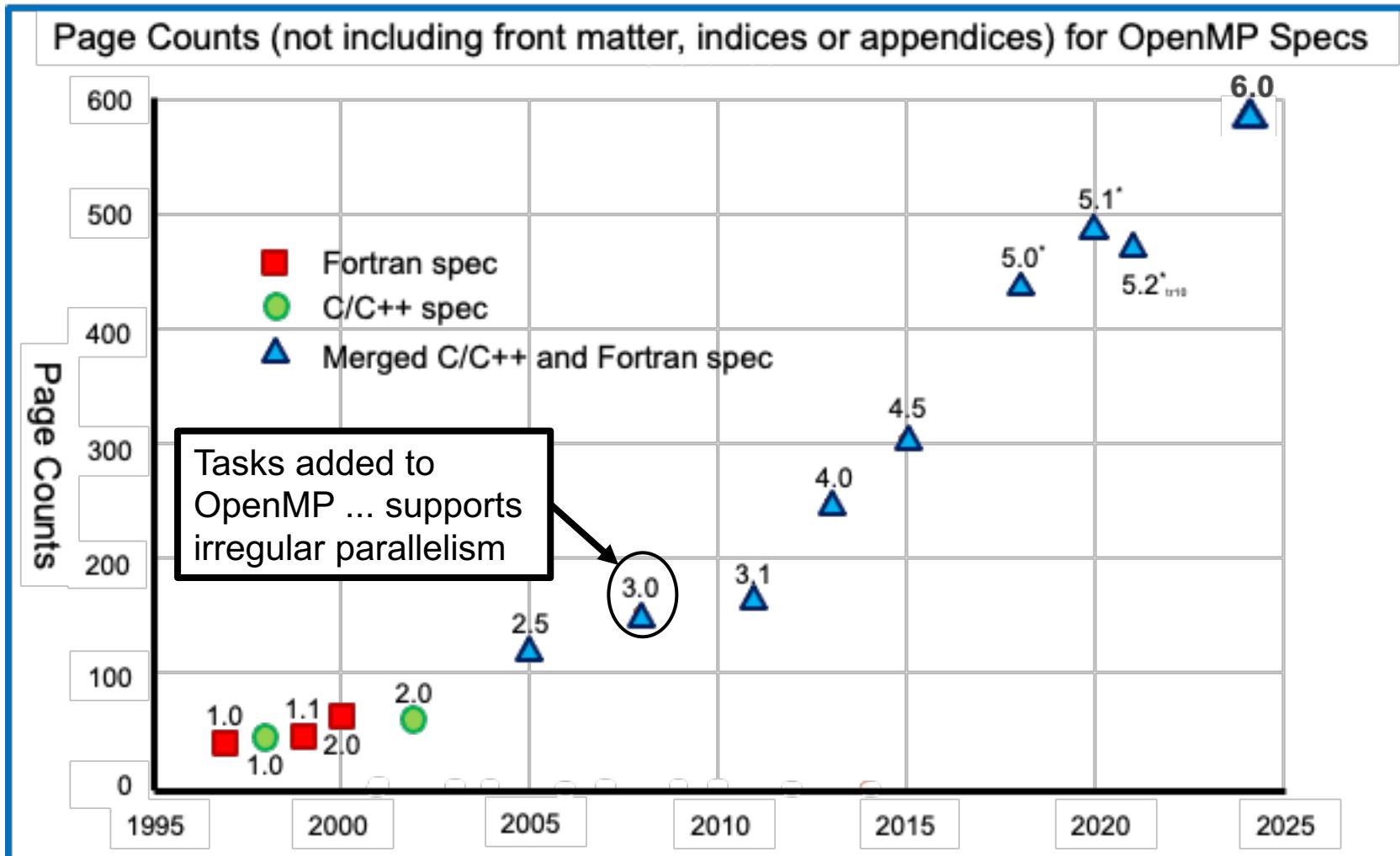
# The growth of complexity in OpenMP

- OpenMP started out in 1997 as a simple interface for the application programmers more versed in their area of science than computer science.
- The complexity has grown considerably over the years!



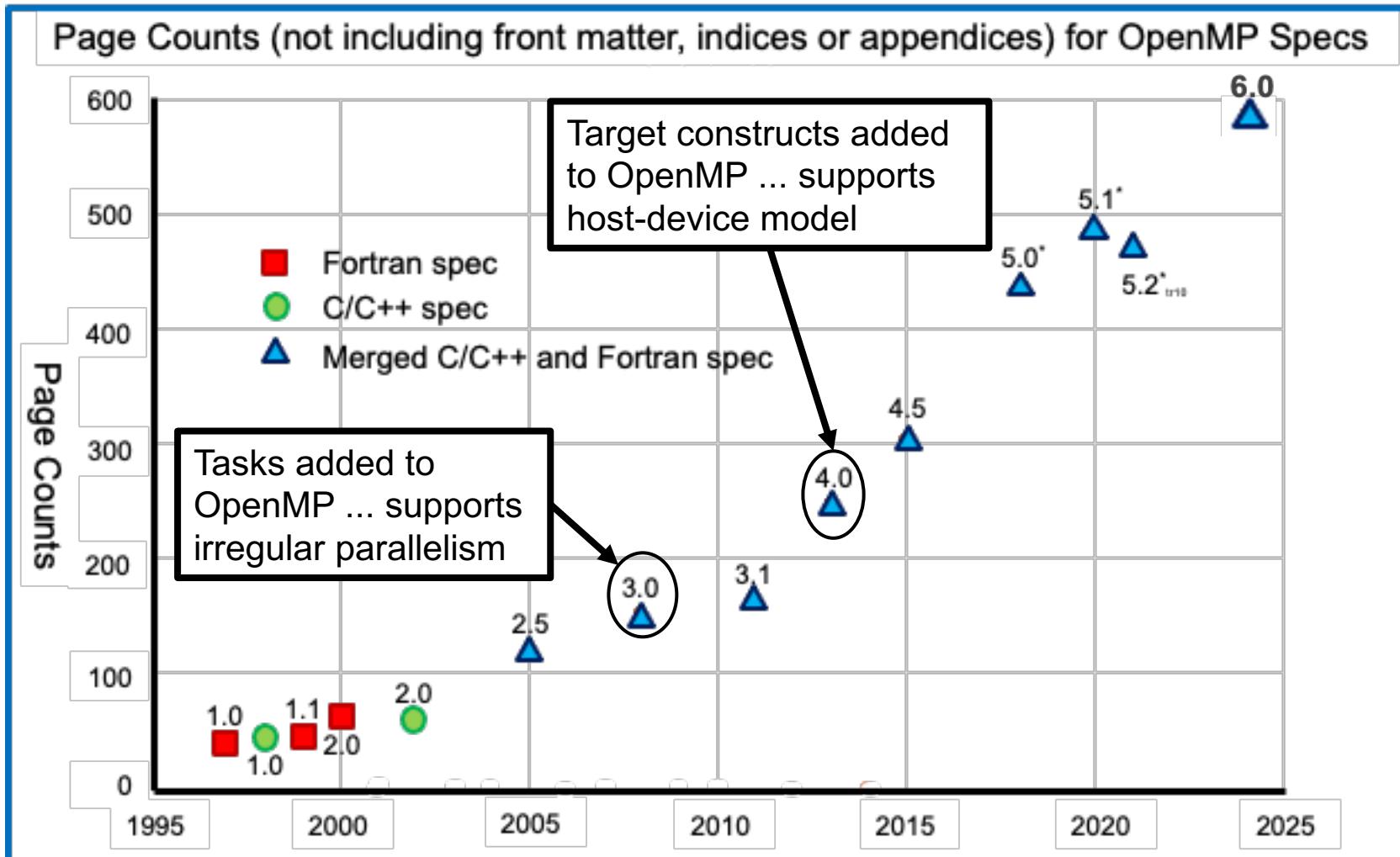
# The growth of complexity in OpenMP

- OpenMP started out in 1997 as a simple interface for the application programmers more versed in their area of science than computer science.
- The complexity has grown considerably over the years!

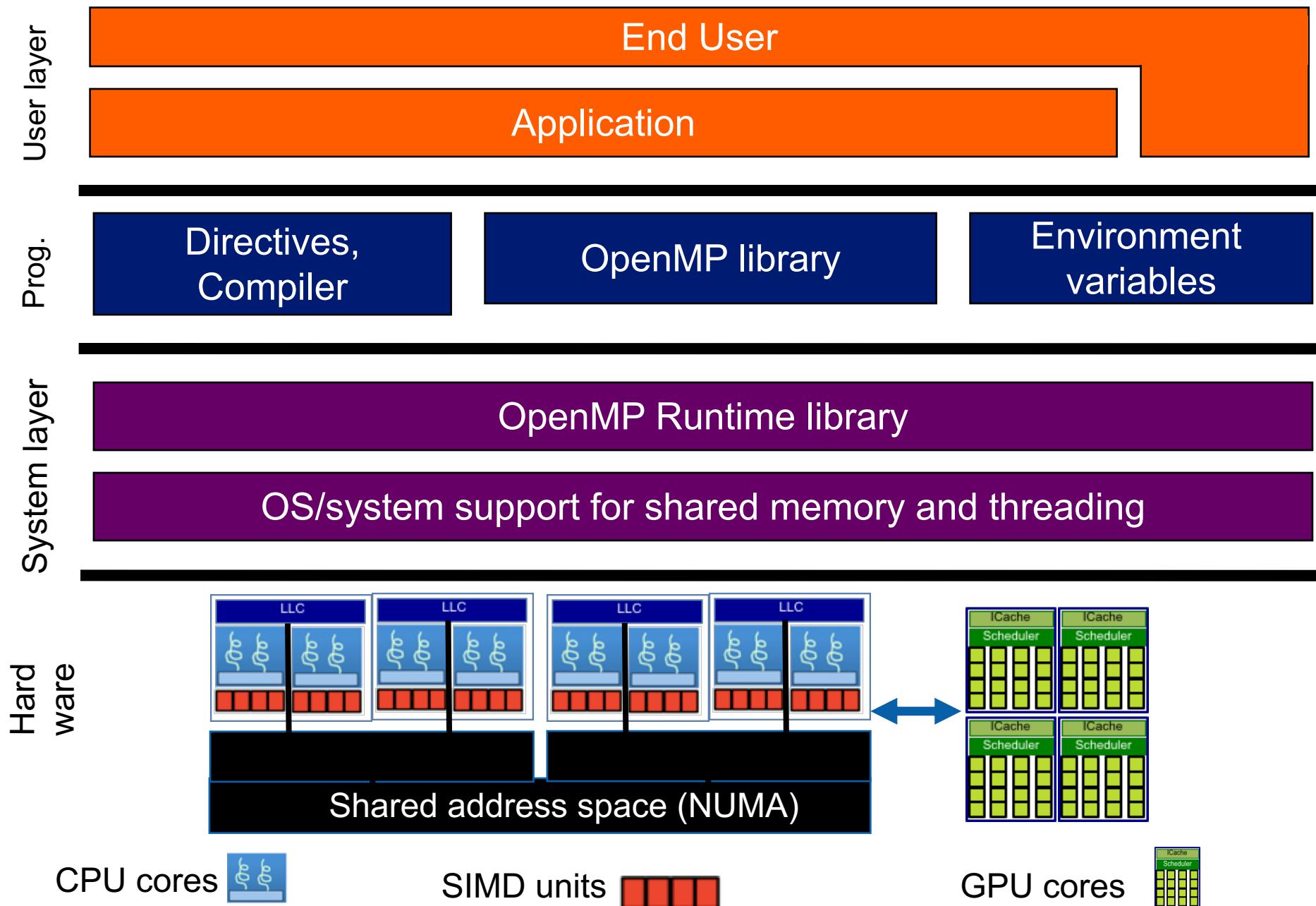


# The growth of complexity in OpenMP

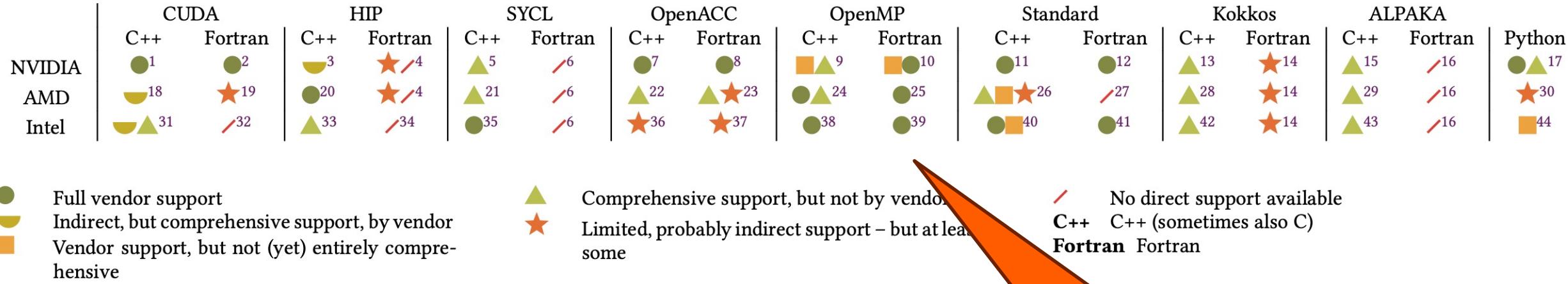
- OpenMP started out in 1997 as a simple interface for the application programmers more versed in their area of science than computer science.
- The complexity has grown considerably over the years!



# OpenMP basic definitions: the solution stack

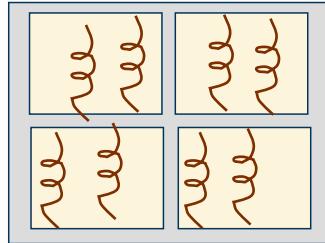


# Why is OpenMP so important?



OpenMP is the **only** model  
with **full** support  
from **all** vendors  
for C/C++ and Fortran

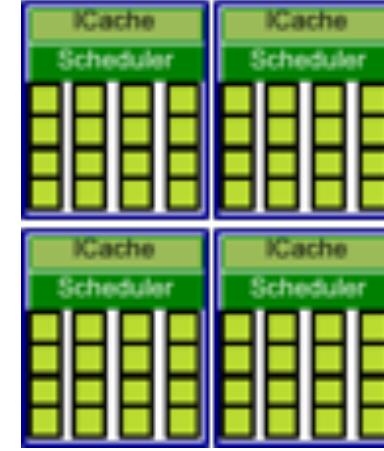
# Hardware is diverse ... and its only getting worse!!!



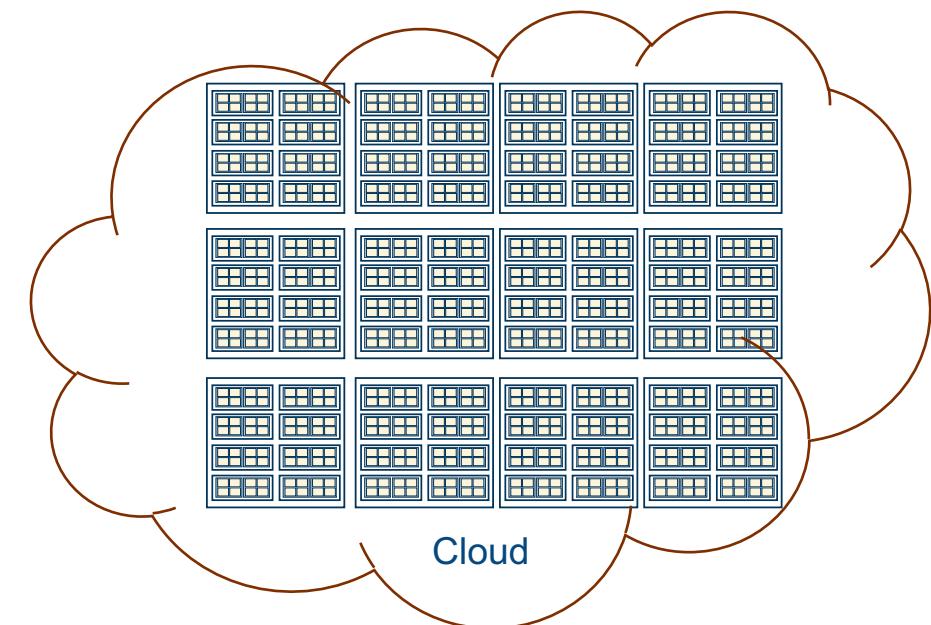
CPU



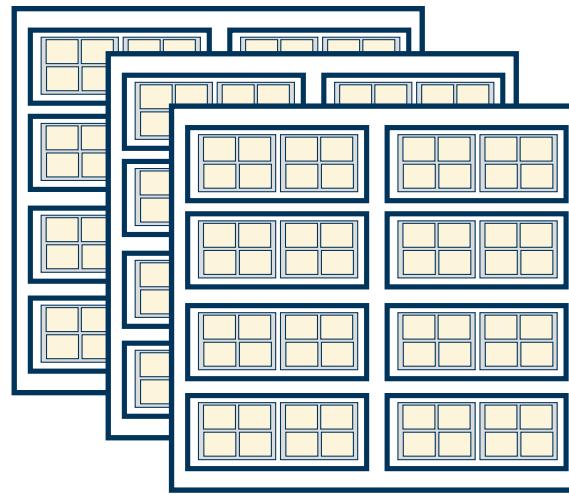
SIMD/Vector



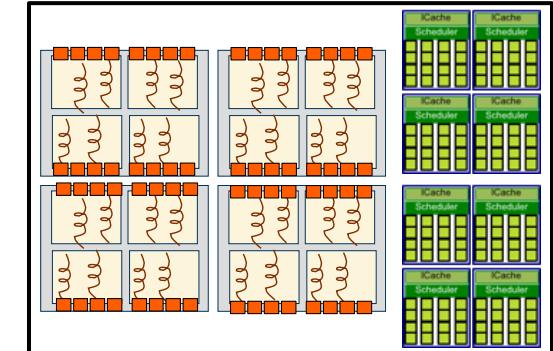
GPU



Cloud



Cluster



Heterogeneous node

# The Big Three

- In HPC, 3 programming environments dominate ... covering the major classes of hardware.
  - **MPI**: distributed memory systems ... though it works nicely on shared memory computers.
  - **OpenMP**: Shared memory systems ... more recently, GPGPU too.
  - **CUDA, OpenCL, Sycl, OpenACC, OpenMP**... : GPU programming (use CUDA if you don't mind locking yourself to a single vendor ... it is a really nice programming model)
- Even if you don't plan to spend much time programming with these systems ... a well rounded HPC programmer should know what they are and how they work.

# How do we execute code on a GPU: The SIMT model (Single Instruction Multiple Thread)

1. Turn kernel code into a scalar work-item

```
// Compute sum of order-N matrices: C = A + B
void __global__
matAdd (float* a, float* b, float* c, int N) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;
    if (i < N && j < N) c[i][j] = a[i][j] + b[i][j];
}

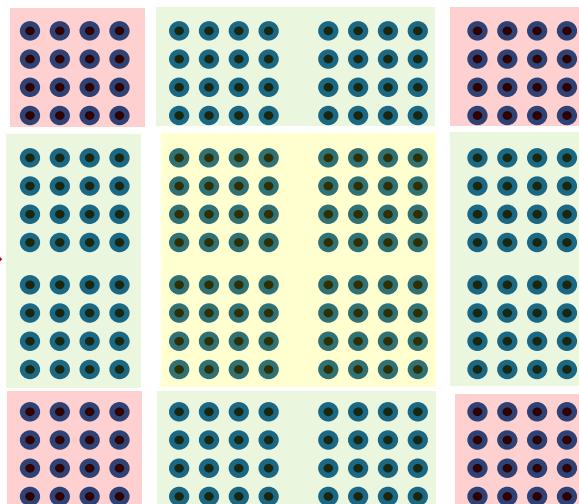
int main () {
    int N = ... ;
    float *a, *b, *c;
    cudaMalloc (&a, sizeof(float) * N);
    // ... allocate other arrays (b and c)
    // and fill with data

    // define threadBlocks and the Grid
    dim3 dimBlock(4,4);
    dim3 dimGrid(4,4);

    // Launch kernel on Grid
    matAdd <<< digGrid, dimBlock>>> (a, b, c, N);
}
```

This is CUDA code

2. Map work-items onto an N dim index space.



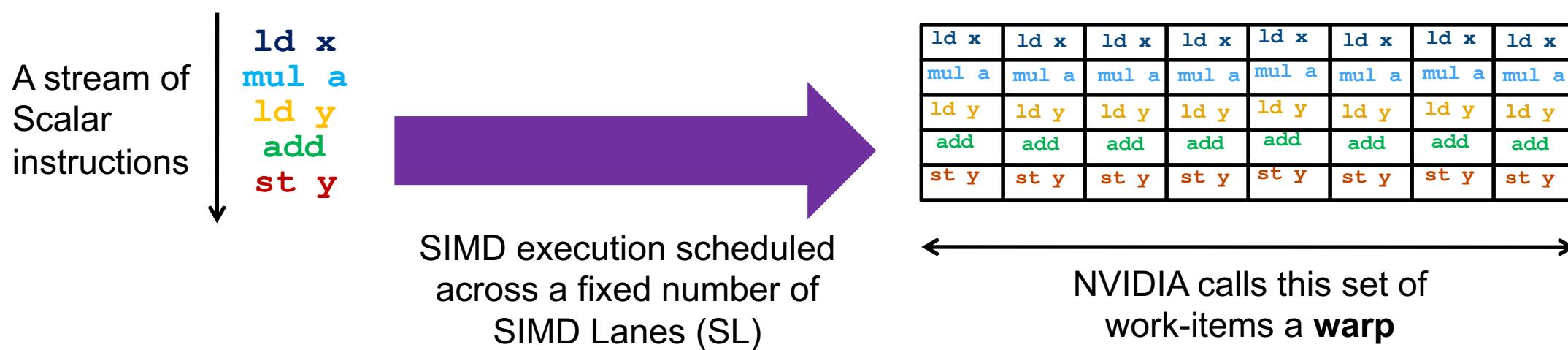
3. Map data structures onto the same index space

