# Lecture: Manycore GPU Architectures and Programming, Part 4
## -- Introducing OpenMP and HOMP for Accelerators

## CSCE 569 Parallel Computing

Department of Computer Science and Engineering
Yonghong Yan
yanyh@cse.sc.edu
https://passlab.github.io/CSCE569/

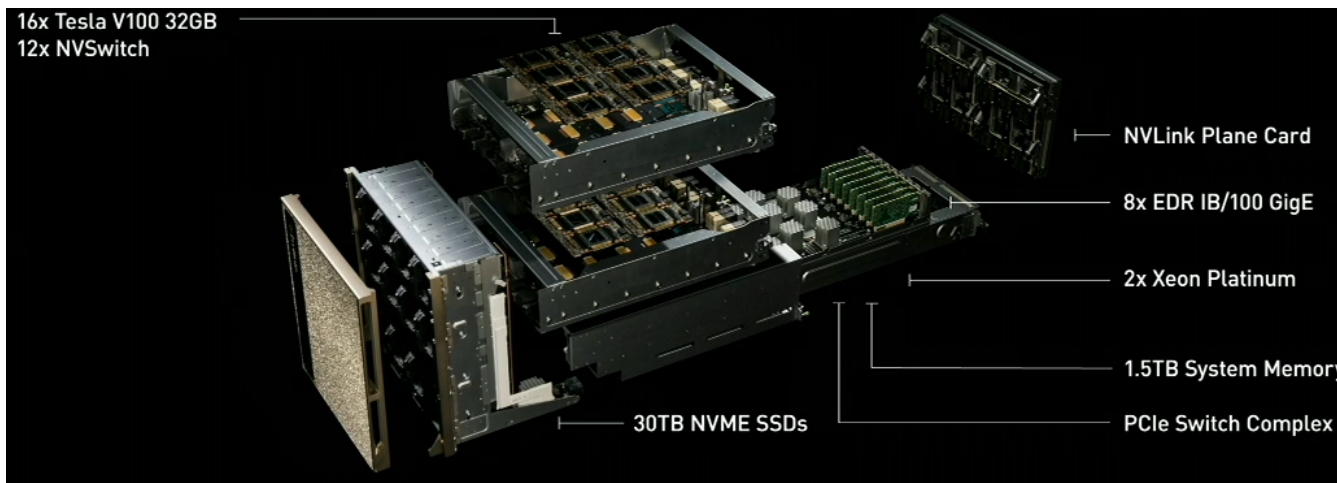# Manycore GPU Architectures and Programming: Outline

- Introduction
  - GPU architectures, GPGPUs, and CUDA
- GPU Execution model
- CUDA Programming model
- Working with Memory in CUDA
  - Global memory, shared and constant memory
- Streams and concurrency
- CUDA instruction intrinsic and library
- Performance, profiling, debugging, and error handling
- ☛ **Directive-based high-level programming model**
  - **OpenMP and OpenACC**

# HPC Systems with Accelerators

- Accelerator architectures become popular
  - GPUs and Xeon Phi
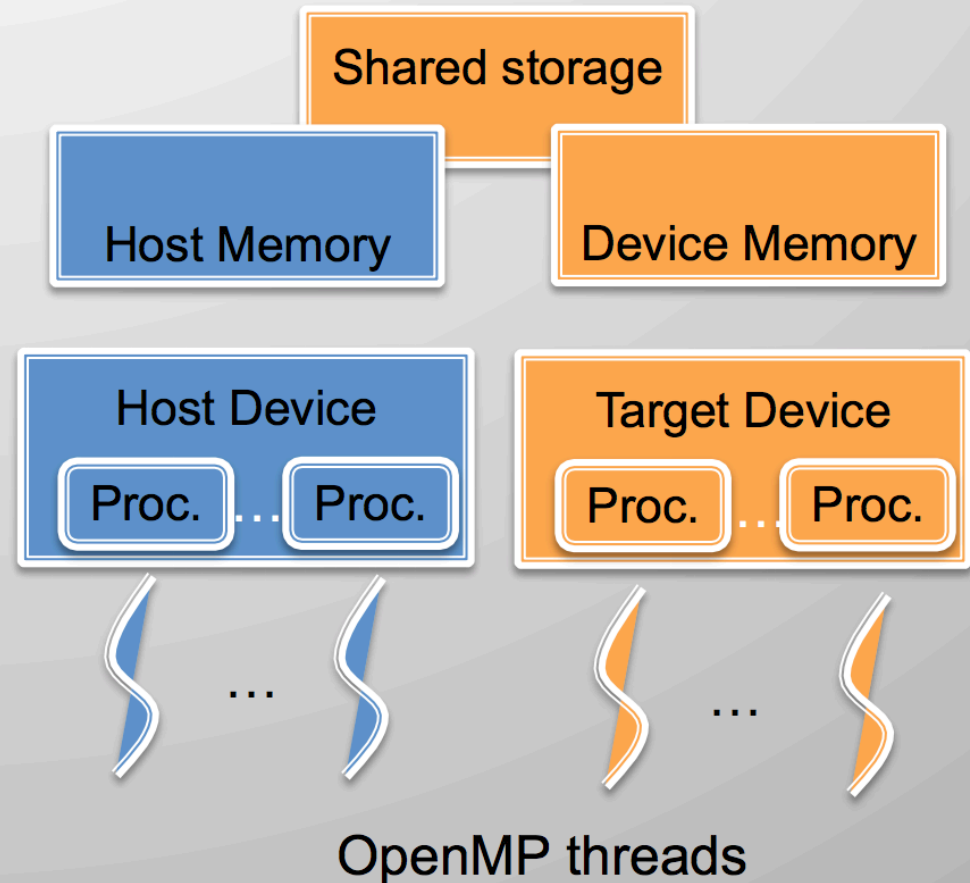- Multiple accelerators are common
  - 2, 4, or 8