4. Run on hardware designed around the same SIMT execution model

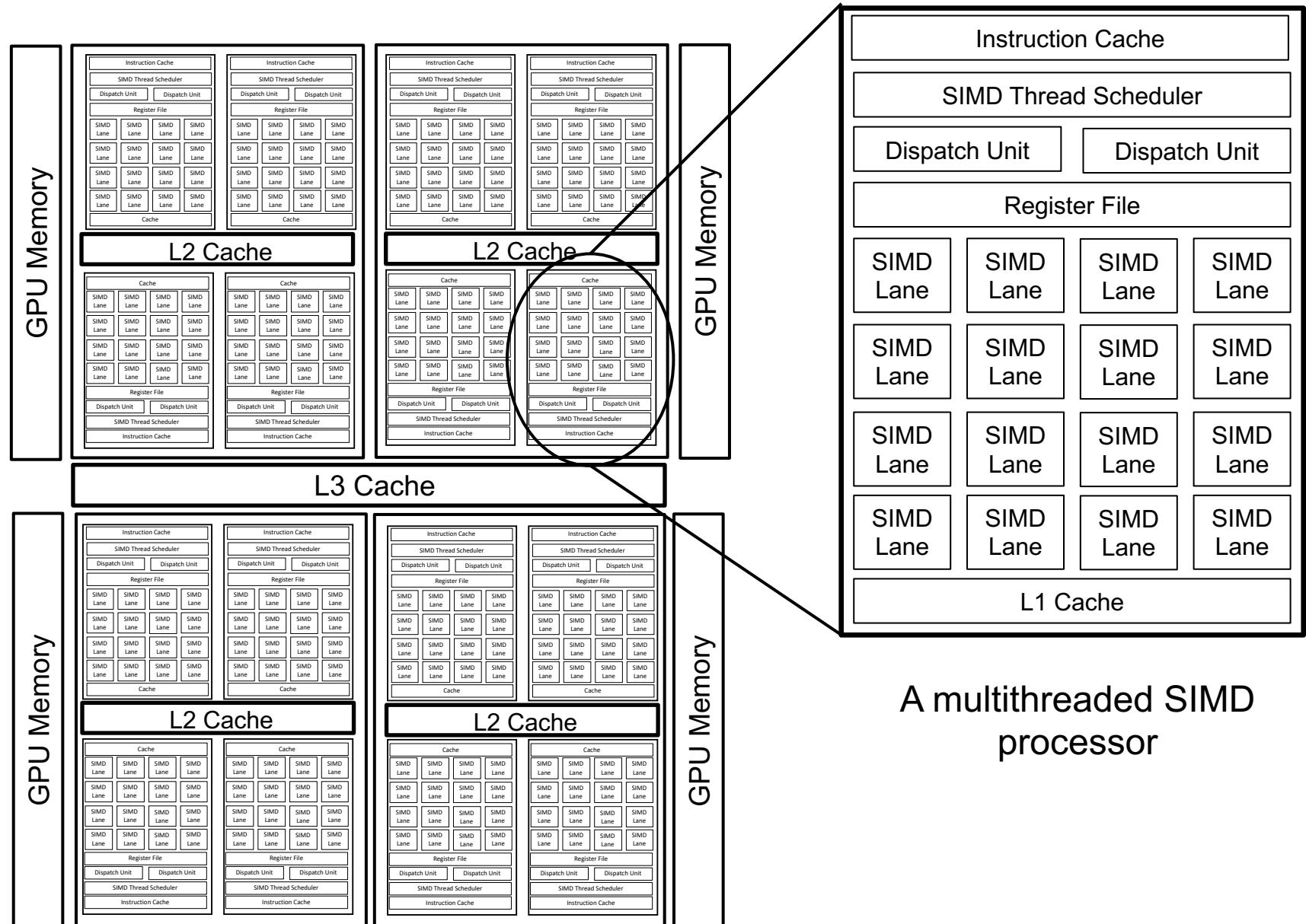


# SIMT: One instruction stream maps onto many SIMD lanes

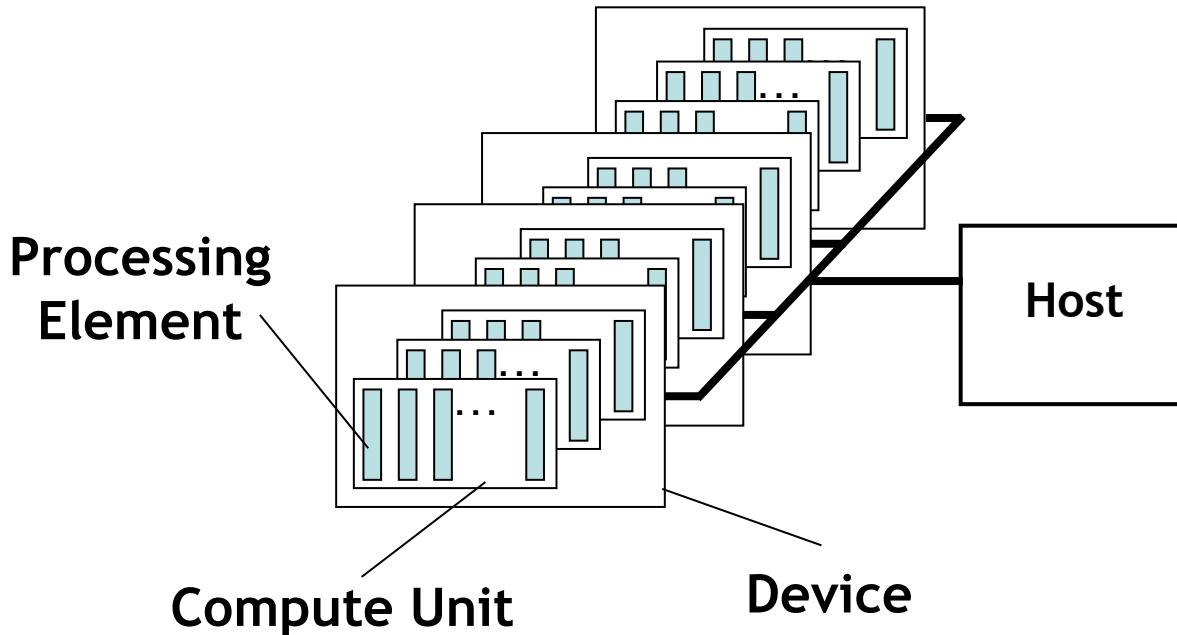
- SIMT model: Individual scalar instruction streams are grouped together for SIMD execution on hardware



# A Generic GPU (following Hennessy and Patterson)

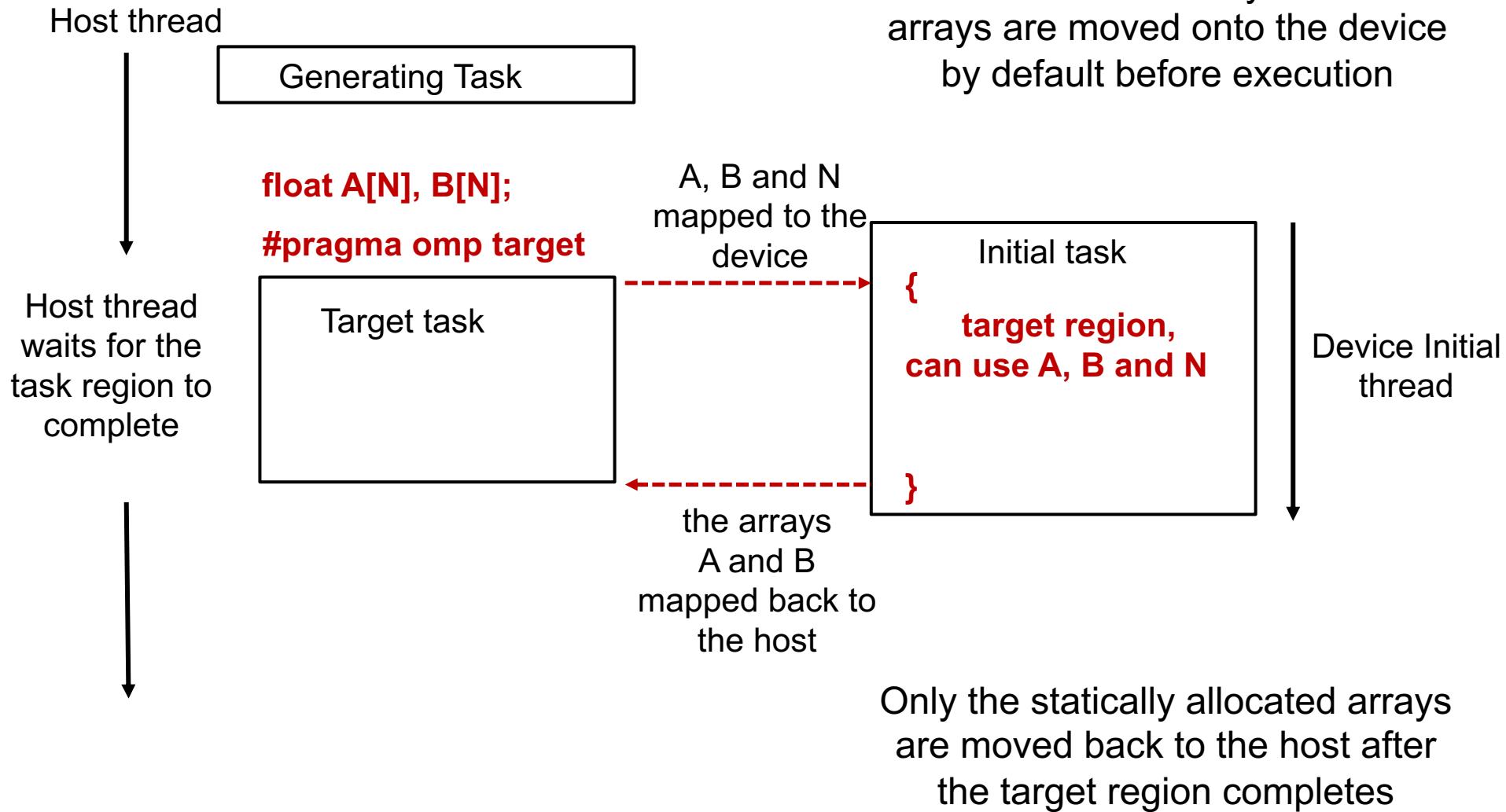


# A Generic Host/Device Platform Model



- One **Host** and one or more **Devices**
  - Each Device is composed of one or more Compute Units
  - Each Compute Unit is divided into one or more **Processing Elements**
- Memory divided into **host memory** and **device memory**

# Running code on the GPU: The target construct and default data movement



# Default Data Sharing: example

```
int main(void) {  
    int N = 1024;  
    double A[N], B[N];  
  
    #pragma omp target  
    {  
  
        for (int ii = 0; ii < N; ++ii) {  
  
            A[ii] = A[ii] + B[ii];  
  
        }  
    } // end of target region  
}
```

1. Variables created in host memory.

2. Scalar **N** and stack arrays **A** and **B** are copied *to* device memory. Execution transferred to device.

3. **ii** is **private** on the device as it's declared within the target region

4. Execution on the device.

5. stack arrays **A** and **B** are copied *from* device memory back to the host. Host resumes execution.

# Now let's run code in parallel on the device

```
int main(void) {  
    int N = 1024;  
    double A[N], B[N];  
  
    #pragma omp target  
    {  
        #pragma omp loop  
        for (int ii = 0; ii < N; ++ii) {  
  
            A[ii] = A[ii] + B[ii];  
  
        }  
    } // end of target region  
}
```

The loop construct tells the compiler:  
*"this loop will execute correctly if the loop iterations run in any order. You can safely run them concurrently. And the loop-body doesn't contain any OpenMP constructs. So do whatever you can to make the code run fast"*

The loop construct is a declarative construct. You tell the compiler what you want done but you DO NOT tell it how to "do it". This is new for OpenMP

# Exercise: Parallel vector addition on a GPU

- Make a copy of your parallel vadd.c program for a CPU (i.e. save the CPU version)
  - vadd.c Adds together two arrays, element by element:  $\text{for}(i=0;i<N;i++) c[i]=a[i]+b[i];$
- Parallelize your vadd program for a GPU
- Time it for large N and save the result. How does it compare to the CPU version?
  - double omp\_get\_wtime();
  - #pragma omp target
  - #pragma omp loop

For tiny little programs, OpenMP may opt to run the code on the host. You can force the OpenMP runtime to use the GPU by setting the OMP\_TARGET\_OFFLOAD environment variable

> OMP\_TARGET\_OFFLOAD=MANDATORY ./a.out

git clone <https://github.com/tgmattso/Bonn2025.git>

# Solution: Simple vector add in OpenMP on GPU

```
int main()
{
    float a[N], b[N], c[N], res[N];
    int err=0;

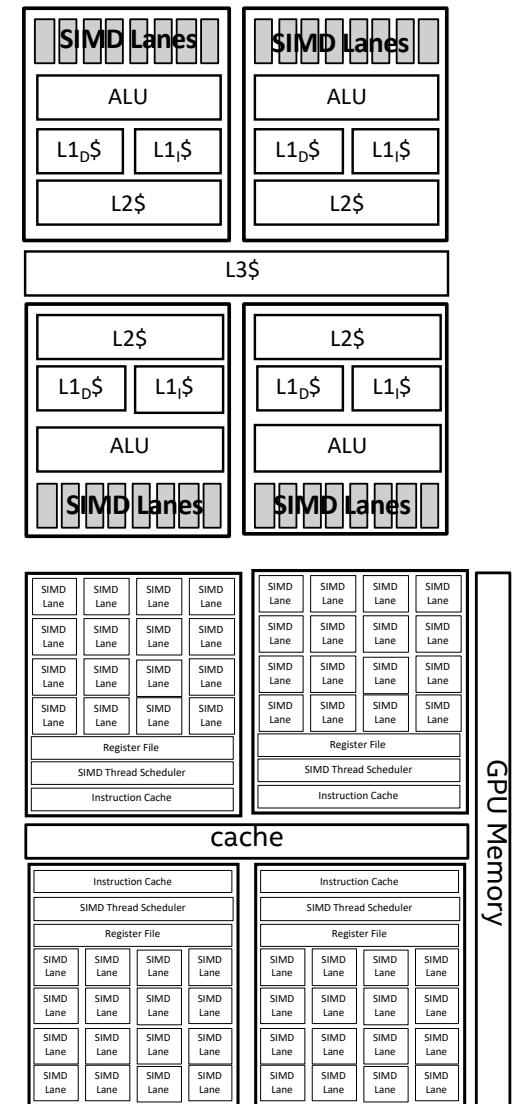
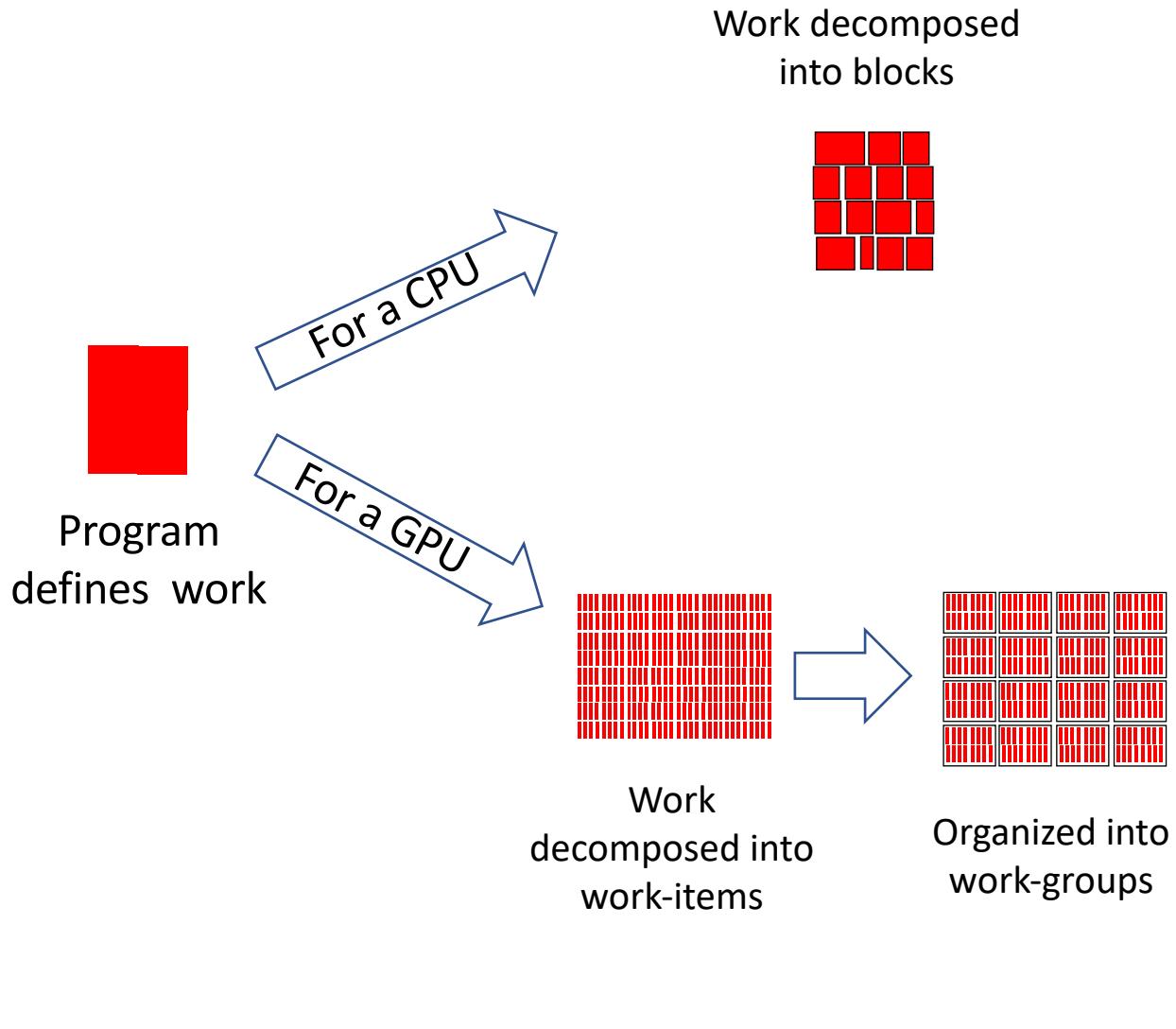
    // fill the arrays
    #pragma omp parallel for
    for (int i=0; i<N; i++) {
        a[i] = (float)i;
        b[i] = 2.0*(float)i;
        c[i] = 0.0;
        res[i] = i + 2*i;
    }

    // add two vectors
    #pragma omp target
    #pragma omp loop
    for (int i=0; i<N; i++) {
        c[i] = a[i] + b[i];
    }

    // test results
    #pragma omp parallel for reduction(+:err)
    for(int i=0;i<N;i++) {
        float val = c[i] - res[i];
        val = val*val;
        if(val>TOL) err++;
    }
    printf("vectors added with %d errors\n", err);
    return 0;
}
```

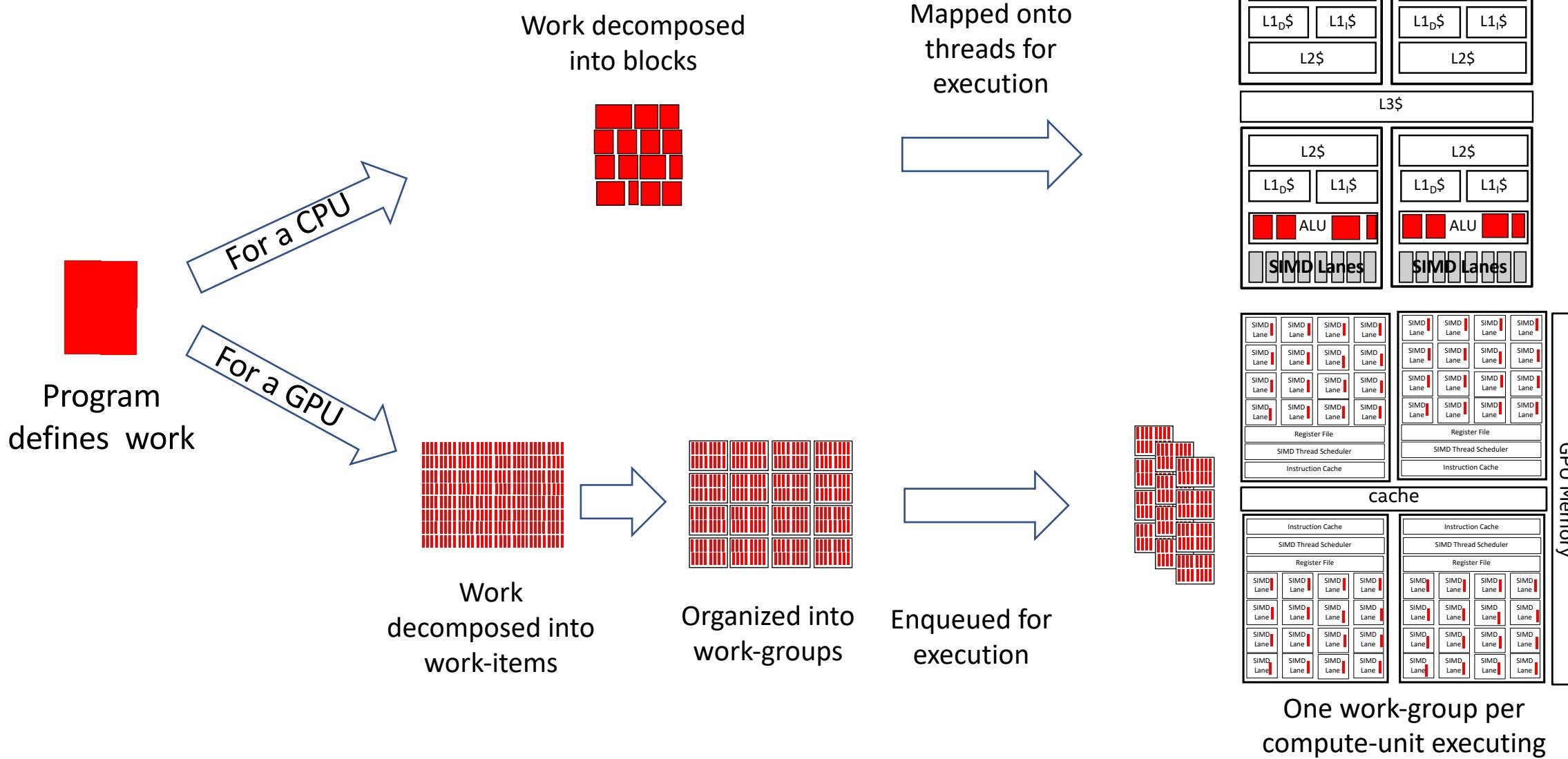
**Let's compare/contrast concurrency on a  
CPU and a GPU**

# Executing a program on CPUs and GPUs



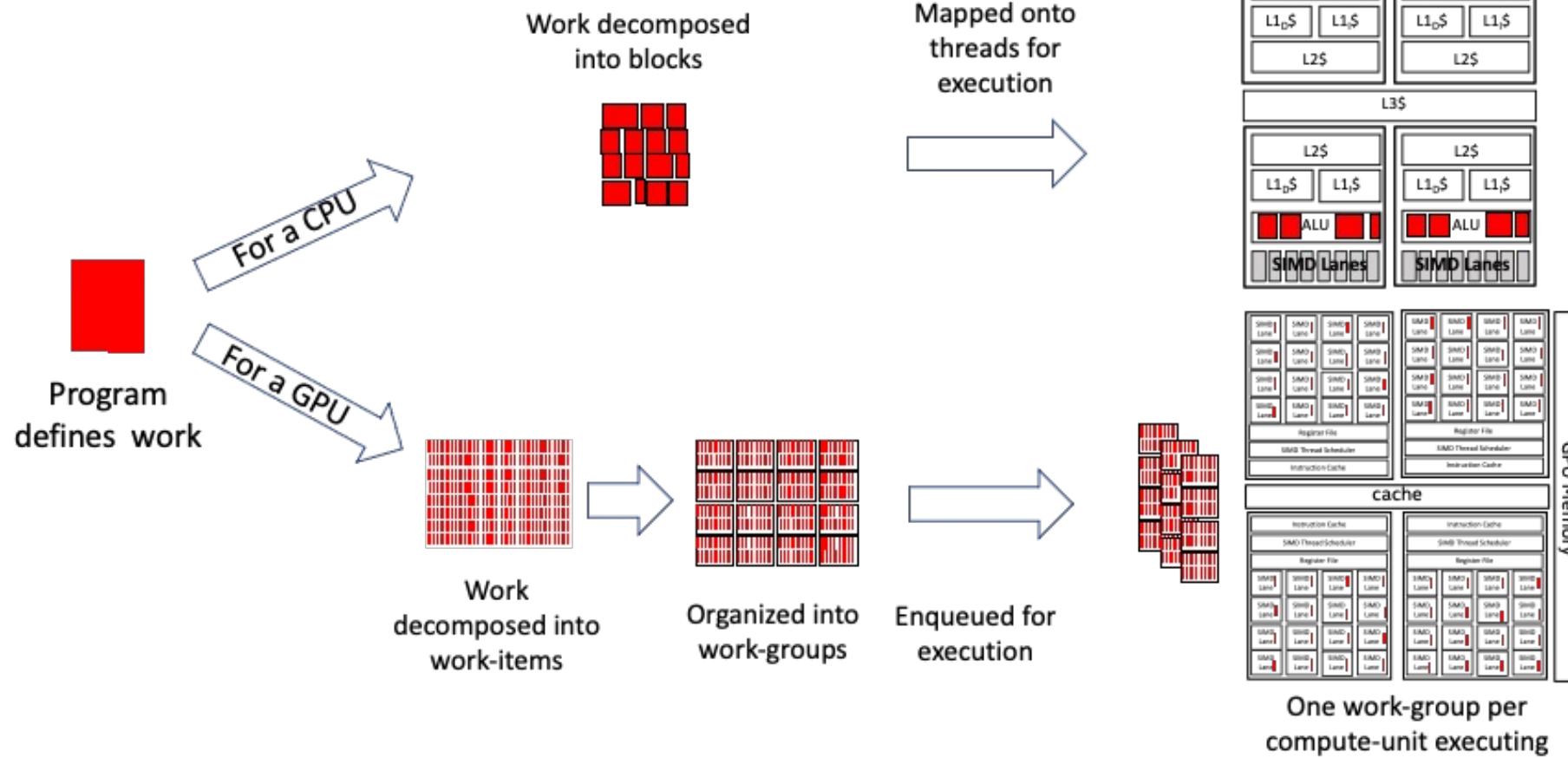
One work-group per  
compute-unit executing

# Executing a program on CPUs and GPUs



# CPU/GPU execution models

## Executing a program on CPUs and GPUs

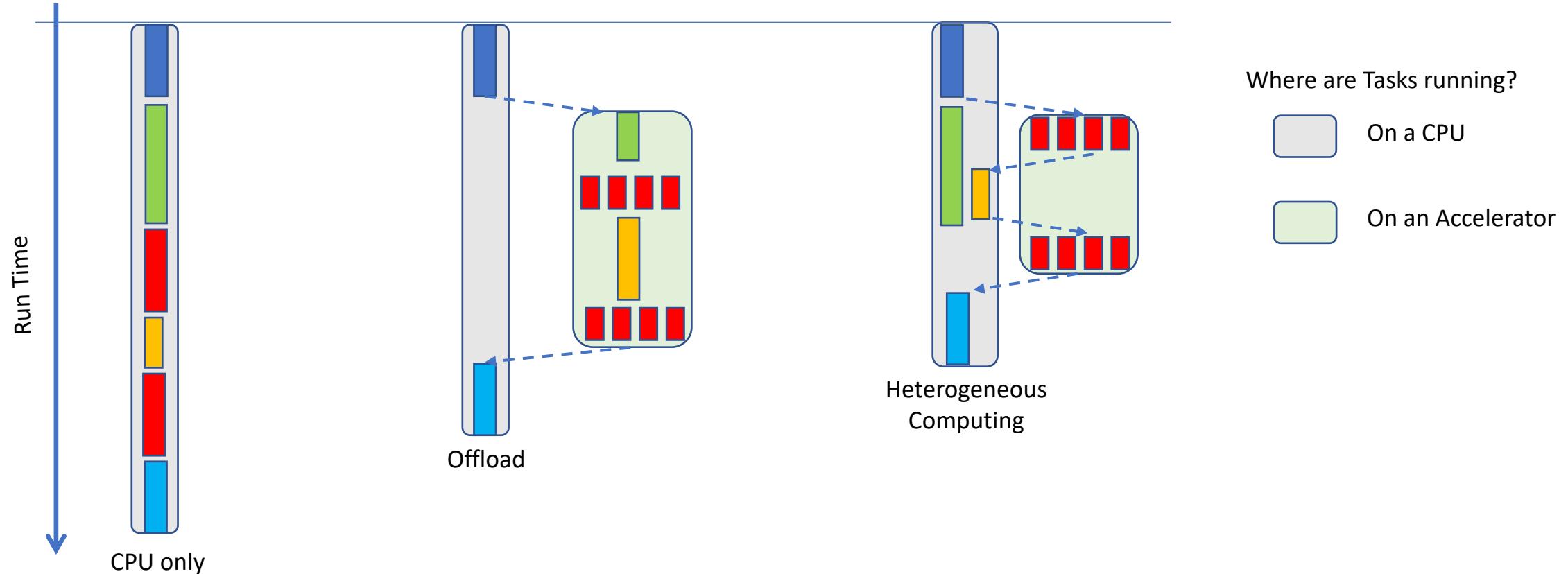


For a CPU, the threads are all active and able to make forward progress.

For a GPU, any given work-group might be in the queue waiting to execute.

# No single processor is best at everything

- The idea that you should move everything to the GPU makes no sense
- **Heterogeneous Computing:** Run sub-problems in parallel on the hardware best suited to them.



**Implicit data movement covers a small subset of  
the cases you need in a real program.**

**To be more general ... we need to manage data  
movement explicitly**

# Explicit data movement

- Previously, we described the rules for *implicit* data movement.
- We can *explicitly* control the movement of data using the **map** clause.
- **Data allocated on the heap needs to be explicitly copied to/from the device:**

```
int main(void) {  
    int ii=0, N = 1024;  
    int* A = (int *)malloc(sizeof(int)*N);  
  
#pragma omp target  
{  
    // N, ii and A all exist here  
    // The data that A points to (*A , A[ii]) DOES NOT exist here!  
}
```

# Moving data with the map clause

```
int main(void) {  
    int N = 1024;  
    int* A = malloc(sizeof(int)*N);  
  
    #pragma omp target map(A[0:N])  
    {  
        // N, ii and A all exist here  
        // The data that A points to DOES exist here!  
    }  
}
```

Default mapping  
**map(tofrom: A[0:N])**

Copy at start and end of  
**target** region.

# Controlling data movement

```
int i, a[N], b[N], c[N];  
#pragma omp target map(to:a,b) map(tofrom:c)
```

Data movement  
defined from the  
*host* perspective.

- The various forms of the map clause
  - **map(to:list)**: On entering the region, variables in the list are initialized on the device using the original values from the host (host to device copy).
  - **map(from:list)**: At the end of the target region, the values from variables in the list are copied into the original variables on the host (device to host copy). On entering the region, the initial value of the variables on the device is not initialized.
  - **map(tofrom:list)**: the effect of both a map-to and a map-from (host to device copy at start of region, device to host copy at end).
  - **map(alloc:list)**: On entering the region, data is allocated and uninitialized on the device.
  - **map(list)**: equivalent to **map(tofrom:list)**.

# Vector add with dynamic memory on GPU

```
int main()
{
    float *a      = malloc(sizeof(float) * N);
    float *b      = malloc(sizeof(float) * N);
    float *c      = malloc(sizeof(float) * N);
    float *res = malloc(sizeof(float) * N);
    int err=0;

    // fill the arrays <<<code not shown>>>

    // add two vectors
#pragma omp target map(to: a[0:N],b[0:N]) map (tofrom: c[0:N])
#pragma omp loop
for (int i=0; i<N; i++) {
    c[i] = a[i] + b[i];
}

    // test results <<<code not shown>>>

    printf("vectors added with %d errors\n", err);
    return 0;
}
```

# Commonly used clauses on target and loop constructs

- The basic construct\* is:

`#pragma omp target [clause[,]clause]...]`

`#pragma omp loop [clause[,]clause]...]`

*for-loops*

- The most commonly used clauses are:

- **map(to | from | tofrom list)** ← default is tofrom

- **private(list)** **firstprivate(list)** **lastprivate(list)** **shared(list)**

- behave as data environment clauses in the rest of OpenMP, but note values are only created or copied into the region, not back out “at the end”.

- **reduction(reduction-identifier : list)**

- behaves as in the rest of OpenMP

- **collapse(n)**

- Combines loops before the distribute directive splits up the iterations between teams

**Going beyond simple vector addition ...**

**Using OpenMP for GPU application  
programming ... the heat diffusion problem**

# 5-point stencil: the heat program

- The heat equation models changes in temperature over time.

$$\frac{\partial u}{\partial t} - \alpha \nabla^2 u = 0$$

- We'll solve this numerically on a computer using an explicit **finite difference** discretisation.
- $u = u(t, x, y)$  is a function of space and time.
- Partial differentials are approximated using diamond difference formulae:

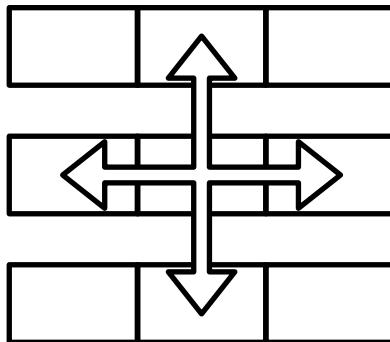
$$\frac{\partial u}{\partial t} \approx \frac{u(t+1, x, y) - u(t, x, y)}{dt}$$

$$\frac{\partial^2 u}{\partial x^2} \approx \frac{u(t, x+1, y) - 2u(t, x, y) + u(t, x-1, y)}{dx^2}$$

- Forward finite difference in time, central finite difference in space.

# 5-point stencil: the heat program

- Given an initial value of  $u$ , and any boundary conditions, we can calculate the value of  $u$  at time  $t+1$  given the value at time  $t$ .
- Each update requires values from the north, south, east and west neighbours only:



- Computation is essentially a weighted average of each cell and its neighbouring cells.
- If on a boundary, look up a boundary condition instead.

# Method of Manufactured Solution

- Stencil codes are notoriously difficult to **know** if the answer is “correct”.
- Analytic solutions hard to come by:
  - It’s why you’re using a computer to solve the equation approximately after all!
- Method of Manufactured Solution (MMS) is a way to help determine if the code does the correct thing.
- An approach often used to find errors in CFD codes and check convergence properties.

# Method of Manufactured Solution

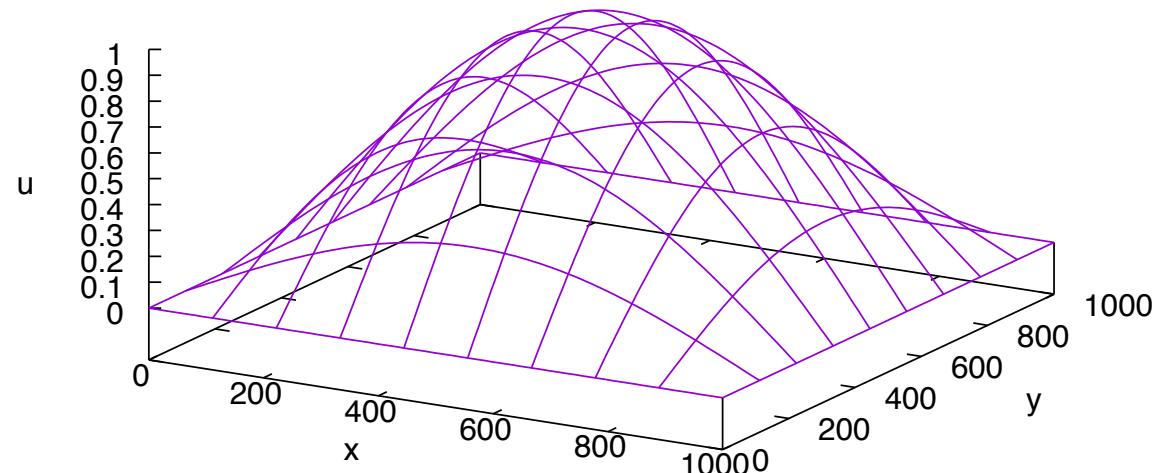
- **Choose** a function for  $u(t, x, y)$ , substitute into the equation and work through the algebra.
- Ideally like the equation to evaluate to zero so don't need to consider a right-hand side to the equation.
- $u(0, x, y)$  gives the initial conditions.
- Can evaluate boundary conditions, e.g. bottom boundary  $u(0,0, y)$
- Because  $u$  is **known** for all timesteps (it was chosen!), the exact solution is **known**.
- Compare the **computed** solution to the known  $u$  to compute an error.
- Any differences come from approximations in the method, or a bug in your code.

# Method of Manufactured Solution

- For the problem of length  $l$ , choose  $u$ :

$$u(t, x, y) = e^{\frac{-2\alpha\pi^2 t}{l^2}} \sin \frac{\pi x}{l} \sin \frac{\pi y}{l}$$

- Boundary conditions:  $u$  is always zero on the boundaries
- Initial value of grid is then  $u(0, x, y) = \sin \frac{\pi x}{l} \sin \frac{\pi y}{l}$



# Heat program ...

```
// Loop over time steps
for (int t = 0; t < nsteps; ++t) {

    // solve over spatial domain for step t
    solve(n, alpha, dx, dt, u, u_tmp);

    // Pointer swap to get ready for next step
    tmp = u;
    u = u_tmp;
    u_tmp = tmp;
}
```

- Takes two optional command line arguments: <ncells> <nsteps>
  - E.g. ./heat 1000 10
  - 1000x1000 cells, 10 timesteps (the default problem size).
- If no command line arguments are provided, it uses a default:
  - These two commands both run the default problem size of 1000x1000 cells, 10 timesteps.
  - ./heat
  - ./heat 1000 10
- A sensible bigger problem is 8000 x 8000 cells and 10 timesteps.

# Heat diffusion problem: 5-point stencil code

```

const double r = alpha * dt / (dx * dx);
const double r2 = 1.0 - 4.0*r;

// malloc and initialize u_tmp and u (code not shown)

for (int t = 0; t < nsteps; ++t) {           Loop over time steps

    for (int i = 0; i < n; ++i) {
        for (int j = 0; j < n; ++j) {           Loop over NxN spatial domain
            u_tmp[i+j*n] = r2 * u[i+j*n] +  

                r * ((i < n-1) ? u[i+1+j*n] : 0.0) +  

                r * ((i > 0)   ? u[i-1+j*n] : 0.0) +  

                r * ((j < n-1) ? u[i+(j+1)*n] : 0.0) +  

                r * ((j > 0)   ? u[i+(j-1)*n] : 0.0);
        }
    }
    // Pointer swap to get ready for next step
    tmp = u;
    u = u_tmp;
    u_tmp = tmp;
}

```

# Heat diffusion problem: 5-point stencil code

```

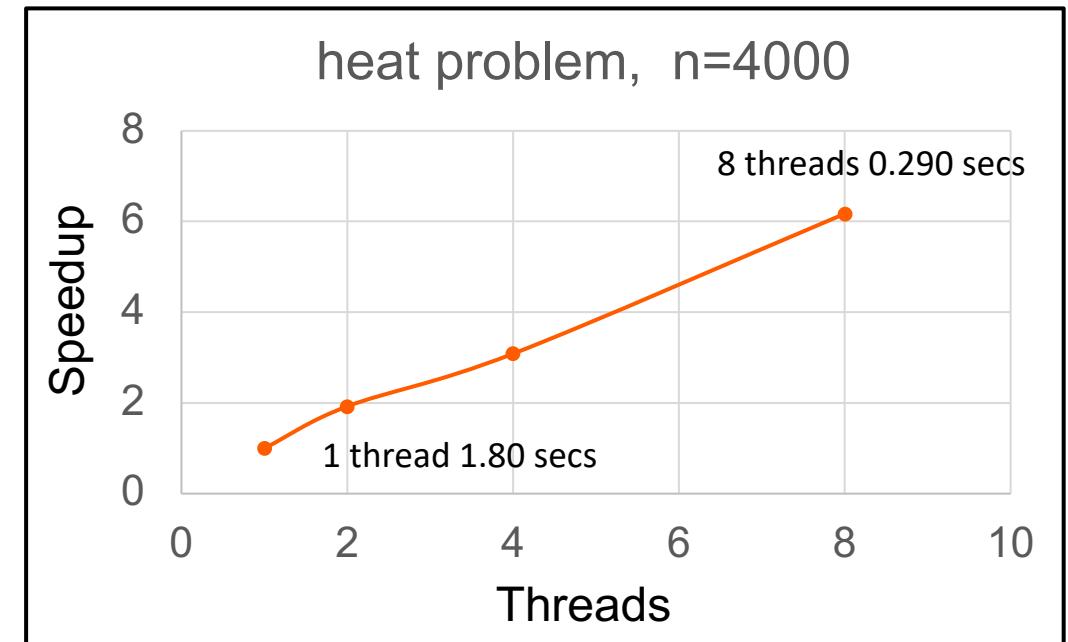
const double r = alpha * dt / (dx * dx);
const double r2 = 1.0 - 4.0*r;

// malloc and initialize u_tmp and u (code not shown)

for (int t = 0; t < nsteps; ++t) {

    #pragma omp parallel for collapse(2)
    for (int i = 0; i < n; ++i) {
        for (int j = 0; j < n; ++j) {
            u_tmp[i+j*n] = r2 * u[i+j*n] + 
                r * ((i < n-1) ? u[i+1+j*n] : 0.0) +
                r * ((i > 0)   ? u[i-1+j*n] : 0.0) +
                r * ((j < n-1) ? u[i+(j+1)*n] : 0.0) +
                r * ((j > 0)   ? u[i+(j-1)*n] : 0.0);
        }
    }
    // Pointer swap to get ready for next step
    tmp = u;
    u = u_tmp;
    u_tmp = tmp;
}

```



Intel® Xeon™ Gold 5218 @ 2.3 Ghz, 8 cores. Nvidia HPC Toolkit compiler  
nvc –fast –fopenmp heat.c

# Exercise: parallel stencil (heat)

- Take the provided heat stencil code (heat.c)
- Add OpenMP directives to parallelize the loops on the GPU
- Most of the runtime occurs in the solve() routine. Focus on that function. The rest of the code is there to just support the work inside solve.

- double omp\_get\_wtime();
- #pragma omp parallel
- #pragma omp for
- #pragma omp parallel for
- #pragma omp task
- #pragma omp taskwait
- #pragma single
- #pragma omp target
- #pragma omp loop
- Plus the clauses
  - private(), firstprivate(), reduction(+:var), collapse(n)
  - map(to:vptr[Lower:Count]) map(from:vptr[Lower:Count]) map(tofrom:vptr[Lower:Count])

Note: Swapping pointer on the host before entering the target region on the next iteration works on a GPU.

When you map pointers between the host and the device, OpenMP remembers the address.

Swapped addresses on the hosts swaps addresses on the device

Default is tofrom: map(vptr[Lower:Count])



# Heat diffusion problem: 5-point stencil code

```

const double r = alpha * dt / (dx * dx);
const double r2 = 1.0 - 4.0*r;

// malloc and initialize u_tmp and u (code not shown)

for (int t = 0; t < nsteps; ++t) {
    #pragma omp target map(tofrom: u[0:n*n], u_tmp[0:n*n])
    #pragma omp loop
    for (int i = 0; i < n; ++i) {
        for (int j = 0; j < n; ++j) {
            u_tmp[i+j*n] = r2 * u[i+j*n] +
                r * ((i < n-1) ? u[i+1+j*n] : 0.0) +
                r * ((i > 0) ? u[i-1+j*n] : 0.0) +
                r * ((j < n-1) ? u[i+(j+1)*n] : 0.0) +
                r * ((j > 0) ? u[i+(j-1)*n] : 0.0);
        }
    }
    // Pointer swap to get ready for next step
    tmp = u;
    u = u_tmp;
    u_tmp = tmp;
}

```

When you map pointers between the host and the device, OpenMP remembers the address.

Swapped addresses on the hosts swaps addresses on the device

GPU Solver time = 1.40 secs

This isn't much better than the runtime for a single CPU (1.8 secs) and worse than 8 cores on a CPU (0.29 secs).

**Why is the performance so bad?**

NVIDIA T4 GPU, 16 Gbyte, Turing Arch.  
 Nvidia HPC Toolkit compiler  
 nvc -fast -mp=gpu -gpu=cc75 heat.c

# Heat diffusion problem: 5-point stencil code

```

const double r = alpha * dt / (dx * dx);
const double r2 = 1.0 - 4.0*r;

// malloc and initialize u_tmp and u (code not shown)

for (int t = 0; t < nsteps; ++t) {
    #pragma omp target map(tofrom: u[0:n*n], u_tmp[0:n*n])
    #pragma omp loop
    for (int i = 0; i < n; ++i) {
        for (int j = 0; j < n; ++j) {
            u_tmp[i+j*n] = r2 * u[i+j*n] +
                r * ((i < n-1) ? u[i+1+j*n] : 0.0) +
                r * ((i > 0) ? u[i-1+j*n] : 0.0) +
                r * ((j < n-1) ? u[i+(j+1)*n] : 0.0) +
                r * ((j > 0) ? u[i+(j-1)*n] : 0.0);
        }
    }
    // Pointer swap to get ready for next step
    tmp = u;
    u = u_tmp;
    u_tmp = tmp;
}

```

At the end of each iteration, **copy**  
 $(2*N^2)*\text{sizeof}(\text{TYPE})$  bytes  
**from** the device

With a runtime of 1.4 secs (worse than the CPU time) we see that Data Movement dominates performance.

At the beginning of each iteration, **copy**  
 $(2*N^2)*\text{sizeof}(\text{TYPE})$  bytes **to** the device

We need to create a **data region** on the GPU that is distinct from the target region.

That way, we can keep the data on the device between target constructs

# Finer control over data movement

- Recall that data is mapped to/from the device at start/end of target region
  - `#pragma omp target map(tofrom: A[0:N]) { ... }`
- Inefficient to move data around all the time
- We need to keep data resident on the device *between* target regions.
- We need to manage the device data environment independently from the target regions that launch kernels.

# Target enter/exit data constructs

- Create a data region on the target device (a device data environment) with two standalone directives:

```
#pragma omp target enter data map(...)  
#pragma omp target exit data map(...)
```

- The **target enter data** maps variables to the device data environment.
- The **target exit data** unmaps variables from the device data environment.
- Once created, subsequent **target** regions inherit the existing data environment.