**Programming on NVIDIA GPUs**

1. CUDA and OpenCL
   - **Low-level**
2. Library, e.g. cublas, cufft, cuDNN
3. **OpenMP, OpenACC, and others**
   - **Rely on compiler support**
4. Application framework
   - **TensorFlow, etc**



16x Tesla V100 32GB
12x NVSwitch

NVLink Plane Card

8x EDR IB/100 GigE

2x Xeon Platinum

1.5TB System Memory

PCIe Switch Complex

30TB NVME SSDs

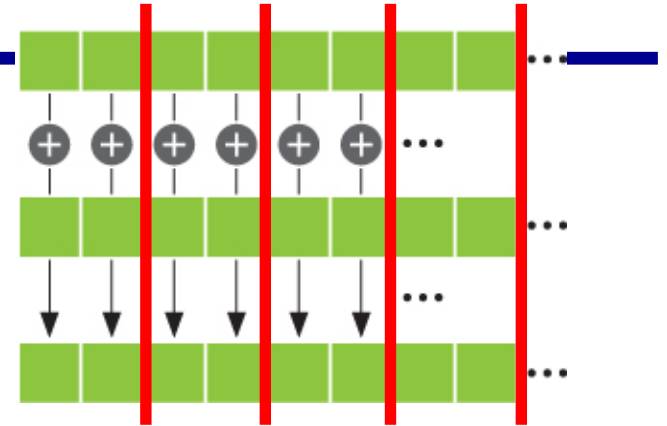https://www.anandtech.com/show/12587/nvidias-dgx2-sixteen-v100-gpus-30-tb-of-nvme-only-400k

# OpenMP 4.0 for Accelerators

- Device: a logical execution engine
  - Host device: where OpenMP program begins, one only
  - Target devices: **1 or more** accelerators

- Memory model
  - Host data environment: one
  - Device data environment: one or more
  - Allow shared host and device memory

- Execution model: Host-centric
  - Host device : "offloads" code regions and data to accelerators/target devices
  - Target Devices: still fork-join model
  - Host waits until devices finish
  - Host executes device regions if no accelerators are available /supported



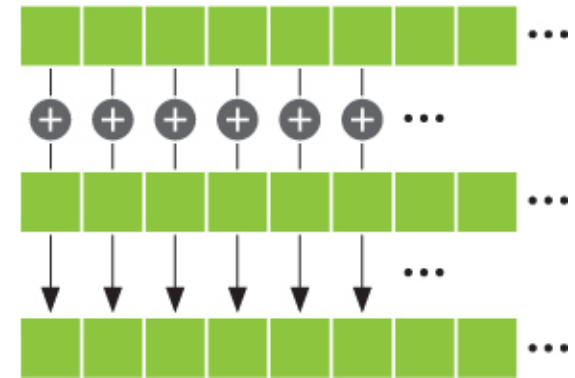OpenMP threads

# AXPY Example with OpenMP: Multicore

- $y = \alpha \cdot x + y$
  - x and y are vectors of size n
  - $\alpha$ is scalar

```
1  void axpy(REAL *x, REAL *y, long n, REAL a) {
2      #pragma omp parallel for shared(x, y, n, a)
3      for (int i = 0; i < n; ++i)
4          y[i] += a * x[i];
5  }
```