# Target enter/exit data example

```
void init_array(int *A, int *B, int N) {  
    for (int i = 0; i < N; ++i) { A[i] = i;  B[i]=2*i; }  
  
    #pragma omp target enter data map(to: A[0:N], B[0:N])  
}  
  
int main(void) {  
  
    int N = 1024;  
    int *A = malloc(sizeof(int) * N);  
    int *B = malloc(sizeof(int) * N);  
    init_array(A, B, N);  
  
    #pragma omp target  
    #pragma omp loop  
    for (int i = 0; i < N; ++i)  
        A[i] = A[i] * B[i];  
  
    #pragma omp target exit data map(from: A[0:N])  
}
```

# Exercise

- Modify your parallel heat code from the last exercise.
- Use the ‘target data’ constructs to control the device data environment.
- Think about the map clauses and how to use them to minimize data movement.
- Question ... will the pointer swap on the host still work?
  - `#pragma omp target`
  - `#pragma omp target enter data map(to: list)`
  - `#pragma omp target exit data map(from: list)`
  - `#pragma omp target update`
  - `map(to:list) map(from:list) map(tofrom:list)`
  - `#pragma omp teams distribute parallel for simd`

# Heat diffusion problem: 5-point stencil code

```

const double r = alpha * dt / (dx * dx);
const double r2 = 1.0 - 4.0*r;
// malloc and initialize u_tmp and u (code not shown)
#pragma omp target enter data map(to: u[0:n*n], u_tmp[0:n*n])

for (int t = 0; t < nsteps; ++t) {
    #pragma omp target
    #pragma omp loop
    for (int i = 0; i < n; ++i) {
        for (int j = 0; j < n; ++j) {
            u_tmp[i+j*n] = r2 * u[i+j*n] +
                r * ((i < n-1) ? u[i+1+j*n] : 0.0) +
                r * ((i > 0) ? u[i-1+j*n] : 0.0) +
                r * ((j < n-1) ? u[i+(j+1)*n] : 0.0) +
                r * ((j > 0) ? u[i+(j-1)*n] : 0.0);
        }
    }
    // Pointer swap to get ready for next step
    tmp = u;
    u = u_tmp;
    u_tmp = tmp;
}
#pragma omp target exit data map(from: u[0:n*n])
}

```

Create a data region and map indicated data on entry

GPU Solver time\* = **0.42** secs

This is a general principle ...  
if you want performance, you  
must optimize data  
movement.

\*includes time for target enter/exit data

Exit the data  
region and map  
indicated data

NVIDIA T4 GPU, 16 Gbyte, Turing Arch.  
Nvidia HPC Toolkit compiler  
nvc -fast -mp=gpu -gpu=cc75 heat.c

**Getting the data movement between host memory and device memory is key.**

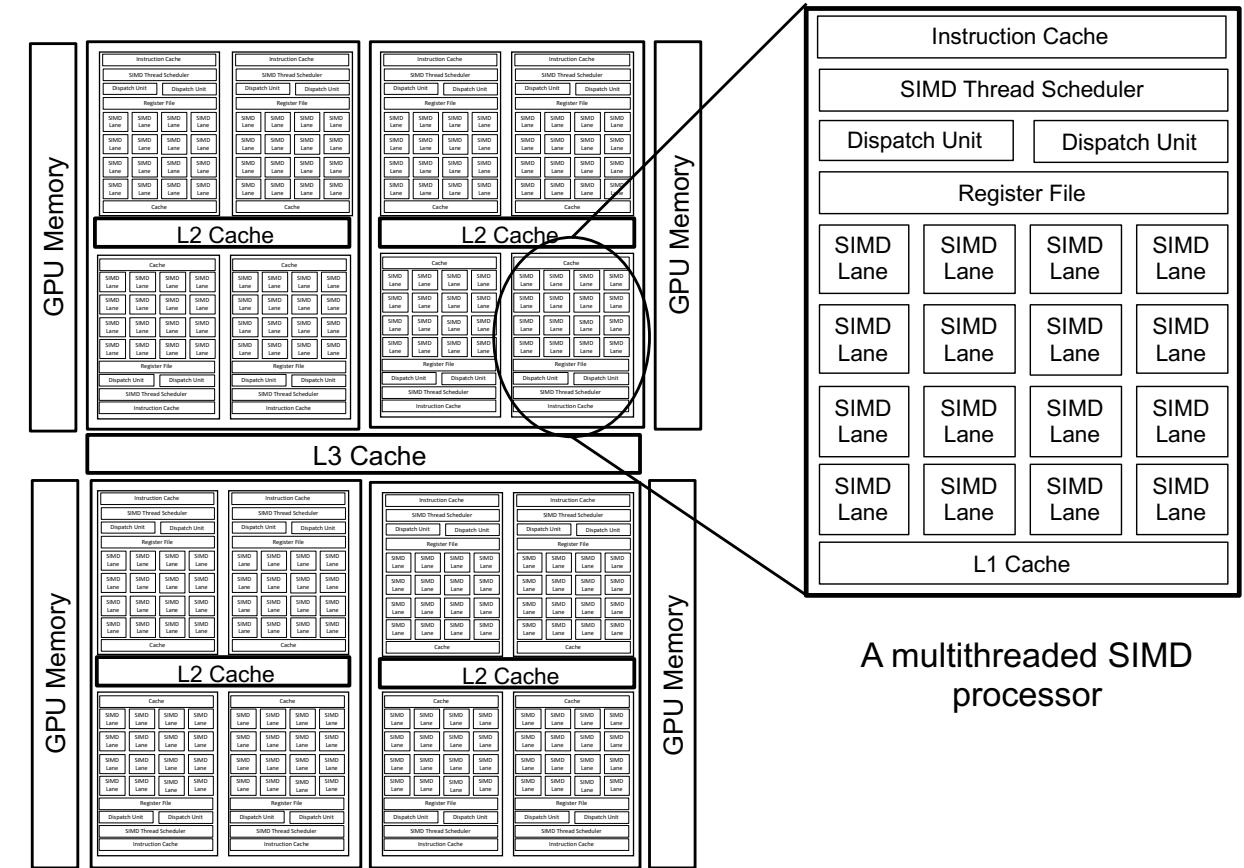
**What are the other major issues to consider when optimizing performance?**

# Occupancy: Keep all the GPU resources busy

- In our “GPU cartoon” we have 16 multithreaded SIMD processors each with 16 SIMD lanes .... For a total of  $16^2=256$  processing elements.
- You want all resources busy at all times. You do that by keeping excess work for the multithreaded SIMD processors ... if they are other busy on some high latency operation, you want a new work-group is ready to be scheduled for execution.
- Occupancy having enough work-groups to keep the GPU busy. To support high occupancy, you need many more work-items than SIMD-lanes.

```
#pragma omp parallel for  
for(int i=0;i<N;i++)  
    for(int j=0;j<N;j++)  
        for(int k=0;k<N;k++)  
            *(C+(i*N+j)) += *(A+(i* N +k)) * *(B+(k* N +j));
```

Parallelize i-loop  
parallelism O(N)

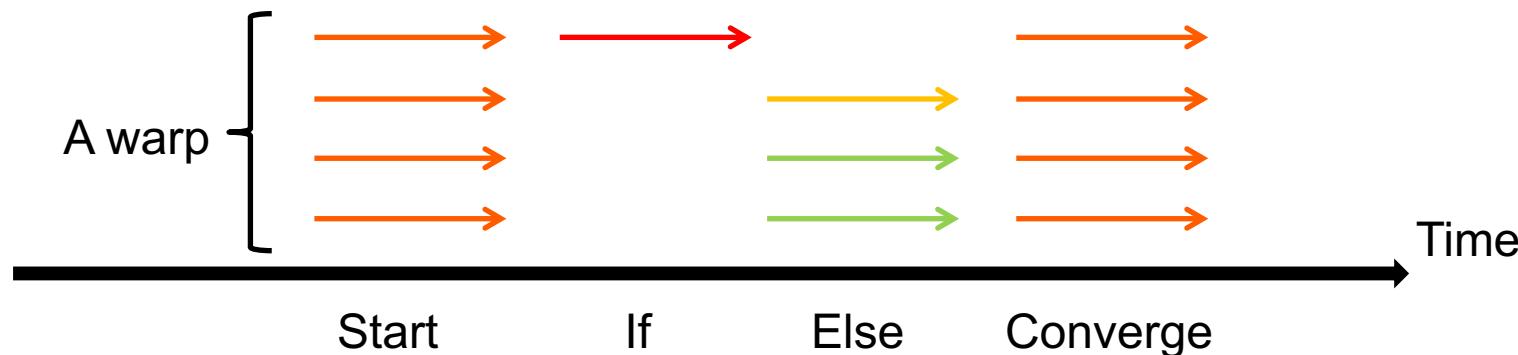


```
#pragma omp parallel for collapse(2)  
for(int i=0;i<N;i++)  
    for(int j=0;j<N;j++)  
        for(int k=0;k<N;k++)  
            *(C+(i*N+j)) += *(A+(i* N +k)) * *(B+(k* N +j));
```

Parallelize combined i/j-loops  
parallelism O(N<sup>2</sup>)

# Converged Execution: Single Instruction Multiple Data

- Individual work-items of a warp start together at the same program address
- Each work-item has its own instruction address counter and register state
  - Each work-item is free to branch and execute independently
  - Supports the Single Program Multiple Data (SPMD) pattern.
- Branch behavior
  - Each branch will be executed serially
  - Work-items not following the current branch will be disabled



# Converged Execution: Branching

- GPUs tend not to support speculative execution, which means that branch instructions have high latency
- This latency can be hidden by switching to alternative work-items/work-groups, but avoiding branches where possible is still a good idea to improve performance
- When different work-items executing within the same SIMD ALU array take different paths through conditional control flow, we have ***divergent branches*** (vs. ***uniform branches***)
- Divergent branches are bad news: some work-items will stall while waiting for the others to complete
- We can use predication, selection and masking to convert conditional control flow into straight line code and significantly improve the performance of code that has lots of conditional branches

# Branching

## Conditional execution

```
// Only evaluate expression  
// if condition is met  
if (a > b)  
{  
    acc += (a - b*c);  
}
```

## Selection and masking

```
// Always evaluate expression  
// and mask result  
temp = (a - b*c);  
mask = (a > b ? 1.f : 0.f);  
acc += (mask * temp);
```

# Coalesced memory accesses

- **Coalesced memory accesses** are key for high performance code, especially on GPUs
- In principle, it's very simple, but frequently requires transposing or transforming data on the host before sending it to the GPU
- Sometimes this is an issue of Array of Structures vs. Structure of Arrays (AoS vs. SoA)

# Memory layout is critical to performance

- Structure of Arrays vs. Array of Structures

- Array of Structures (AoS) more natural to code:

```
struct Point{ float x, y, z, a; };
```

```
Point *Points;
```



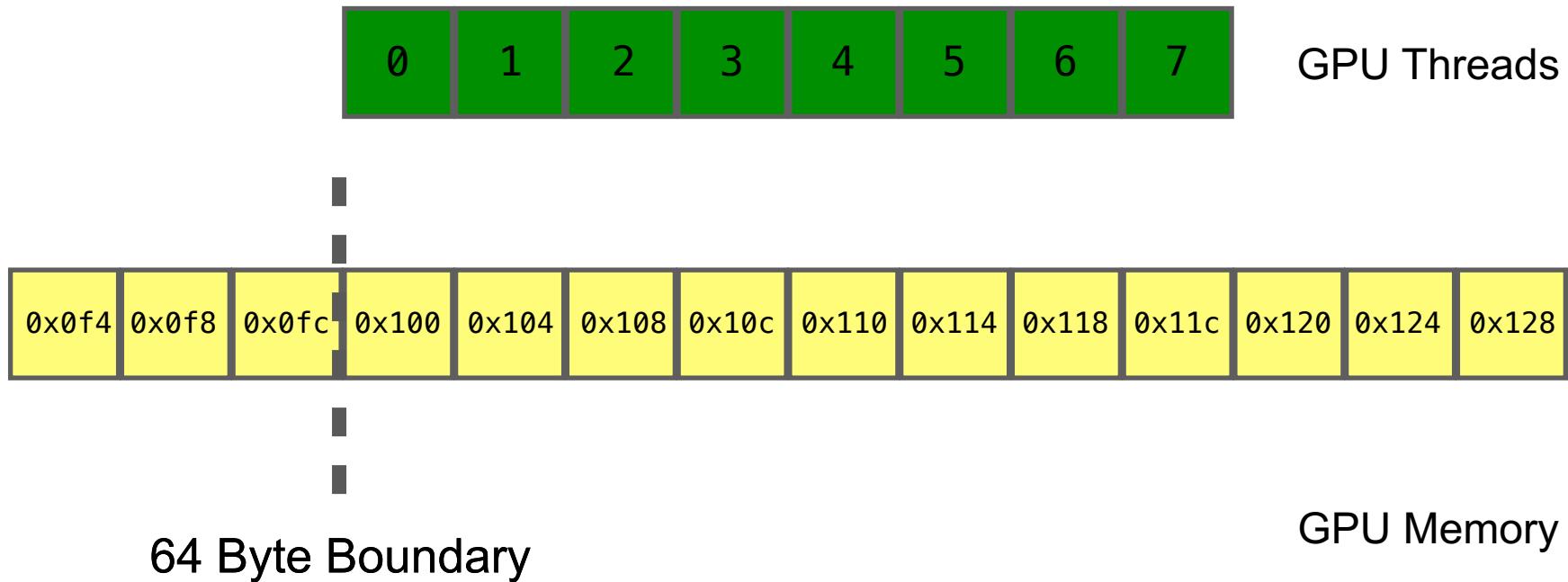
- Structure of Arrays (SoA) suits memory coalescence in vector units

```
struct { float *x, *y, *z, *a; } Points;
```



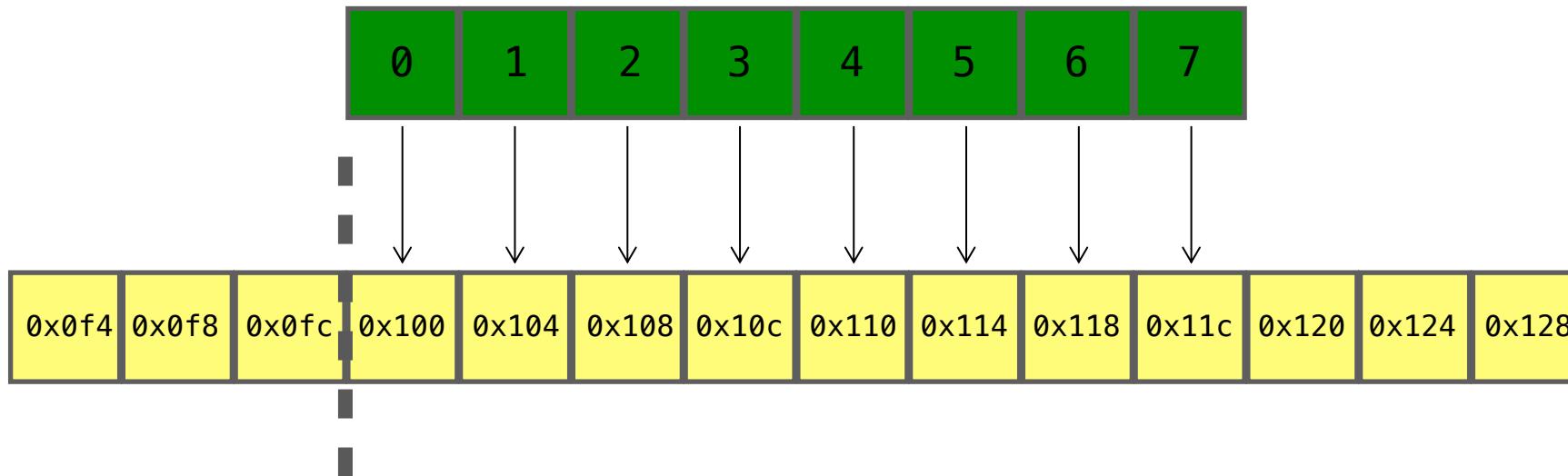
Adjacent work-items/vector-lanes like to access adjacent memory locations

# Memory access patterns



# Memory access patterns

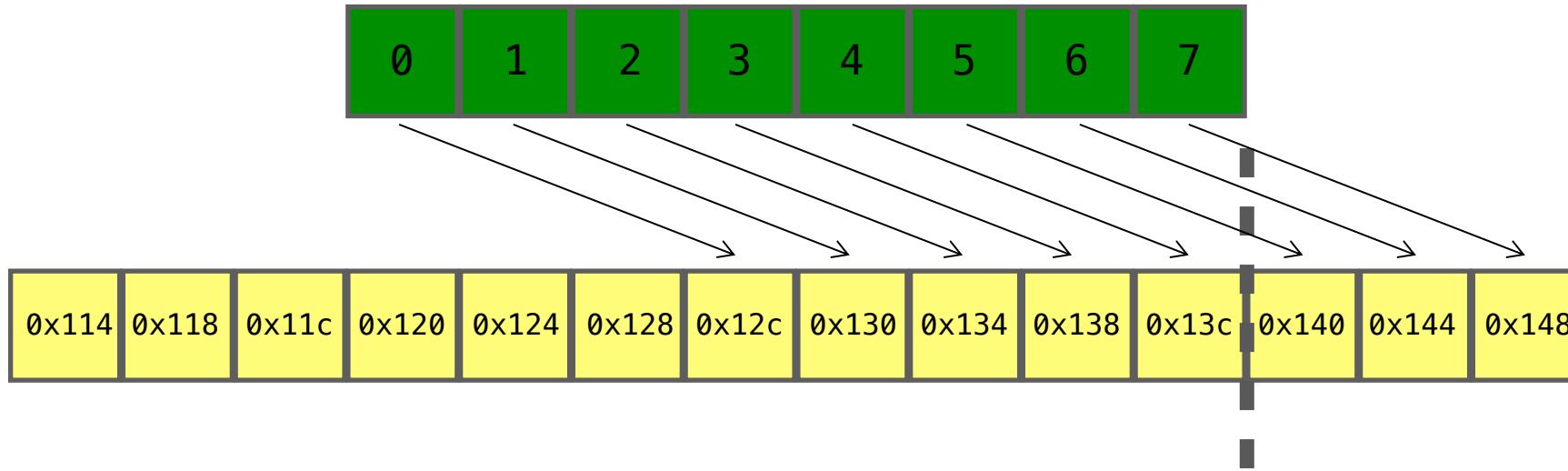
```
float val1 = memA[id];
```



64 Byte Boundary

# Memory access patterns

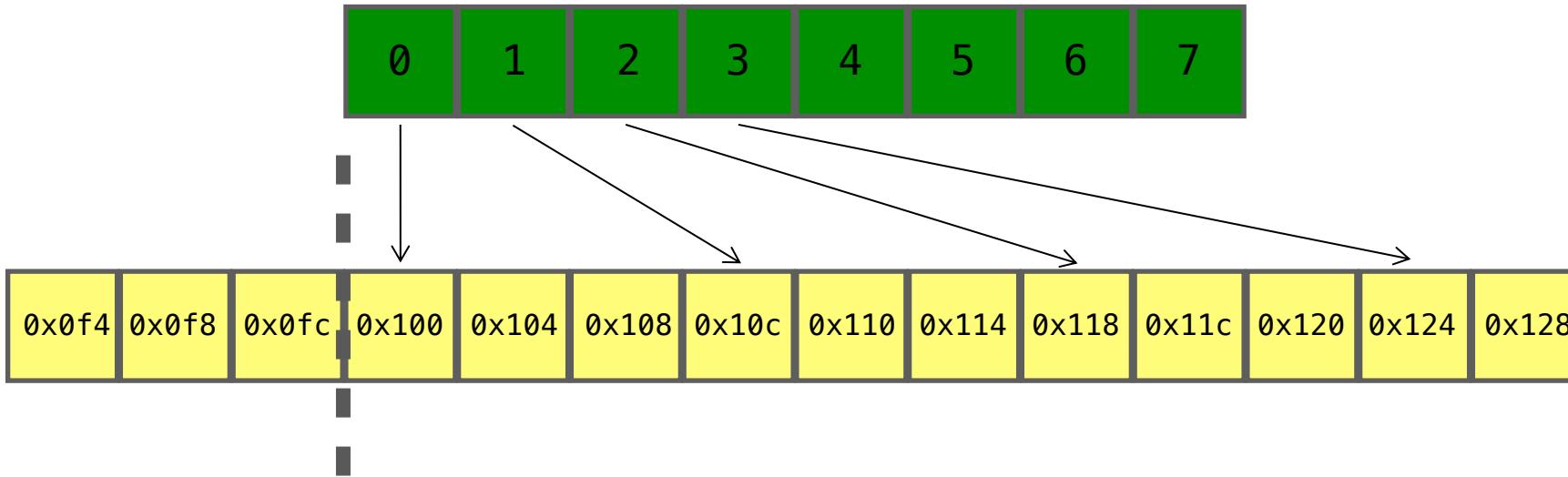
```
const int c = 3;  
float val2 = memA[id + c];
```



64 Byte Boundary

# Memory access patterns

```
float val3 = memA[3*id];
```



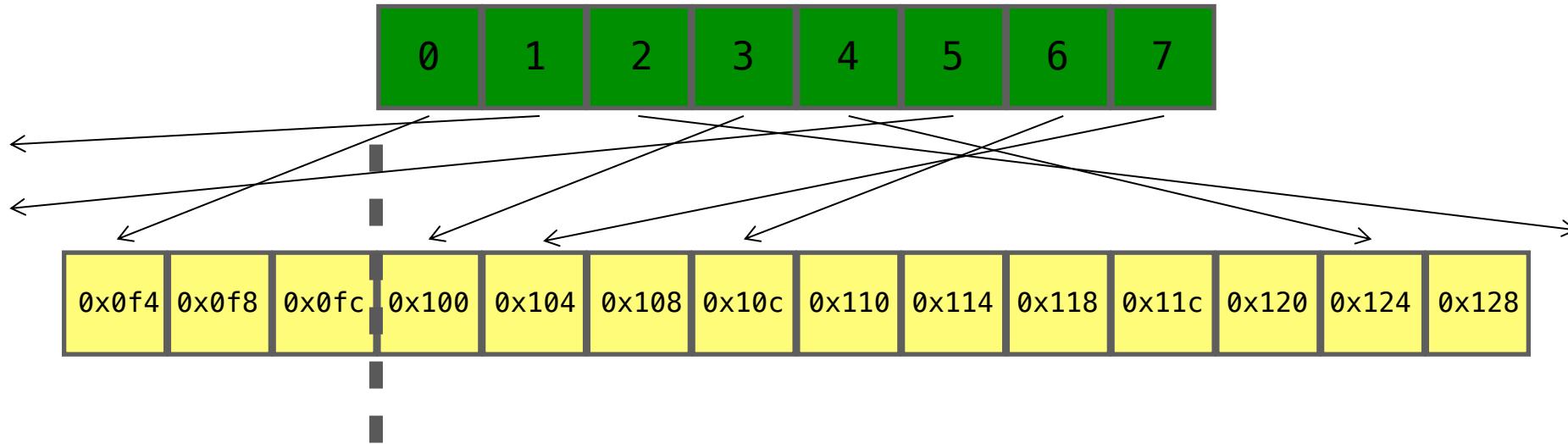
64 Byte Boundary

Strided access results in multiple  
memory transactions (and  
kills throughput)

# Memory access patterns

```
const int loc =  
    some_strange_func(id);
```

```
float val4 = memA[loc];
```



64 Byte Boundary

# Coalescence

- **Coalesce** - to combine into one
- Coalesced memory accesses are key for high bandwidth
- Simply, it means, if thread  $i$  accesses memory location  $n$  then thread  $i+1$  accesses memory location  $n+1$
- In practice, it's not quite as strict...

```
for (int id = 0; id < size; id++)
{
    // ideal
    float val1 = memA[id];

    // still pretty good
    const int c = 3;
    float val2 = memA[id + c];

    // stride size is not so good
    float val3 = memA[c*id];

    // terrible
    const int loc =
        some_strange_func(id);

    float val4 = memA[loc];
}
```

# Exercise

- Optimize the stencil ‘solve’ kernel.
- Start with your code with optimized memory movement from the last exercise.
- Experiment with the optimizations we’ve discussed.
- Focus on the memory access pattern.
- Try different input sizes to see the effect of the optimizations.
- Keep an eye on the solve time as reported by the application.

# Heat diffusion problem: 5-point stencil code

Parallel CPU and GPU results, n=4000

```
const double r = alpha * dt / (dx * dx);  
const double r2 = 1.0 - 4.0*r;  
// malloc and initialize u_tmp and u (code not shown)
```

```
for (int t = 0; t < nsteps; ++t) {
```

```
#pragma omp parallel for
```

```
for (int i = 0; i < n; ++i) {
```

```
    for (int j = 0; j < n; ++j) {
```

```
        u_tmp[i+j*n] = r2 * u[i+j*n] +  
            r * ((i < n-1) ? u[i+1+j*n] : 0.0) +  
            r * ((i > 0) ? u[i-1+j*n] : 0.0) +  
            r * ((j < n-1) ? u[i+(j+1)*n] : 0.0) +  
            r * ((j > 0) ? u[i+(j-1)*n] : 0.0);
```

```
}
```

```
}
```

```
// Pointer swap for next step
```

```
tmp = u;
```

```
u = u_tmp;
```

```
u_tmp = tmp;
```

```
}
```

This is the ij loop order.  
Swap these loops to get  
the ji order.

GPU

ij without timing  
enter and exit data

0.056830

Num threads	ij loop order	ji loop order
1	1.512849	0.262260
2	0.776229	0.132453
4	0.400822	0.064220
8	0.227317	0.046586

Intel® Xeon™ Gold 5218 @ 2.3 Ghz, 8 cores.  
Nvidia HPC Toolkit compiler nvc-fast-fopenmp heat.c

All times in seconds

This GPU code used the  
target enter data and  
target exit data

NVIDIA T4 GPU, 16 Gbyte, Turing Arch.  
Nvidia HPC Toolkit compiler. nvc -fast -mp=gpu heat.c

# Solution: collapse + swap loop order

```
// Compute the next timestep, given the current timestep
void solve(const int n, const double alpha, const double dx, const double dt, const double * restrict u,
double * restrict u_tmp) {
    // Finite difference constant multiplier
    const double r = alpha * dt / (dx * dx);
    const double r2 = 1.0 - 4.0*r;

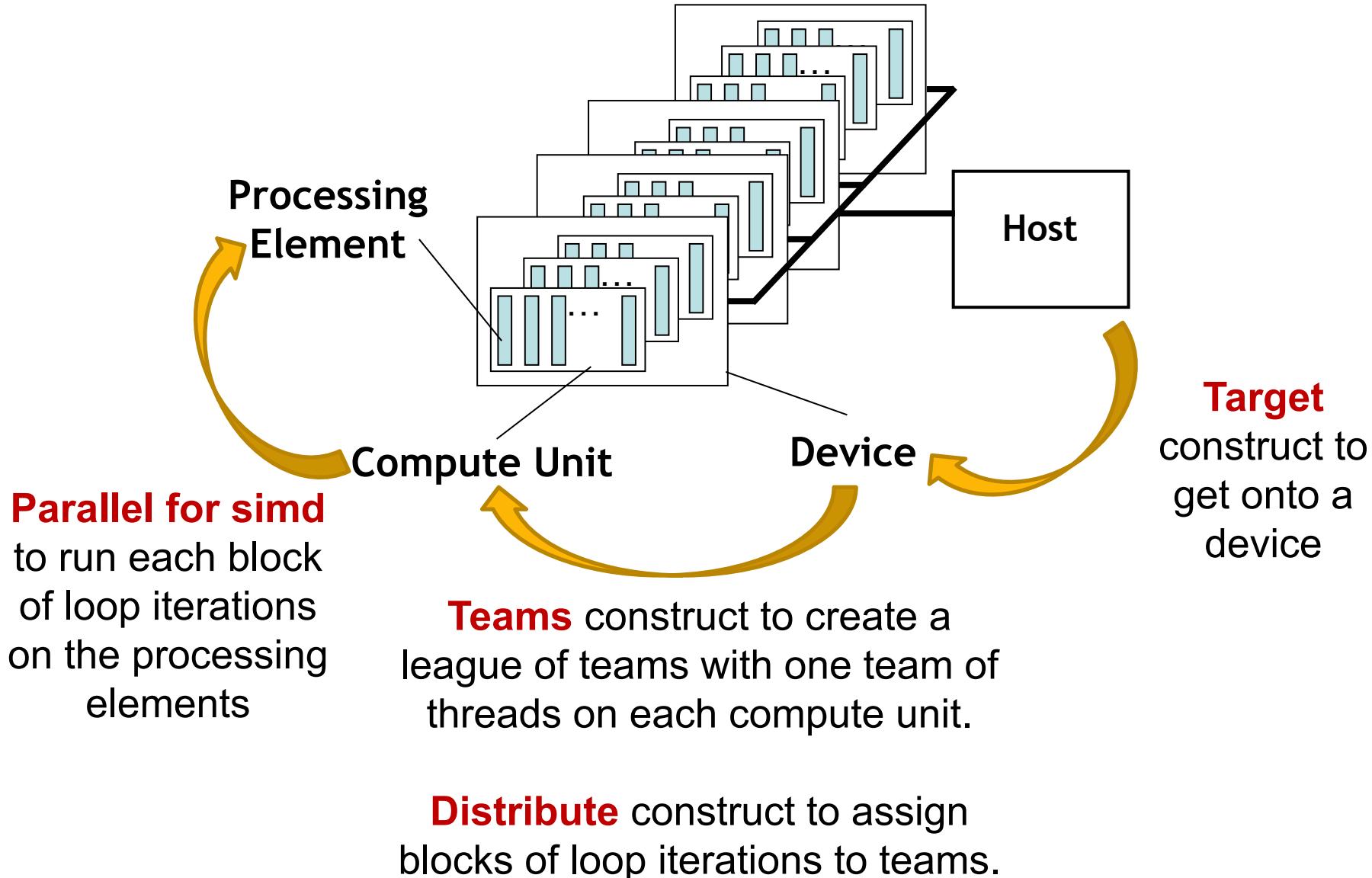
    // Loop over the nxn grid
    #pragma omp target
    #pragma omp loop collapse(2)
    for (int j = 0; j < n; ++j) {
        for (int i = 0; i < n; ++i) {
            // Update the 5-point stencil, using boundary conditions on the edges of the domain.
            // Boundaries are zero because the MMS solution is zero there.
            u_tmp[i+j*n] = r2 * u[i+j*n] +
                           r * ((i < n-1) ? u[i+1+j*n] : 0.0) +
                           r * ((i > 0) ? u[i-1+j*n] : 0.0) +
                           r * ((j < n-1) ? u[i+(j+1)*n] : 0.0) +
                           r * ((j > 0) ? u[i+(j-1)*n] : 0.0);
        }
    }
}
```

Create more work ... to better fill the processing elements of the GPU

Swap the i and j loops so that the  $i+j \cdot n$  memory accesses are contiguous

**The loop construct is great, but sometimes you want more control.**

# Our host/device Platform Model and OpenMP



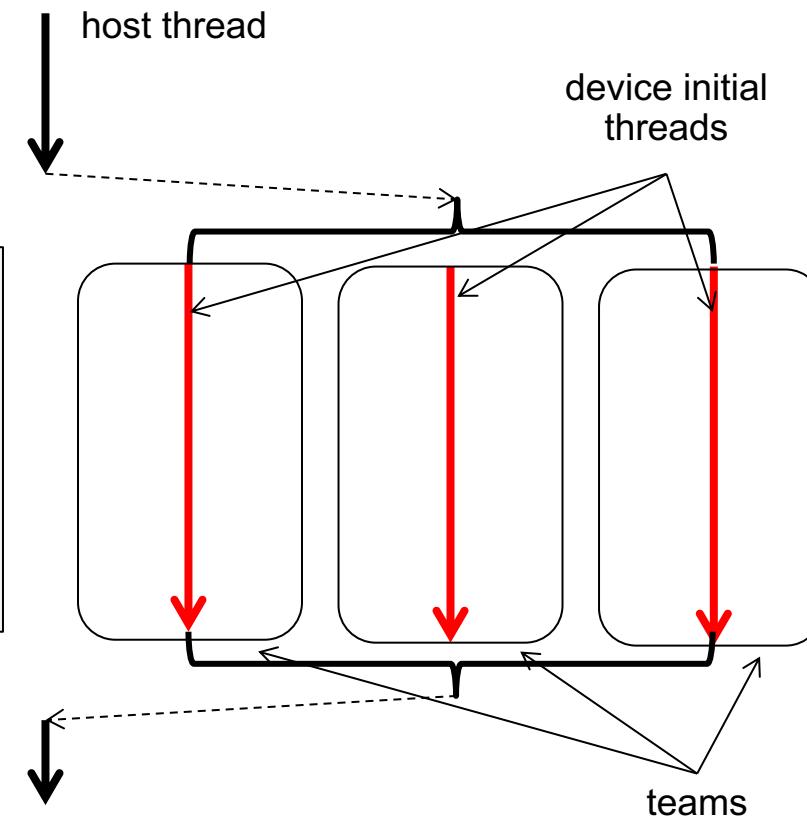
# teams and distribute constructs

- The **teams** construct
  - Similar to the **parallel** construct
  - It starts a league of thread teams
  - Each team in the league starts as one initial thread – a team of one
  - Threads in different teams cannot synchronize with each other
  - The construct must be “perfectly” nested in a **target** construct
- The **distribute** construct
  - Similar to the **for** construct
  - Loop iterations are workshared across the initial threads in a league
  - No implicit barrier at the end of the construct
  - **dist\_schedule(*kind*[, *chunk\_size*])**
    - If specified, scheduling kind must be static
    - Chunks are distributed in round-robin fashion in chunks of size ***chunk\_size***
    - If no chunk size specified, chunks are of (almost) equal size; each team receives at least one chunk

# Create a league of teams and distribute a loop among them

- teams construct
- distribute construct

```
#pragma omp target
#pragma omp teams
#pragma omp distribute
for (i=0;i<N;i++)
...
...
```

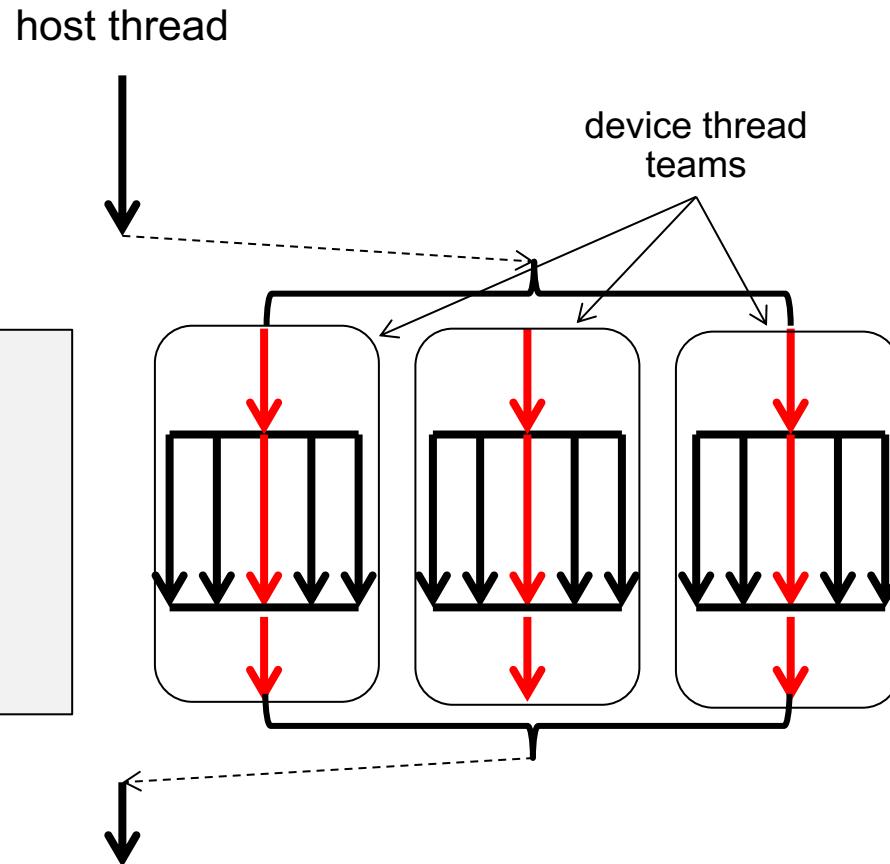


- Transfer execution control to **MULTIPLE** device initial threads
- Workshare loop iterations across the initial threads.

# Create a league of teams and distribute a loop among them and run each team in parallel with its partition of the loop

- teams distribute
- parallel for

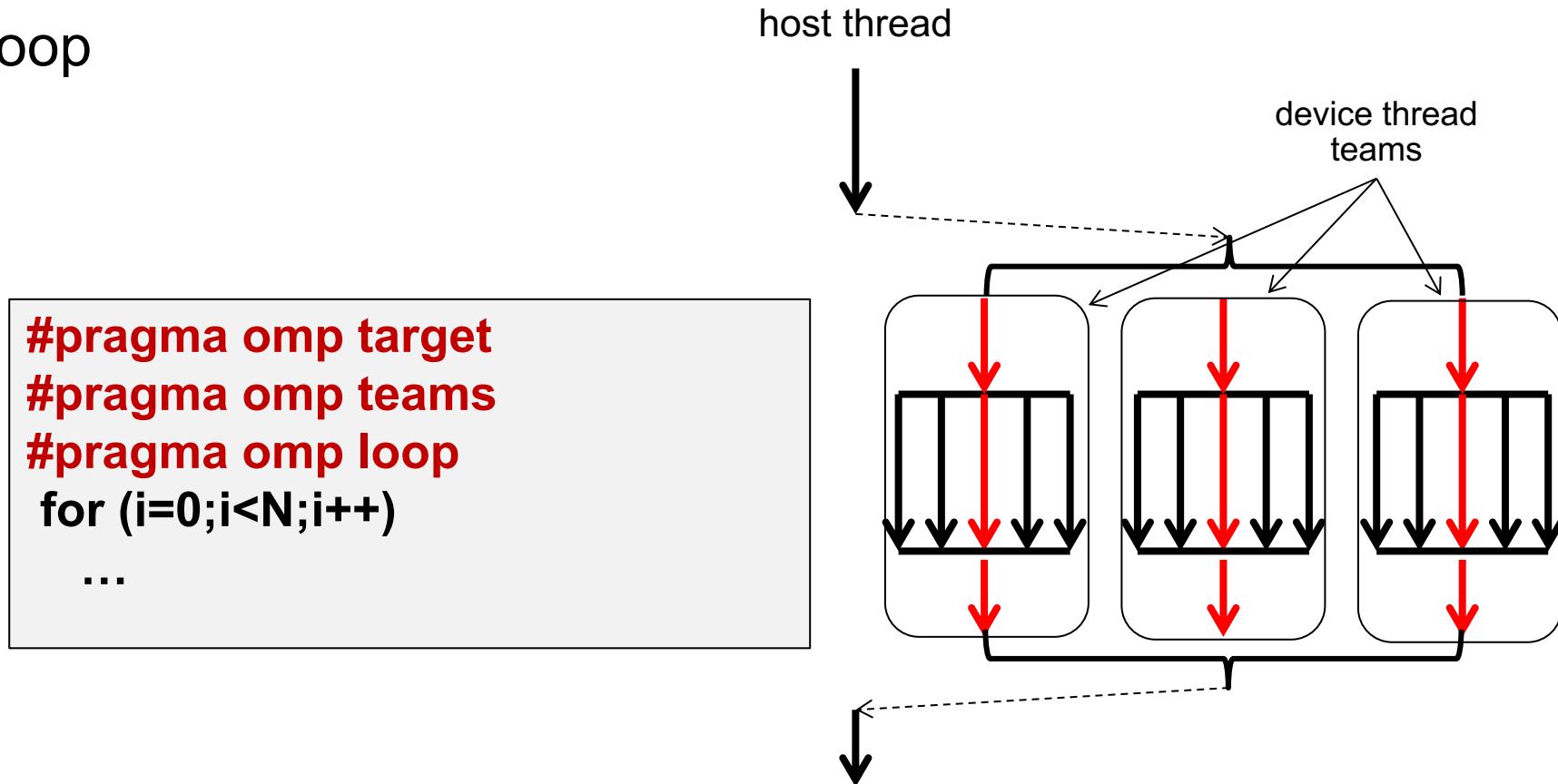
```
#pragma omp target
#pragma omp teams
#pragma omp distribute
#pragma omp parallel for
for (i=0;i<N;i++)
...
```



- Transfer execution control to **MULTIPLE** device initial threads
  - Workshare loop iterations across the initial threads (teams distribute)
- Each initial thread becomes the primary\* thread in a thread team
  - Workshare loop iterations across the threads in a team (parallel for)

# Create a league of teams and distribute a loop among them and run each team in parallel with its partition of the loop

- loop



- Transfer execution control to **MULTIPLE** device initial threads
  - Workshare loop iterations across the initial threads (teams distribute)
- Each initial thread becomes the primary\* thread in a thread team
  - Workshare loop iterations across the threads in a team (parallel for)

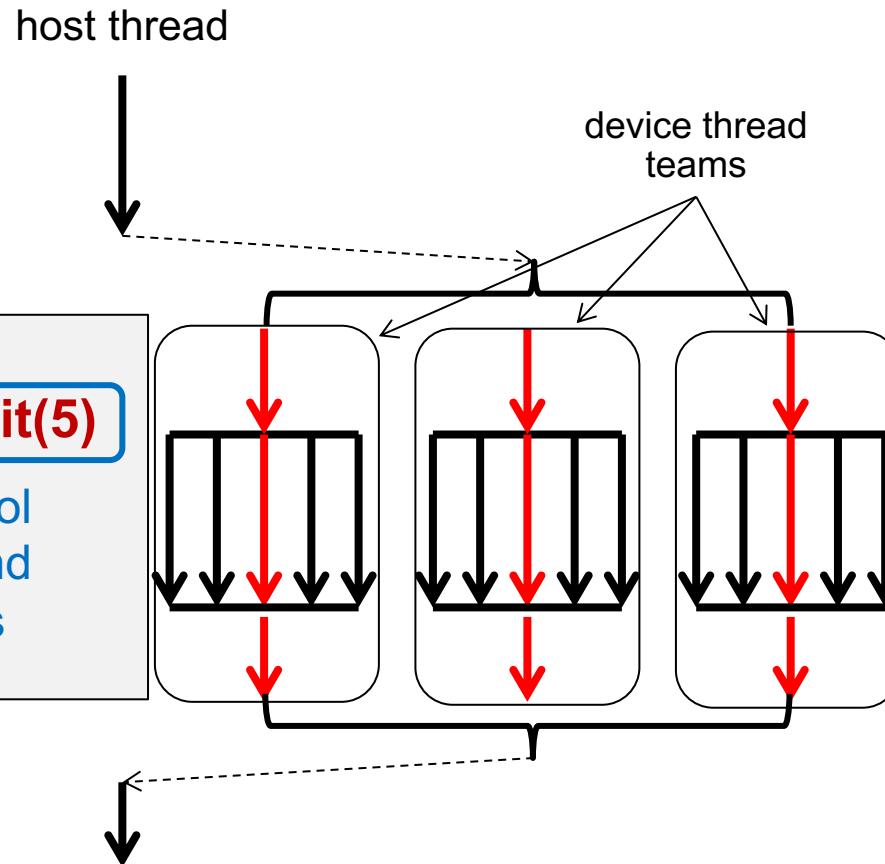
# Create a league of teams and distribute a loop among them and run each team in parallel with its partition of the loop

- teams distribute
- parallel for simd

```
#pragma omp target
#pragma omp teams num_teams(3) thread_limit(5)
#pragma omp distribute
#pragma omp parallel for
for (i=0;i<N;i++)
...

```

Explicit control  
of number and  
size of teams

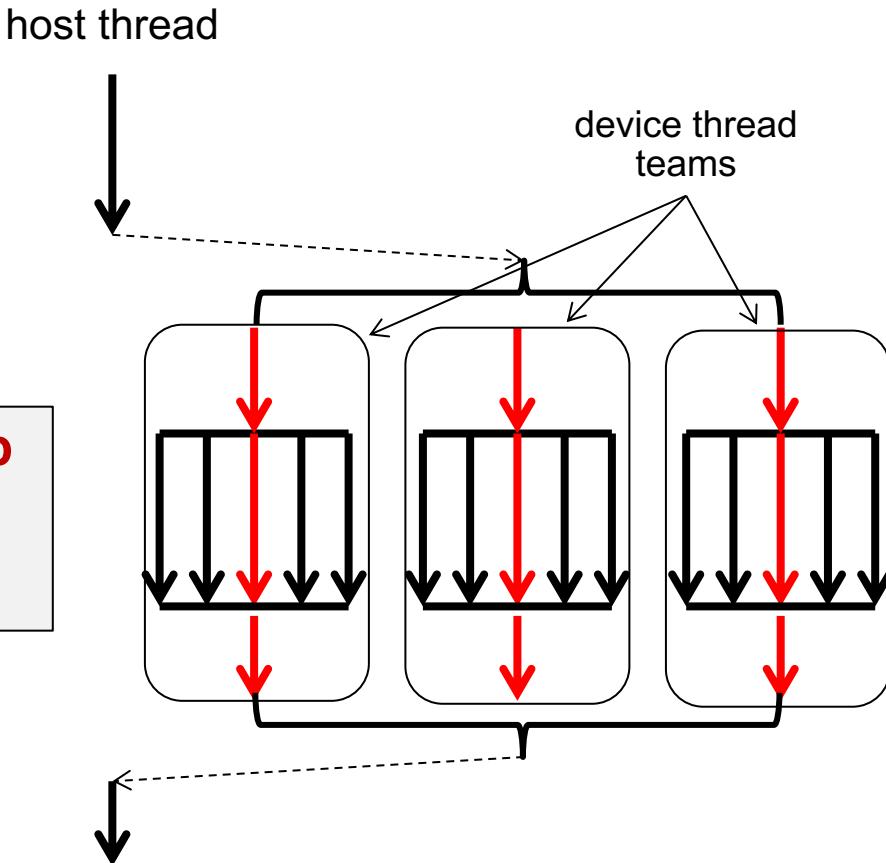


- Transfer execution control to **MULTIPLE** device initial threads
  - Workshare loop iterations across the initial threads (teams distribute)
- Each initial thread becomes the primary\* thread in a thread team
  - Workshare loop iterations across the threads in a team (parallel for)

# Create a league of teams and distribute a loop among them and run each team in parallel with its partition of the loop

- Combined construct

```
#pragma omp target teams loop  
for (i=0;i<N;i++)  
...
```



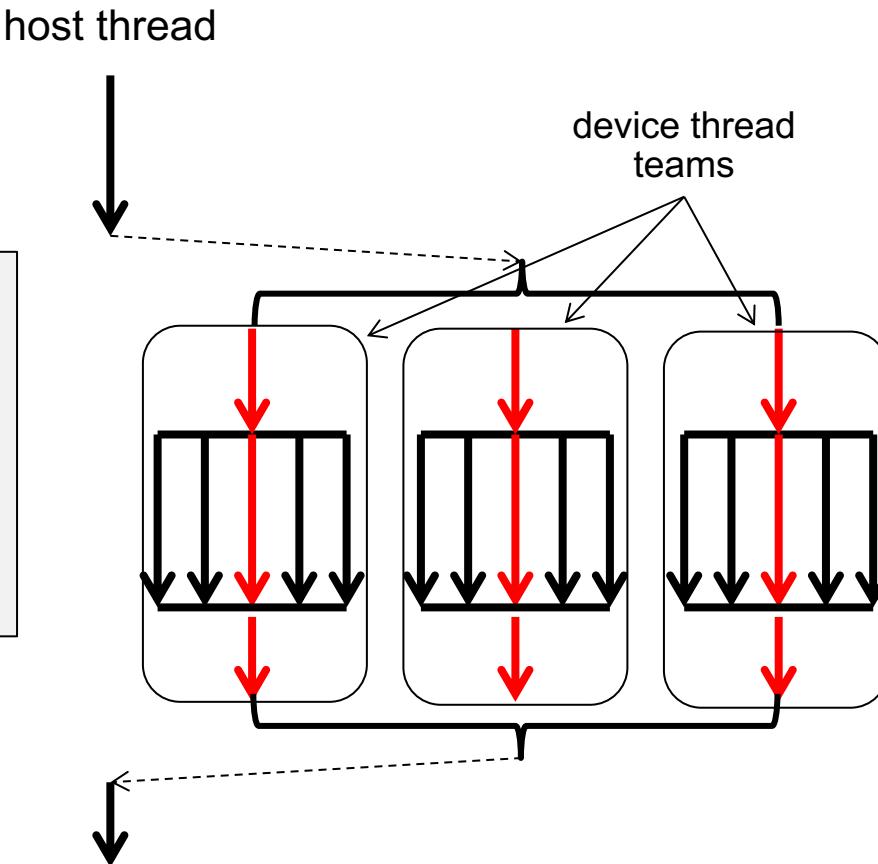
- Transfer execution control to **MULTIPLE** device initial threads
  - Workshare loop iterations across the initial threads (teams distribute)
- Each initial thread becomes the primary\* thread in a thread team
  - Workshare loop iterations across the threads in a team (parallel for)

# Create a league of teams and distribute a loop among them and run each team in parallel with its partition of the loop

- teams distribute
- parallel for simd

Works with nested loops as well

```
#pragma omp target
#pragma omp teams distribute
for (i=0;i<N;i++)
#pragma omp parallel for
for (j=0;j<M;j++)
...
...
```



- Transfer execution control to **MULTIPLE** device initial threads
  - Workshare loop iterations across the initial threads (teams distribute)
- Each initial thread becomes the primary\* thread in a thread team
  - Workshare loop iterations across the threads in a team (parallel for)

# There is MUCH more ... beyond what have time to cover

- Do as much as you can with a simple loop construct. It's portable and as compilers improve over time, it will keep up with compiler driven performance improvements.
- But sometimes you need more:
  - Control over number of teams in a league and the size of the teams
  - Explicit scheduling of loop iterations onto the teams
  - Management of data movement across the memory hierarchy: global vs. shared vs. private ...
  - Calling optimized math libraries (such as cuBLAS)
  - Multi-device programming
  - Asynchrony
- Ultimately, you may need to master all those advanced features of GPU programming. But start with loop. Start with how data on the host maps onto the device (i.e. the GPU). Master that level of GPU programming before worrying about the complex stuff.

**This is the end ... well almost the end.**

**Let's wrap up with a few high-level comments  
about the state of GPU programming more  
generally**

# SIMT Programming models: it's more than just OpenMP

- CUDA:
  - Released ~2006. Made GPGPU programming “mainstream” and continues to drive innovation in SIMT programming.
    - Downside: proprietary to NVIDIA
- OpenCL:
  - Open Standard for SIMT programming created by Apple, Intel, NVIDIA, AMD, and others. 1<sup>st</sup> release in 2009.
  - Supports CPUs, GPUs, FPGAs, and DSP chips. The leading cross platform SIMT model.
    - Downside: extreme portability means verbose API. Painfully low level especially for the host-program.
- Sycl:
  - C++ abstraction layer implements SIMT model with kernels as lambdas. Closely aligned with OpenCL. 1<sup>st</sup> release 2014
    - Downside: Cross platform implementations only emerging recently.
- Directive driven programming models:
  - **OpenACC**: they split from an OpenMP working group to create a competing directive driven API emphasizing descriptive (rather than prescriptive) semantics.
    - Downside: NOT an Open Standard. Controlled by NVIDIA.
  - **OpenMP**: Mixes multithreading and SIMT. Semantics are prescriptive which makes it more verbose. A truly Open standard supported by all the key GPU players.
    - Downside: Poor compiler support so far ... but that will change over the next couple years.

# Vector addition with CUDA

```
// Compute sum of length-N vectors: C = A + B
void __global__
vecAdd (float* a, float* b, float* c, int N) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    if (i < N) c[i] = a[i] + b[i];
}

int main () {
    int N = ... ;
    float *a, *b, *c;
    cudaMalloc (&a, sizeof(float) * N);
    // ... allocate other arrays (b and c), fill with data
    // Use thread blocks with 256 threads each
    vecAdd <<< (N+255)/256, 256 >>> (a, b, c, N);
}
```

Enqueue the kernel  
to execute on the  
Grid

CUDA kernel as  
function

Unified shared  
memory ... allocate  
on host, visible on  
device too

Note: Not all GPUs  
support a unified shared  
memory. We assume its  
available here to we don't  
need to clutter the code  
with explicit data  
movement.

# Vector addition with SYCL

```
// Compute sum of length-N vectors: C = A + B
#include <CL/sycl.hpp>

int main () {
    int N = ... ;
    float *a, *b, *c;
    sycl::queue q;
    *a = (float *)sycl::malloc_shared(N * sizeof(float), q);
    // ... allocate other arrays (b and c), fill with data

    q.parallel_for(sycl::range<1>{N},
                   [=](sycl::id<1> i) {
                       c[i] = a[i] + b[i];
                   });
    q.wait();
}
```

Create a queue for SYCL commands

Note: Not all GPUs support a unified shared memory. We assume its available here so we don't need to clutter the code with explicit data movement.

Unified shared memory ... allocate on host, visible on device too

Kernel as a C++ Lambda function  
[=] means capture external variables by value.

# Vector addition with OpenACC

- Let's add two vectors together ....  $C = A + B$

Host waits here until the kernel is done. Then the output array c is copied back to the host.

```
void vadd(int n,
          const float *a,
          const float *b,
          float *restrict c)
{
    int i;
#pragma acc parallel loop
    for (i=0; i<n; i++)
        c[i] = a[i] + b[i];
}
int main(){
float *a, *b, *c;  int n = 10000;
// allocate and fill a and b

vadd(n, a, b, c);

}
```

Assure the compiler that c is not aliased with other pointers

Turn the loop into a kernel, move data to a device, and launch the kernel.

# Why so many ways to do the same thing?

- The parallel programming model people have failed you ...
  - It's more fun to create something new in your own closed-community that work across vendors to create a portable API
- The hardware vendors have failed you ...
  - Don't you love my “walled garden”? It's so nice here, programmers, just don't even think of going to some other platform since your code is not portable.
- The standards community has failed you ...
  - Standards are great, but they move too slow. OpenACC stabbed OpenMP in the back and I'm pissed, but their comments at the time were spot-on (OpenMP was moving so slow ... they just couldn't wait).
- The applications community failed themselves ...
  - If you don't commit to a standard and use “the next cool thing” you end up with the diversity of overlapping options we have today. Think about what happened with OpenMP and MPI.

# Summary

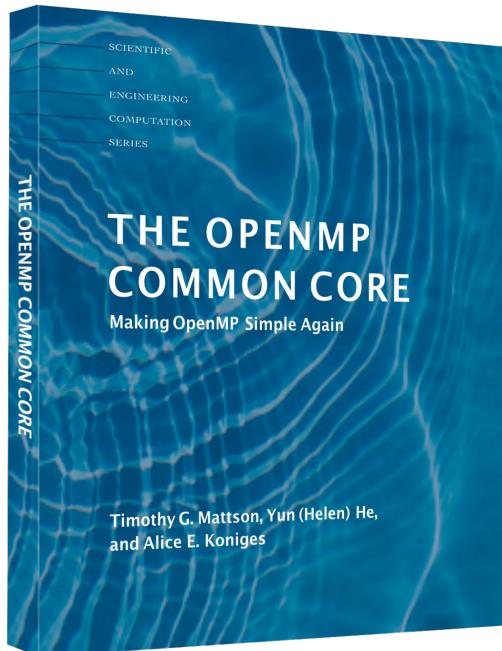
- Parallel computing is fun ... but it can be hard.
- Fortunately, if you stick to the Big-3 and the core patterns of parallel computing for HPC, it's not too overwhelming
  - The big 3: MPI, OpenMP, and “a GPU programming model”
  - Key Patterns: SPMD, loop level parallelism, geometric decomposition, divide and conquer, and SIMD
- Some day we'll automate the hard-parts with Machine Programming, but that may be 10 or 20 years!!!!

# To learn more about OpenMP

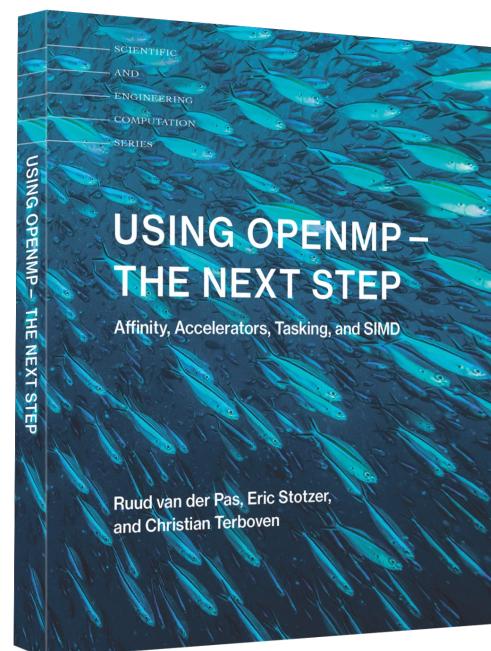
The OpenMP web site has a great deal of material to help you with OpenMP [www.openmp.org](http://www.openmp.org)

Reading the spec is painful ... but each spec has a collection of examples. Study the examples, don't try to read the specs

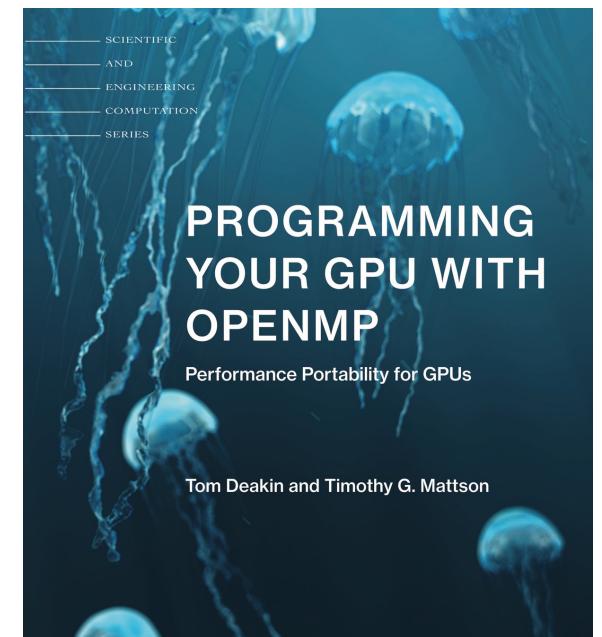
Since the specs are written ONLY for implementors ... programmers need the OpenMP Books to master OpenMP.



Start here ... learn the basics and build a foundation for the future



Learn advanced features in OpenMP including tasking and GPU programming (up to version 4.5)



Learn all the details of GPU programming with OpenMP (up to version 5.2)  
Released in November 2023

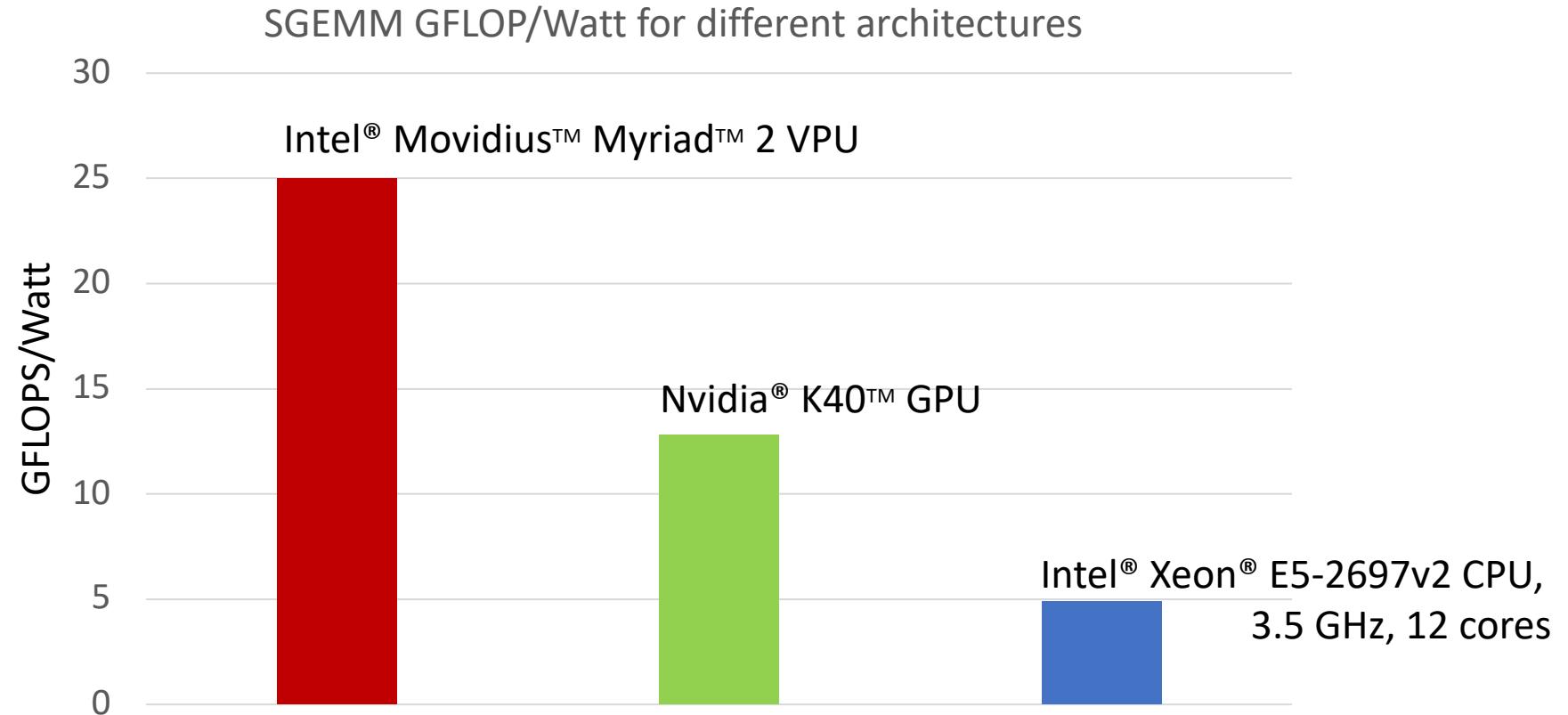
# Backup ... and a bit of extra content



- The future of parallel programming
- The Jacobi solver case study
- Writing functions to call from inside a kernel

# If you care about power, the world is heterogeneous?

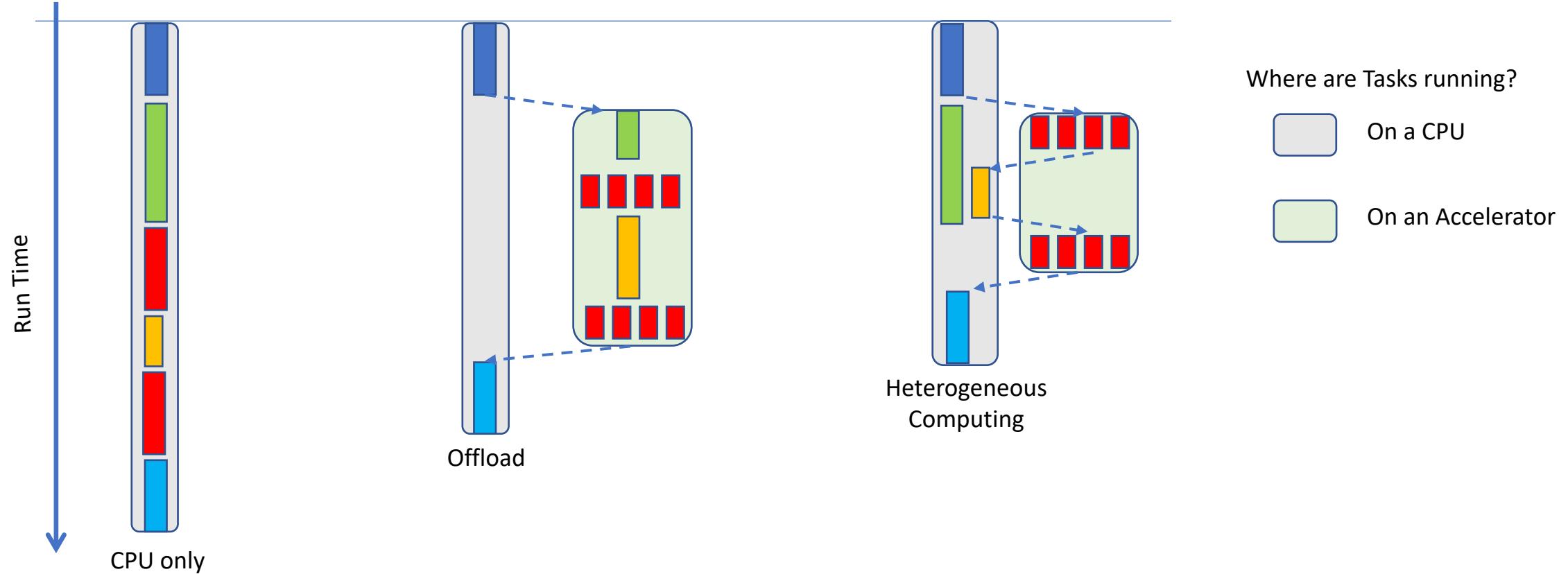
Specialized  
processors doing  
operations suited to  
their architecture  
are more efficient  
than general  
purpose processors.



Hence, future systems will be increasingly heterogeneous ... GPUs, CPUs, FPGAs, and a wide range of accelerators

# Offload vs. Heterogeneous computing

- **Offload:** The CPU moves work to an accelerator and waits for the answer.
- **Heterogeneous Computing:** Run sub-problems in parallel on the hardware best suited to them.



# Example: Single-cell RNA-Seq benchmark (SCANPY)

- SCANPY ... a widely used tool for studying gene expression. All data are elapsed time in seconds
- We started with results from an Nvidia blog (Example 2 from [link](#)), optimized code for one socket of Intel® Xeon® 8380 CPU and then “simulated” heterogeneous computing result by taking the faster of CPU and GPU execution times.

Pipeline stages	64 vCPUs n1-highmem-64 (off-the-shelf Python)	A100 40Gb (Clara Parabricks)	ICX-1s, 40 cores <b>(optimized by Intel)</b>	"Simulated" heterogeneous A100 & ICS-1s 40 cores
Data Loading & Preprocessing	1120	475	15.7	15.7
PCA	44	17.8	5.0	5.0
T-SNE	6509	37	205.6	37
K-means (single iteration)	148	2	7.1	7.1
KNN	154	62	59.8	59.8
UMAP	2571	21	84.5	84.5
Louvain clustering	1153	2.4	6.0	6.0
Leiden clustering	6345	1.7	28.4	28.4
Reanalysis of subgroup	255	17.9	22.5	22.5
Rest	39	49.2	49.0	49.0
<b>End-to-End runtime</b>	<b>18338</b>	<b>686</b>	<b>483.6</b>	<b>211.5</b>

Clara Parabricks: Nvidia solution stack built on RAPIDS for healthcare applications

<https://github.com/clara-parabricks/rapids-single-cell-examples>

github repository as of Dec 16, 2020

## Lessons learned:

- Be careful comparing unoptimized python to hand-tuned CUDA code
- GPUs are great. So are CPUs if you fully utilize all the cores and vector units.
- What you really want is the best of both worlds. **You want heterogeneous computing!**

See Backup for workloads and configurations. Results may vary.

This column shows the potential of heterogeneous computing. We ignored extra communication and synchronization overhead, so actual runtimes would be slightly greater.

Source: Github repository as of Dec 16, 2020 - Example 2: Single-cell RNA-seq of 1.3 Million Mouse Brain Cells comparing CPU (n1-highmem-64 64 vCPUs) vs GPU (n1-highmem-16). <https://github.com/clara-parabricks/rapids-single-cell-examples>. Intel does not control or audit third-party data. You should consult other sources to evaluate accuracy.

1S Ice Lake: See Backup for workloads and configurations. Results may vary.

Third party names are the property of their owners

# Five Epochs of Distributed Computing\*

Epoch starting date	Defining limitations	Application	Interaction time and Network performance	Capability
First 1970	Rare connections to expensive computers	FTP, telnet, email	100 ms Low bandwidth high latency	People to computers
Second 1984	I/O wall, disks can't keep up	RPC, Client Server	10 ms 10 mbps	Computer to computer
Third 1990	Networking wall	MPP HPC, three-tier datacenter networks	1 ms 100 mbs → 1 Gbs	Services to services
Fourth 2000	Dennard scaling wall ... per core plateau	Web search, planet-scale services	100 $\mu$ s 10 Gbps flash	People to people
Fifth 2015	Per socket wall ... accelerators take off	Machine Learning, data centric computing	10 $\mu$ s 200 Gbps → 1 Tbps	People to insights

\*The five Epochs of distributed computing, Amin Vahdat of Google: SIGCOMM Lifetime achievement award keynote, 2020.

# The Eight Fallacies of Distributed Computing

(Peter Deutsch of Sun Microsystems, 1994 ... item 8 added in 1997 by James Gosling)

Essentially everyone, when they first build a distributed application, makes the following eight assumptions. All prove to be false in the long run and all cause *big* trouble and *painful* learning experiences.

1. The network is reliable
2. Latency is zero
3. Bandwidth is infinite
4. The network is secure
5. Topology doesn't change
6. There is one administrator
7. Transport cost is zero
8. The network is homogeneous

# The Eight Fallacies of Distributed Computing

(Peter Deutsch of Sun Microsystems, 1994 ... item 8 added in 1997 by James Gosling)

Essentially everyone, when they first build a distributed application, makes the following eight assumptions. All prove to be false in the long run and all cause *big* trouble and *painful* learning experiences.

1. The network is reliable
2. Latency is **low and fixed**
3. Bandwidth is **high and fixed**
4. The network is secure
5. Topology doesn't change
6. There is one administrator
7. Transport cost is **negligible**
8. The network is homogeneous

# The Eight Fallacies of Distributed Computing

(Peter Deutsch of Sun Microsystems, 1994 ... item 8 added in 1997 by James Gosling)

Essentially everyone, when they first build a distributed application, makes the following eight assumptions. All prove to be false in the long run and all cause *big* trouble and *painful* learning experiences.

## Cloud

- ~~X. The network is reliable~~
- ~~X. Latency is **low and fixed**~~
- ~~X. Bandwidth is **high and fixed**~~
- ~~X. The network is secure~~
- ~~X. Topology doesn't change~~
- ~~X. There is one administrator~~
- ~~X. Transport cost is **negligible**~~
- ~~X. The network is homogeneous~~

## HPC Cluster

- ✓. The network is reliable**
- ✓. Latency is **low and fixed****
- ✓. Bandwidth is **high and fixed****
- ✓. The network is secure**
- ✓. Topology doesn't change**
- ✓. There is one administrator**
- ~~X. Transport cost is **negligible**~~
- ✓. The network is homogeneous**

# The four domains of parallel programming

Platform	GPU	Laptop or server	HPC Cluster	Cloud
Execution Agent	Work-items mapped onto an index space, execute as work-groups	Threads	Processes	Microservices
Memory	Hierarchical local/shared/global memory mapped to the index space	Single Address Space	Distributed memory, local memory owned by individual processes	Distributed object store (in memory) backed by a persistent storage system
Typical Execution Pattern	Streaming*	Fork-join	SPMD	Event driven tasks, FaaS, and Actors

GPU, Laptop/server, and cluster models work well together.

An impenetrable wall separates them from the cloud-native world

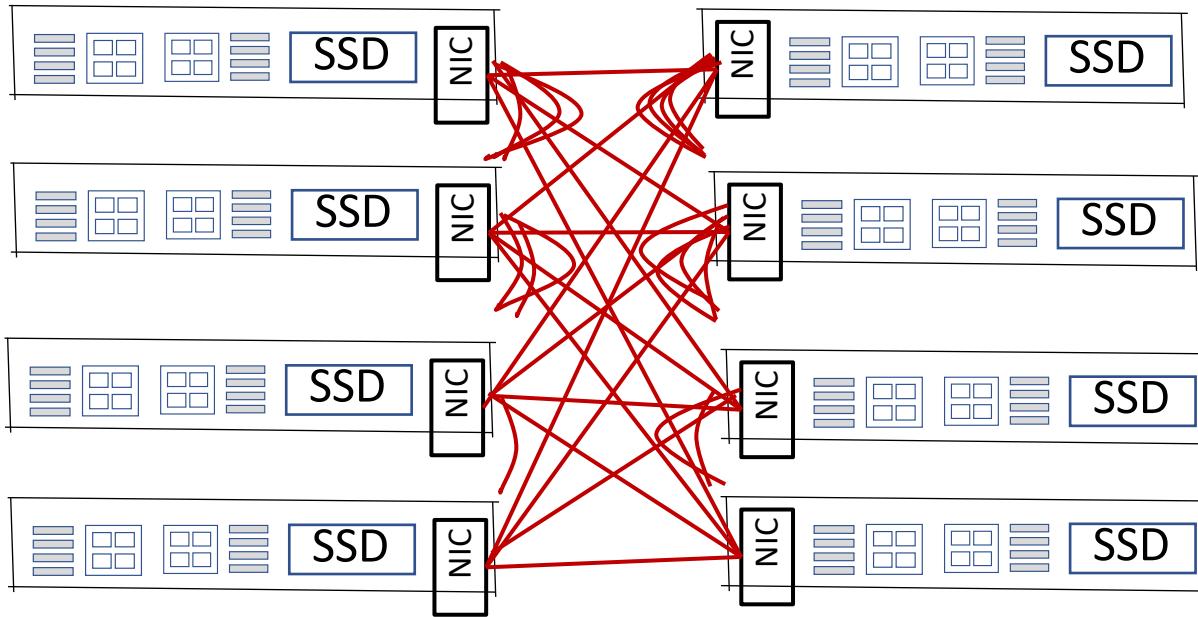
\* Often called by a newer and less accurate name, SIMT ... Single Instruction Multiple Thread

# The sixth Epoch of Distributed Computing

Epoch starting date	Defining limitations	Application	Interaction time and Network performance	Capability
First 1970	Rare connections to expensive computers	FTP, telnet, email	100 ms Low bandwidth high latency	People to computers
Second 1984	I/O wall, disks can't keep up	RPC, Client Server	10 ms 10 mbps	Computer to computer
Third 1990	Networking wall	MPP HPC, three-tier datacenter networks	1 ms 100 mbs → 1 Gbs	Services to services
Fourth 2000	Dennard Scaling Wall ... per core plateau	Web search, planet-scale services	100 $\mu$ s 10 Gbps flash	People to people
Fifth 2015	Per socket wall ... accelerators take off	Machine Learning, data centric computing	10 $\mu$ s 200 Gbps → 1 Tbps	People to insights
Sixth 2025	Speed of light	Dynamic, real-time AI, integrated from data-center to the edge with SDE*	100 ns 10 Tbs	People to experiences

\* SDE: Software defined Everything, i.e. software defined networking, software defined infrastructure, software defined servers ... All at the same time ... to dynamically construct systems to meet the needs of workloads.

# Networking technology... replace generic data center network with a cluster of cliques



A clique: A graph where every vertex is connected to every other vertex

A Clique: a network of diameter one with

$O(\frac{1}{4}N^2)$  bisection bandwidth

Combine with next generation optical networks to hit latencies of 100 ns

# Latencies every engineer should know ...

L1 cache reference 1.5 ns

L2 cache reference 5 ns

Branch misprediction 6 ns

Uncontented mutex lock/unlock 20 ns

L3 cache reference 25 ns

Main memory reference 100 ns

“Far memory”/Fast NVM reference 1,000 ns (1us)

Read 1 MB sequentially from memory 12,000 ns (12 us)

SSD Random Read 100,000 ns (100 us)

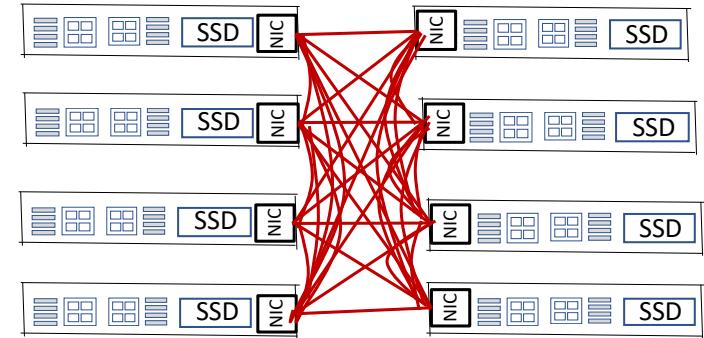
Read 1 MB bytes sequentially from SSD 500,000 ns (500 us)

Read 1 MB sequentially from 10Gbps network 1,000,000 ns (1 ms)

Read 1 MB sequentially from disk 10,000,000 ns (10 ms)

Disk seek 10,000,000 ns (10 ms)

Send packet California → Netherlands → California (150 ms)



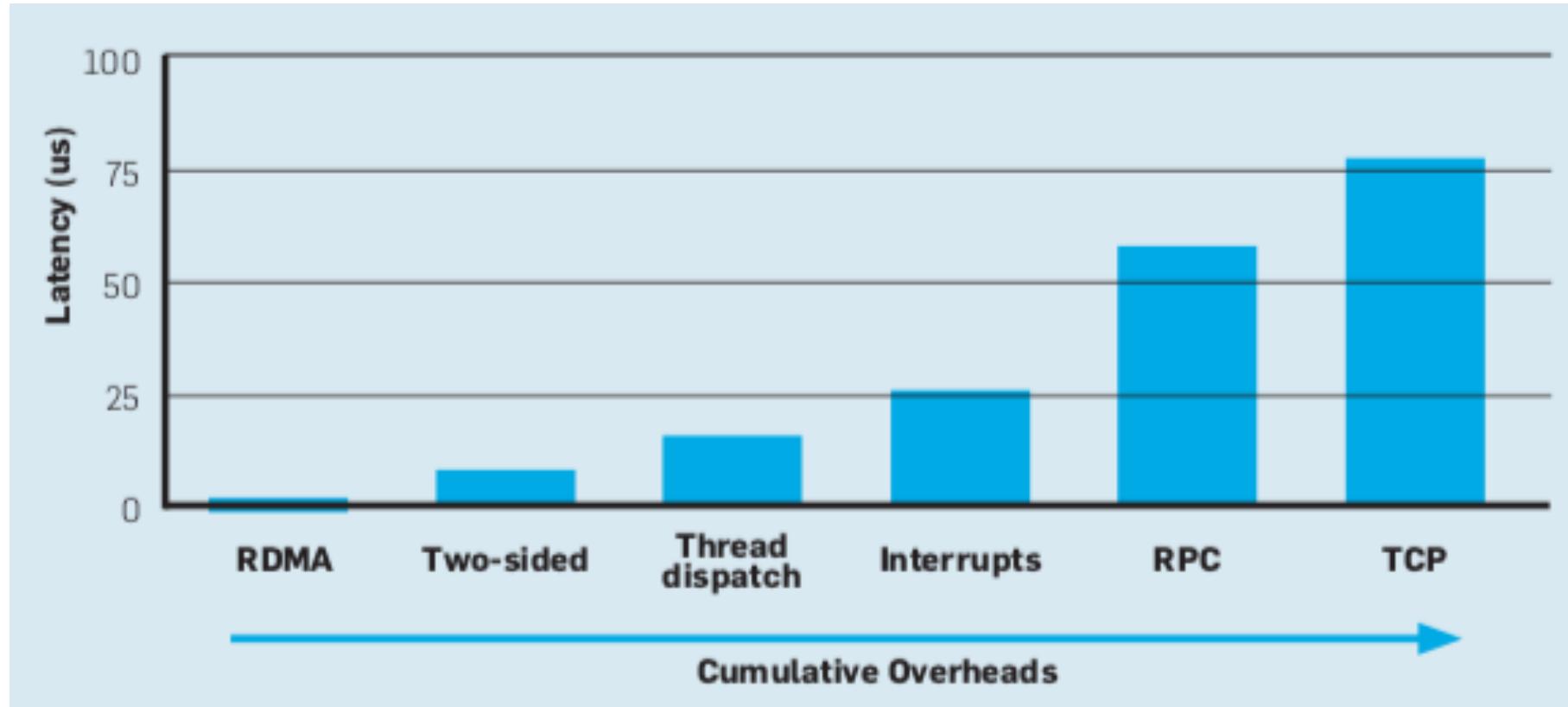
A cluster of nodes with a Clique network topology and low latency optical network...

Yields one hop network latencies on par with DRAM access latencies.

Source: **The Datacenter as a Computer: Designing Warehouse-Scale Machines**, Luiz Andre Barroso, Urs Holzle, Parthasarathy Ranganathan, 3<sup>rd</sup> edition, Morgan & Claypool, 2019.

# Take out the big stuff & you're left with lots of $\mu$ s overheads

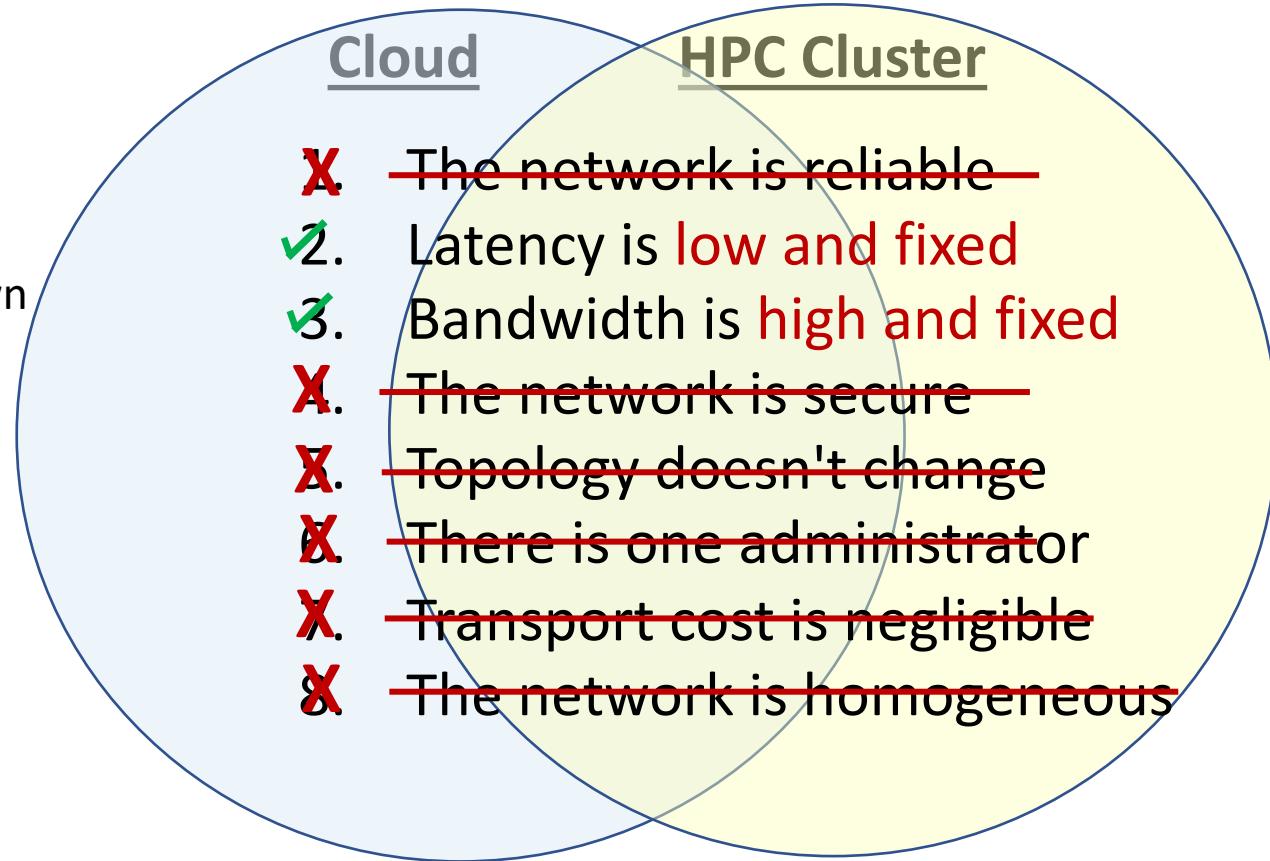
All those SW overheads add up ... like bricks that combine to build a networking-wall ...  
turning a 2  $\mu$ s network into a 100  $\mu$ s network...



Computer Scientists need to rethink system SW stacks to minimize latencies ... fast RDMA, reduce sync contention, low latency interrupt handlers, and more .... All to hit  $O(\mu\text{s})$  latencies.

# In the sixth Epoch of Distributed Computing, cloud and cluster overlap ... or even merge!

Chip-to-chip optical networks push latency down and bandwidth up



Data Streaming Accelerator reduces tail latency.

P4/P5/P6 + Infrastructure Processing Units drive down latency and reduces jitter

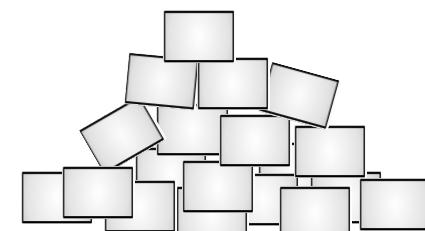
With Low Latencies, high bandwidths and stable performance, we can do loosely synchronous and synchronous applications in the cloud. The economics of the cloud vs dedicated HPC clusters means the cloud will dominate HPC

HPC applications will need to change to deal with reliability and network inhomogeneities.

# The four domains of parallel programming

Platform	GPU	Laptop or server	HPC Cluster	Cloud
Execution Agent	Work-items mapped onto an index space, execute as work-groups	Threads	Processes	Microservices
Memory	Hierarchical local/shared/global memory mapped to the index space	Single Address Space	Distributed memory, local memory owned by individual processes	Distributed object store (in memory) backed by a persistent storage system
Typical Execution Pattern	Streaming*	Fork-join	SPMD	Event driven tasks, FaaS, and Actors

Advances in networking technology plus low-overhead software stacks optimized to reduce tail-latency will shatter the wall



\* Often called by a newer and less accurate name, SIMD ... Single Instruction Multiple Thread

# The four domains of parallel programming

Platform	GPU	Laptop or server	HPC Cluster	Cloud
Execution Agent	Work-items mapped onto an index space, execute as work-groups	Threads	Processes	Microservices
Memory	Hierarchical local/shared/global memory mapped to the index space	Single Address Space	Distributed memory, local memory owned by individual processes	Distributed object store (in memory) backed by a persistent storage system
Typical Execution Pattern	Streaming*	Fork-join	SPMD	Event driven tasks, FaaS, and Actors

There will always be a need for top-end scalable systems in supercomputer centers, but economics will push the bulk of scientific computing into the cloud.

\* Often called by a newer and less accurate name, SIMT ... Single Instruction Multiple Thread

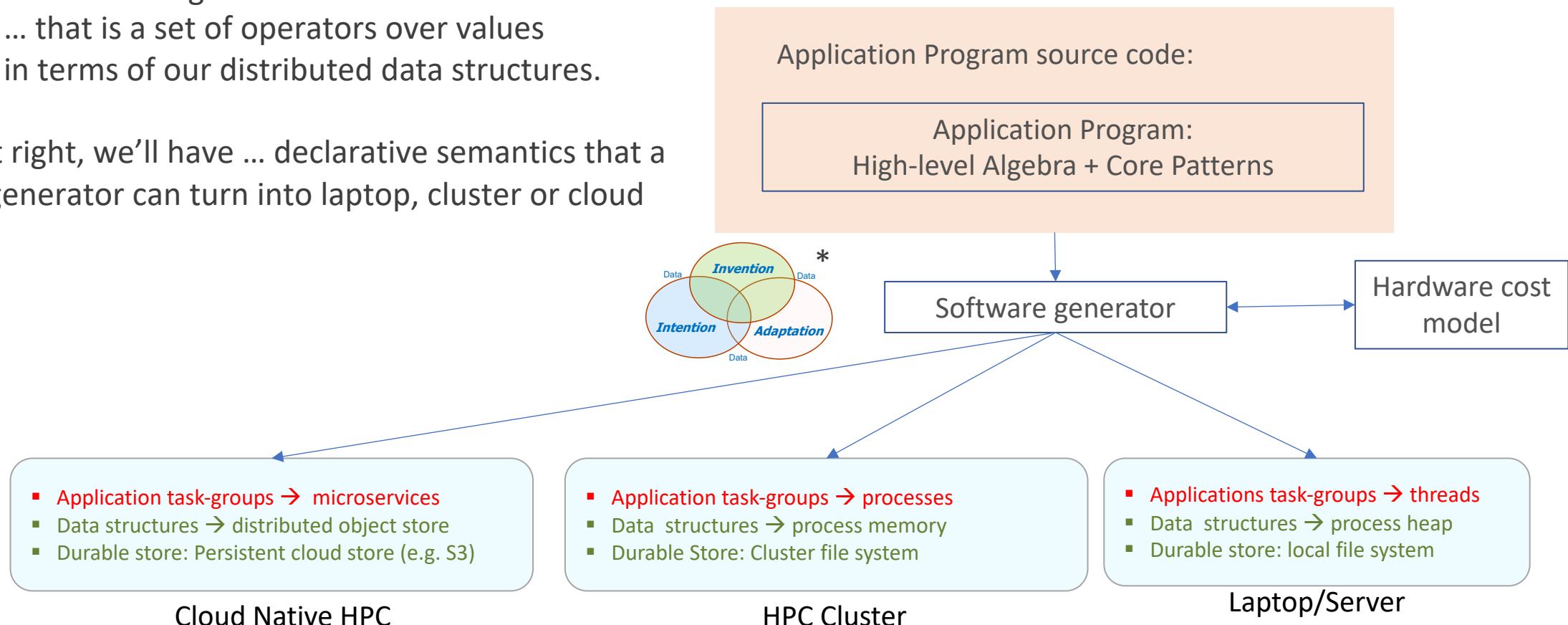
# The three domains of parallel programming

Platform*	Laptop or server	HPC Cluster	Cloud
Execution Agent	Threads	Processes	Microservices
Memory	Single Address Space	Distributed memory, local memory owned by individual processes	Distributed object store (in memory) backed by a persistent storage system
Typical Execution Pattern	Fork-join	SPMD	Event driven tasks, FaaS, and Actors

There will always be a need for top-end scalable systems in supercomputer centers, but economics will push the bulk of scientific computing into the cloud.

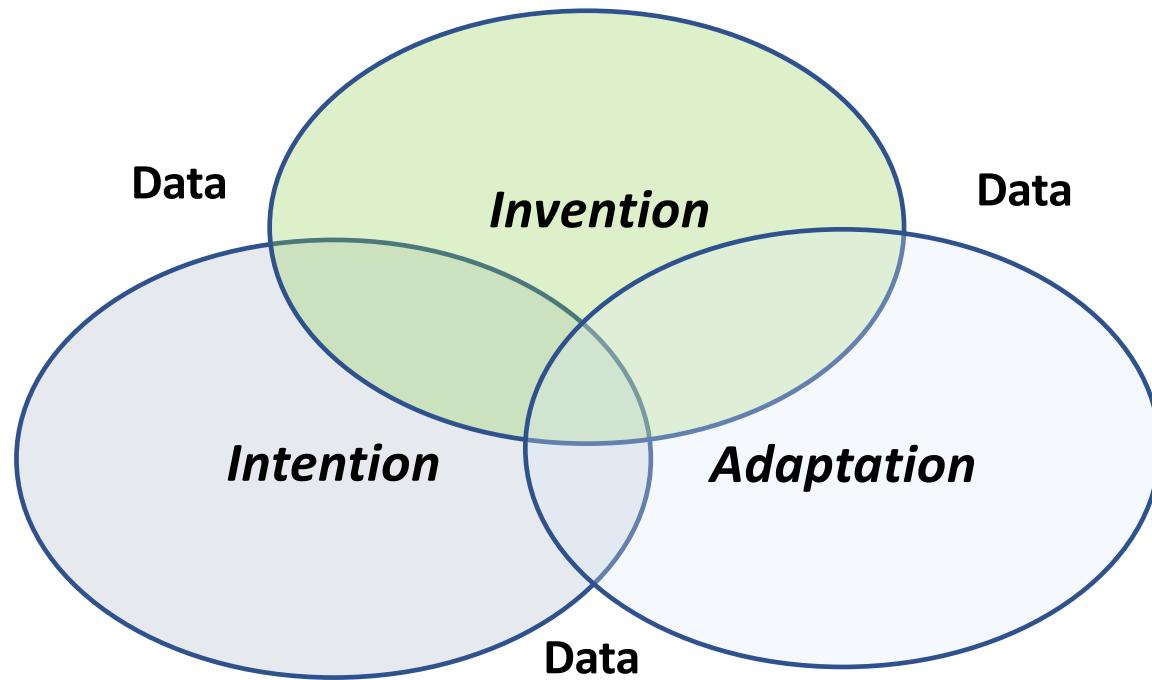
# One codebase → many systems

- Performance, Productivity AND Portability ... the database people “did it” with relational algebras and SQL.
- We can do it too with algebras over distributed data structures ... that is a set of operators over values expressed in terms of our distributed data structures.
- If we get it right, we’ll have ... declarative semantics that a software generator can turn into laptop, cluster or cloud programs.



\*This is the logo of the machine programming research program I help lead inside Intel Labs

# The Three Pillars of Machine Programming (MP)



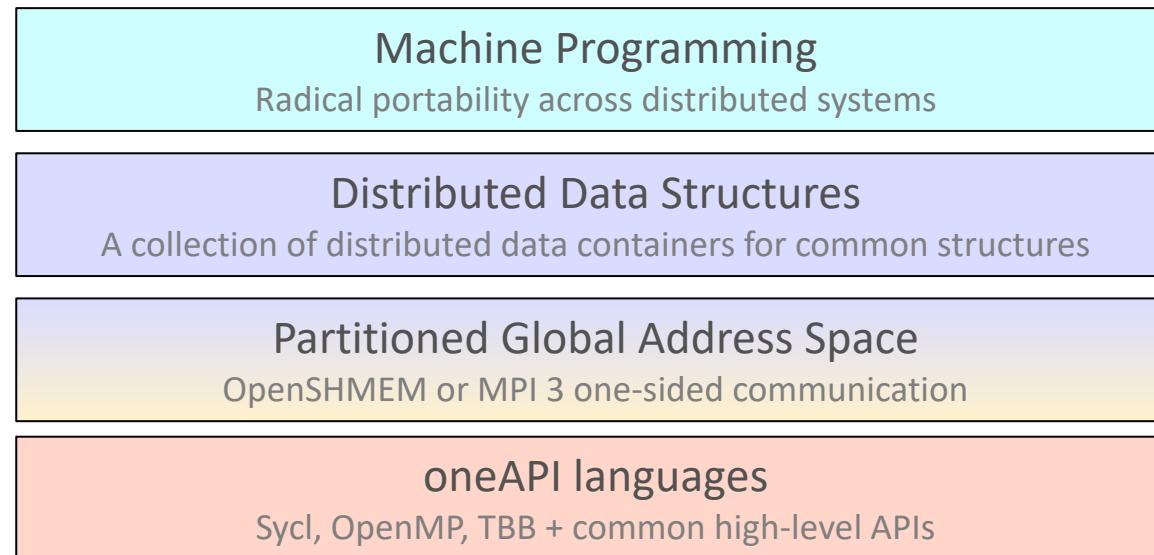
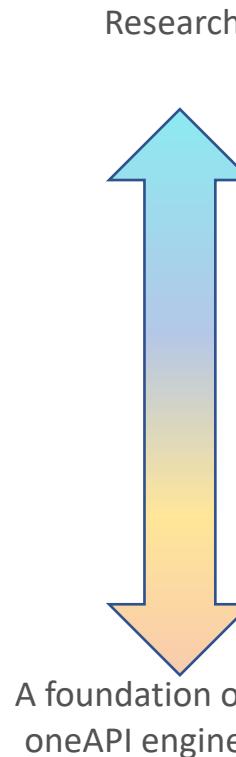
- **MP is the automation of software development**
  - **Intention:** Discover the intent of a programmer
  - **Invention:** Create new algorithms and data structures
  - **Adaptation:** Evolve in a changing hardware/software world

Justin Gottschlich, Intel Labs  
Armando Solar-Lezama, MIT  
Nesime Tatbul, Intel Labs  
Michael Carbin, MIT  
Martin Rinard, MIT  
Regina Barzilay, MIT  
Saman Amarasinghe, MIT  
Joshua B Tenenbaum, MIT  
Tim Mattson, Intel Labs

Summarized ~90 works.  
Key efforts by Berkeley,  
Google, Microsoft, MIT,  
Stanford, UW and others.

# oneAPI: A bridge to our heterogeneous/Distributed Future

My vision for how we bring oneAPI into a future dominated by power-optimized heterogenous chips organized into distributed systems.



The key to making this work ... the programmer is in control and chooses the level of abstraction based on the programming task.

# Summary

- Parallel computing is fun ... but it can be hard.
- Fortunately, if you stick to the Big-3 and the core patterns of parallel computing for HPC, it's not too overwhelming
  - The big 3: MPI, OpenMP, and “a GPU programming model”
  - Key Patterns: SPMD, loop level parallelism, geometric decomposition, divide and conquer, and SIMD
- Some day we'll automate the hard-parts with Machine Programming, but that may be 10 years!!!!

# SCANPY workload details and system configuration

Name	Intel® Xeon® Platinum 8380
Time	Jan 20, 2022
Manufacturer	Intel Corporation
Product Name	Intel® Xeon® Platinum 8380
BIOS Version	SE5C6200.86B.0020.P23.21032613 09
OS	Rocky Linux release 8.5 (Green Obsidian)
Kernel	4.18.0-240.22.1.el8_3.crt6.x86_64
Microcode	0xd000270
IRQ Balance	enabled
CPU Model	Intel(R) Xeon(R) Platinum 8380 CPU @ 2.30GHz
Base Frequency	2.3GHz
Maximum Frequency	3.4GHz
All-core Maximum Frequency	2.5GHz
CPU(s)	40
Thread(s) per Core	2
Core(s) per Socket	40

Socket(s)	1
NUMA Node(s)	1
Prefetchers	
Turbo	Enabled
PPIN(s)	
Power & Perf Policy	Performance
TDP	270 watts
Frequency Driver	
Frequency Governer	Performance
Frequency (MHz)	
Max C-State	
Installed	Intel® Xeon® Platinum 8380 40c D1 DDR4 16*16GB@3200MHz - Mellanox HDR
Huge Pages Size	2048 kB
Transparent Huge Pages	Always
Automatic NUMA Balancing	Enabled

- The following was done to optimize the SCANPY benchmark
  - Data preprocessing - used warm file cache and multi-threaded using Numba JIT
  - PCA, K-means, KNN – Used the Intel extension for scikit-learn.
  - t-SNE - Used optimized version from Intel's oneDAL Library.
  - Parallelized quadtree building, sorting and summarization steps using Morton codes.
  - UMAP - optimized the UMAP code using AVX512/AVX2. Used MKL for eigenvalue computation.
  - Louvain and Leiden algorithms – collaborated with Katana Graph to get well optimized versions and integrated them into SCANPY.

# Backup ... and a bit of extra content

- The future of parallel programming
- • The Jacobi solver case study
- Writing functions to call from inside a kernel

# Our running example: Jacobi solver

- An iterative method to solve a system of linear equations
  - Given a matrix A and a vector b find the vector x such that  $Ax=b$
- The basic algorithm:
  - Write A as a lower triangular (L), upper triangular (U) and diagonal matrix
$$Ax = (L+D+U)x = b$$
  - Carry out multiplications and rearrange
$$Dx = b - (L+U)x \rightarrow x = (b - (L+U)x)/D$$
  - Iteratively compute a new x using the x from the previous iteration
$$X_{\text{new}} = (b - (L+U)x_{\text{old}})/D$$
- Advantage: we can easily test if the answer is correct by multiplying our final x by A and comparing to b
- Disadvantage: It takes many iterations and only works for diagonally dominant matrices

# Jacobi Solver

Iteratively update xnew until the value stabilizes (i.e. change less than a preset TOL)

```
<<< allocate and initialize the matrix A >>>
<<< and vectors x1, x2 and b           >>>

while((conv > TOL) && (iters<MAX_ITERS))
{
    iters++;
    for (i=0; i<Ndim; i++){
        xnew[i] = (TYPE) 0.0;
        for (j=0; j<Ndim;j++){
            if(i!=j)
                xnew[i]+= A[i*Ndim + j]*xold[j];
        }
        xnew[i] = (b[i]-xnew[i])/A[i*Ndim+i];
    }

    // test convergence
    conv = 0.0;
    for (i=0; i<Ndim; i++){
        tmp  = xnew[i]-xold[i];
        conv += tmp*tmp;
    }
    conv = sqrt((double)conv);

    // swap pointers for next
    // iteration
    TYPE* tmp = xold;
    xold = xnew;
    xnew = tmp;

} // end while loop
```

# Jacobi Solver (Parallel Target/loop, 1/2)

```
while((conv > TOL) && (iters<MAX_ITERS))
{
    iters++;
#pragma omp target map(tofrom:xnew[0:Ndim],xold[0:Ndim]) \
    map(to:A[0:Ndim*Ndim], b[0:Ndim])
#pragma omp loop
for (i=0; i<Ndim; i++){
    xnew[i] = (TYPE) 0.0;
    for (j=0; j<Ndim;j++){
        if(i!=j)
            xnew[i]+= A[i*Ndim + j]*xold[j];
    }
    xnew[i] = (b[i]-xnew[i])/A[i*Ndim+i];
}
```

# Jacobi Solver (Parallel Target/loop, 2/2)

```
//  
// test convergence  
//  
conv = 0.0;  
#pragma omp target map(to:xnew[0:Ndim],xold[0:Ndim]) \  
    map(tofrom:conv)  
#pragma omp loop private(i,tmp) reduction(+:conv)  
for (i=0; i<Ndim; i++){  
    tmp = xnew[i]-xold[i];  
    conv += tmp*tmp;  
}  
conv = sqrt((double)conv);  
TYPE* tmp = xold;  
xold = xnew;  
xnew = tmp;  
} // end while loop
```

This worked but the performance was awful. Why?

System	Implementation	Ndim = 4096
NVIDIA® K20X™ GPU	Target dir per loop	131.94 secs

Cray® XC40™ Supercomputer running Cray® Compiling Environment 8.5.3.  
Intel® Xeon ® CPU E5-2697 v2 @ 2.70GHz with 32 GB DDR3. NVIDIA® Tesla® K20X, 6GB.

# Data movement dominates!!!

```
while((conv > TOLERANCE) && (iters<MAX_ITERS))  
{ iters++;  
xnew = iters % s ? x2 : x1;  
xold = iters % s ? x1 : x2;
```

Typically over 4000 iterations!

```
#pragma omp target map(tofrom:xnew[0:Ndim],xold[0:Ndim]) \  
map(to:A[0:Ndim*Ndim], b[0:Ndim] )
```

```
#pragma omp loop private(i,j)  
for (i=0; i<Ndim; i++){  
    xnew[i] = (TYPE) 0.0;  
    for (j=0; j<Ndim;j++){  
        if(i!=j)  
            xnew[i]+= A[i*Ndim + j]*xold[j];  
    }  
    xnew[i] = (b[i]-xnew[i])/A[i*Ndim+i];  
}
```

```
// test convergence  
conv = 0.0;
```

```
#pragma omp target map(to:xnew[0:Ndim],xold[0:Ndim]) \  
map(tofrom:conv)
```

```
#pragma loop reduction(+:conv)  
for (i=0; i<Ndim; i++){  
    tmp = xnew[i]-xold[i];  
    conv += tmp*tmp;  
}
```

```
conv = sqrt((double)conv);  
}
```

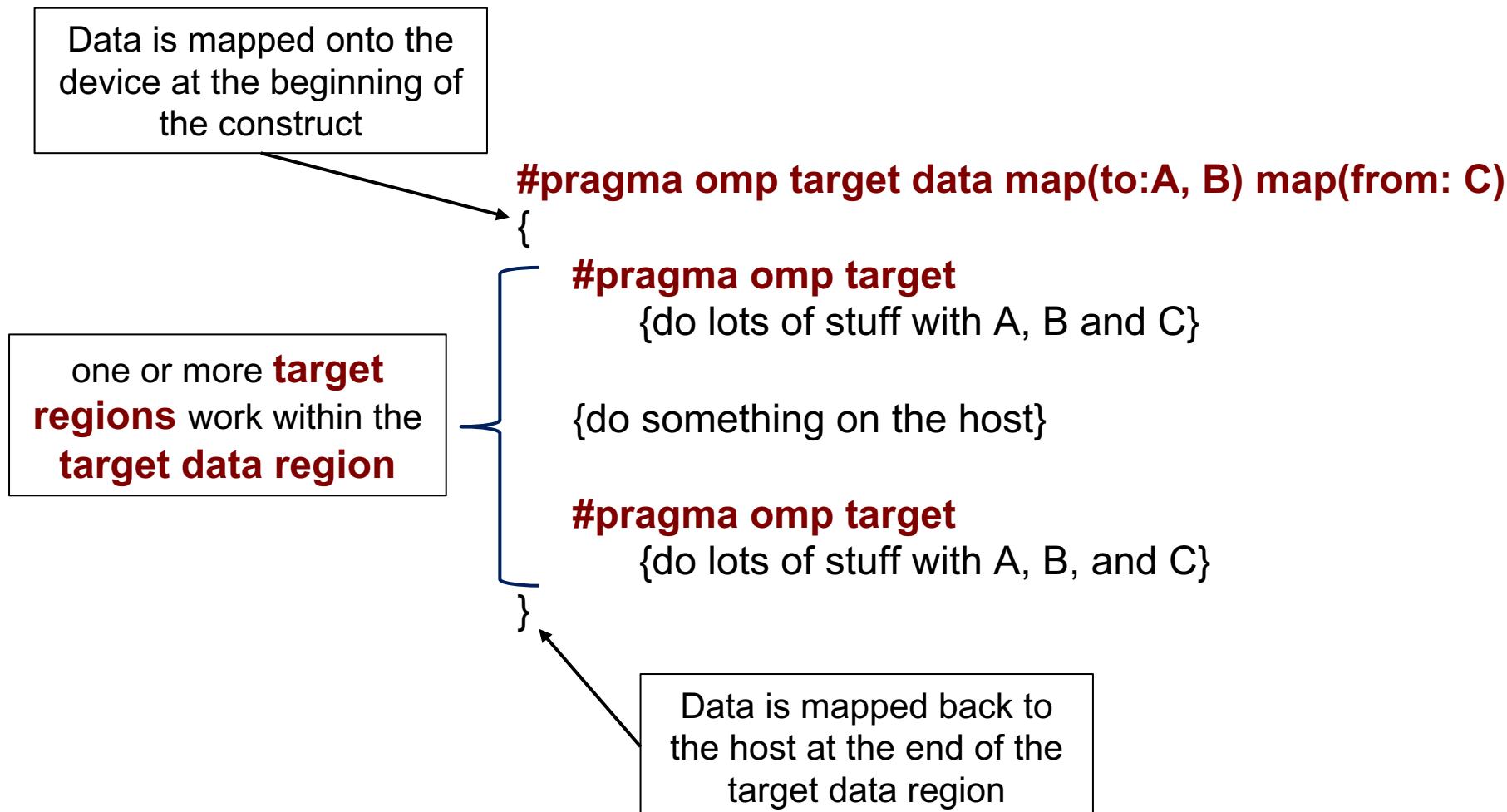
For each iteration, **copy to device**  
 $(3*Ndim+Ndim^2)*\text{sizeof}(\text{TYPE})$  bytes

For each iteration, **copy from device**  
 $2*Ndim*\text{sizeof}(\text{TYPE})$  bytes

For each iteration, **copy to device**  
 $2*Ndim*\text{sizeof}(\text{TYPE})$  bytes

# Target data directive

- The **target data** construct creates a target data region
  - ... use **map** clauses for explicit data management



# Jacobi Solver (Par Target Data, 1/2)

```
#pragma omp target data map(tofrom:xold[0:Ndim],xnew[0:Ndim]) \
    map(to:A[0:Ndim*Ndim], b[0:Ndim] ,Ndim)
while((conv > TOL) && (iters<MAX_ITERS))
{ iters++;
  #pragma omp target
  #pragma omp loop private(j) firstprivate(xnew,xold)
  for (i=0; i<Ndim; i++){
    xnew[i] = (TYPE) 0.0;
    for (j=0; j<Ndim;j++){
      if(i!=j)
        xnew[i]+= A[i*Ndim + j]*xold[j];
    }
    xnew[i] = (b[i]-xnew[i])/A[i*Ndim+i];
  }
}
```

# Jacobi Solver (Par Target Data, 2/2)

```
// test convergence
conv = 0.0;
#pragma omp target map(tofrom: conv)
#pragma omp loop private(tmp) firstprivate(xnew,xold) reduction(+:conv)

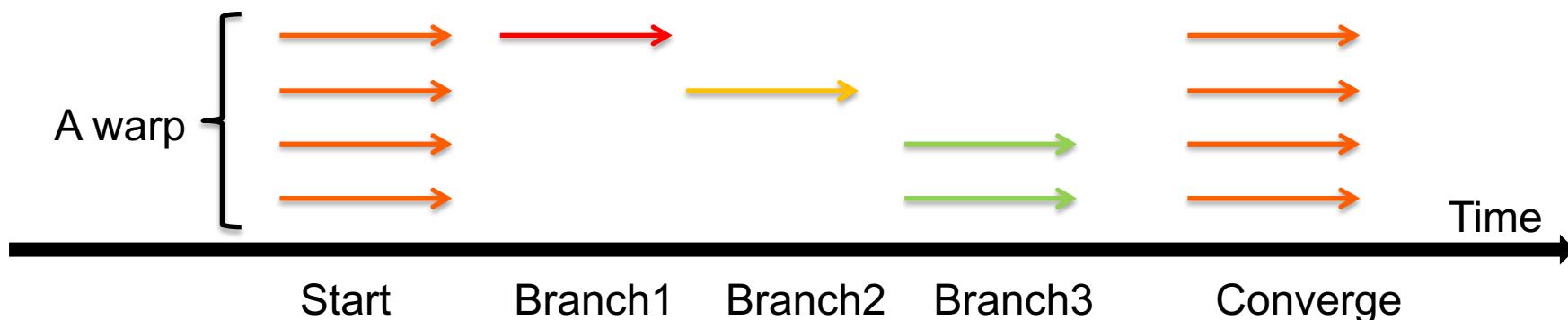
for (i=0; i<Ndim; i++){
    tmp = xnew[i]-xold[i];
    conv += tmp*tmp;
}
// end target region
conv = sqrt((double)conv);

TYPE* tmp = xold;
xold = xnew;
xnew = tmp;
} // end while loop
```

System	Implementation	Ndim = 4096
NVIDIA® K20X™ GPU	Target dir per loop	131.94 secs
	Above plus target data region	18.37 secs

# Single Instruction Multiple Data

- Individual work-items of a warp start together at the same program address
- Each work-item has its own instruction address counter and register state
  - Each work-item is free to branch and execute independently
  - Supports the SPMD pattern.
- Branch behavior
  - Each branch will be executed serially
  - Work-items not following the current branch will be disabled



# Branching

## Conditional execution

```
// Only evaluate expression  
// if condition is met  
if (a > b)  
{  
    acc += (a - b*c);  
}
```

## Selection and masking

```
// Always evaluate expression  
// and mask result  
temp = (a - b*c);  
mask = (a > b ? 1.f : 0.f);  
acc += (mask * temp);
```

# Coalescence

- Coalesce - to combine into one
- Coalesced memory accesses are key for high bandwidth
- Simply, it means, if thread  $i$  accesses memory location  $n$  then thread  $i+1$  accesses memory location  $n+1$
- In practice, it's not quite as strict...

```
for (int id = 0; id < size; id++)
{
    // ideal
    float val1 = memA[id];

    // still pretty good
    const int c = 3;
    float val2 = memA[id + c];

    // stride size is not so good
    float val3 = memA[c*id];

    // terrible
    const int loc =
        some_strange_func(id);

    float val4 = memA[loc];
}
```

# Jacobi Solver (Target Data/branchless/coalesced mem, 1/2)

```
#pragma omp target data map(tofrom:x1[0:Ndim],x2[0:Ndim]) \
    map(to:A[0:Ndim*Ndim], b[0:Ndim] ,Ndim)
while((conv > TOL) && (iters<MAX_ITERS))
{ iters++;
#pragma omp target
    #pragma omp loop private(j)
    for (i=0; i<Ndim; i++){
        xnew[i] = (TYPE) 0.0;
        for (j=0; j<Ndim;j++){
            xnew[i]+= (A[j*Ndim + i]*xold[j])*( (TYPE) (i != j));
        }
        xnew[i] = (b[i]-xnew[i])/A[i*Ndim+i];
    }
}
```

We replaced the original code with a poor memory access pattern  
 $xnew[i]+= (A[i*Ndim + j]*xold[j])$   
With the more efficient  
 $xnew[i]+= (A[j*Ndim + i]*xold[j])$

# Jacobi Solver (Target Data/branchless/coalesced mem, 2/2)

```
//  
// test convergence  
conv = 0.0;  
#pragma omp target map(tofrom: conv)  
#pragma omp loop private(tmp) reduction(+:conv)  
for (i=0; i<Ndim; i++){  
    tmp = xnew[i]-xold[i];  
    conv += tmp*tmp;  
}  
conv = sqrt((double)conv);  
TYPE* tmp = xold;  
xold = xnew;  
xnew = tmp;  
} // end while loop
```

System	Implementation	Ndim = 4096
NVIDIA® K20X™ GPU	Target dir per loop	131.94 secs
	Above plus target data region	18.37 secs
	Above plus reduced branching	13.74 secs
	Above plus improved mem access	7.64 secs

# A more complicated example: Jacobi iteration: OpenACC (GPU)

```
#pragma acc data copy(A), create(Anew)
while (err>tol && iter < iter_max) {
    err = 0.0;
    #pragma acc parallel loop reduction(max:err)
    for(int j=1; j< n-1; j++) {
        for(int i=1; i<M-1; i++) {
            Anew[j][i] = 0.25* (A[j][i+1] + A[j][i-1] +
                                  A[j-1][i] + A[j+1][i]);
            err = max(err,abs(Anew[j][i] - A[j][i]));
        }
    }
    #pragma acc parallel loop
    for(int j=1; j< n-1; j++) {
        for(int i=1; i<M-1; i++) {
            A[j][i] = Anew[j][i];
        }
    }
    iter++;
}
```

Create a data region on the GPU. Copy A once onto the GPU, and create Anew on the device (no copy from host)

Copy A back out to host  
... but only once

# A more complicated example:

## Jacobi iteration: OpenMP target directives

```
#pragma omp target data map(A) map(alloc:Anew)
while (err>tol && iter < iter_max) {
    err = 0.0;
    #pragma target
    #pragma omp teams loop reduction(max:err)
    for(int j=1; j< n-1; j++) {
        for(int i=1; i<M-1; i++) {
            Anew[j][i] = 0.25* (A[j][i+1] + A[j][i-1] +
                                  A[j-1][i] + A[j+1][i]);
            err = max(err,abs(Anew[j][i] - A[j][i]));
        }
    }
    #pragma omp target
    #pragma omp teams loop
    for(int j=1; j< n-1; j++) {
        for(int i=1; i<M-1; i++) {
            A[j][i] = Anew[j][i];
        }
    }
    iter++;
}
```

Create a data region on the GPU. Map A and Anew onto the target device

Copy A back out to host  
... but only once

# Backup ... and a bit of extra content

- The future of parallel programming
- The Jacobi solver case study
- • Writing functions to call from inside a kernel

# Defining a function to be called from inside a kernel

Tell OpenMP  
to compile this  
function for the  
GPU (and the  
CPU)

```
#include mm_utils.h

#pragma omp declare target
void ddot(double *C, double *A, double *B, int i, int j, int Mdim, int Pdim){
    for(int k=0;k<Pdim;k++){
        /* C(i,j) = sum(over k) A(i,k) * B(k,j) */
        *(C+(i*Mdim+j)) += *(A+(i*Pdim+k)) * *(B+(k*Mdim+j));
    }
}
#pragma omp end declare target

void mm_gpu(int Ndim, int Mdim, int Pdim, TYPE *A, TYPE *B, TYPE *C){
    int i, j, k;

#pragma omp target teams map(tofrom:C[0:Ndim*Mdim]) map(to:B[0:Pdim*Mdim],A[0:Ndim*Pdim])
#pragma omp loop collapse(2)
    for (i=0; i<Ndim; i++){
        for (j=0; j<Mdim; j++){
            ddot(C, A, B, i, j, Mdim, Pdim);
        }
    }
}
```

Call inside a  
target region, and  
the GPU version  
is called.

Call on the host,  
and a CPU  
version is called

This is in my file mm\_gpu.c