- **Data (x, y and a) are shared**
  - Parallelization is relatively easy
- Other examples
  - sum: reduction
  - Stencil: halo region exchange and synchronization

# AXPY Offloading To a GPU using CUDA

```
1  // CUDA kernel. Each thread takes care of one element of c
2  __global__ void axpy(REAL *x, REAL *y, int n, REAL a) {
3      int id = blockIdx.x*blockDim.x+threadIdx.x;
4      if (id < n) y[id] += a * x[id];
5  }
6
7  int main( int argc, char* argv[] ) {
8
9      // ... init host a, x and y
10     // Allocate memory for each vector on GPU
11     cudaMalloc(&d_x, size);
12     cudaMalloc(&d_y, size);
13
14     // Copy host vectors to device
15     cudaMemcpy( d_x, h_x, size, cudaMemcpyHostToDevice);
16     cudaMemcpy( d_y, h_y, size, cudaMemcpyHostToDevice);
17
18     int blockSize, gridSize;
19     blockSize = 1024;
20     gridSize = (int)ceil((float)n/blockSize);
21     axpy<<<gridSize, blockSize>>>(d_x, d_y, n, a);
22
23     // Copy array back to host
24     cudaMemcpy( h_y, d_y, size, cudaMemcpyDeviceToHost );
25
26     // Release device memory
27     cudaFree(d_x);
28     cudaFree(d_y);
29 }
```

**Memory allocation on device**

**Memcpy from host to device**

**Launch parallel execution**

**Memcpy from device to host**

**Deallocation of dev memory**

6

# AXPY Example with OpenMP: single device

- y = α·x + y
  - x and y are vectors of size n
  - α is scalar

```
1  void axpy_ompacc(REAL* x, REAL* y, int n, REAL a) {
2    #pragma omp target device (0) map(tofrom: y[0:n]) \
3         map(to: x[0:n],a,n)
4    #pragma omp parallel for shared(x, y, n, a)
5    for (int i = 0; i < n; ++i)
6      y[i] += a * x[i];
7  }
```

- **target** directive: annotate an offloading code region
- **map** clause: map data between host and device → moving data
  - **to|tofrom|from**: mapping directions
  - Use array region

# OpenMP Computation and Data Offloading

- #pragma omp target  *device(id) map() if()*
  - **target**: create a data environment and offload computation on the device
  - **device (int_exp)**: specify a target device
  - **map(to|from|tofrom|alloc:var_list)** : data mapping between the current data environment and a device data environment

- #pragma target data *device (id) map() if()*
  - Create a device data environment: to be reused/inherited

# target and map Examples

```
void vec_mult(int N)
{
    int i;
    float p[N], v1[N], v2[N];
    init(v1, v2, N);
    #pragma omp target map(to: v1, v2) map(from: p)
    #pragma omp parallel for
    for (i=0; i<N; i++)
        p[i] = v1[i] * v2[i];
    output(p, N);
}
```

```
void vec_mult(float *p, float *v1, float *v2, int N)
{
    int i;
    init(v1, v2, N);
    #pragma omp target map(to: v1[0:N], v2[:N]) map(from: p[0:N])
    #pragma omp parallel for
    for (i=0; i<N; i++)
        p[i] = v1[i] * v2[i];
    output(p, N);
}
```

# Accelerator: Explicit Data Mapping

- Relatively small number of truly shared memory accelerators so far
- Require the user to explicitly *map* data to and from the device memory
- Use array region

```
long a = 0x858;
long b = 0;
int anArray[100]

#pragma omp target data map(to:a) \\
   map(tofrom:b,anArray[0:64])
{
   /* a, b and anArray are mapped
    * to the device */

   /* work on the device */
#pragma omp target …
   {
       …
   }|
}
/* b and anArray are mapped
 * back to the host */
```

# target date Example

```
void vec_mult(float *p, float *v1, float *v2, int N)
{
    int i;
    init(v1, v2, N);
    #pragma omp target data map(from: p[0:N])
    {
        #pragma omp target map(to: v1[:N], v2[:N])
        #pragma omp parallel for
        for (i=0; i<N; i++)
            p[i] = v1[i] * v2[i];
        init_again(v1, v2, N);
        #pragma omp target map(to: v1[:N], v2[:N])
        #pragma omp parallel for
        for (i=0; i<N; i++)
            p[i] = p[i] + (v1[i] * v2[i]);
    }
    output(p, N);
}
```

Note mapping inheritance

# Accelerator: Hierarchical Parallelism

- Organize massive number of threads
  - teams of threads, e.g. map to CUDA grid/block
- Distribute loops over teams

```
#pragma omp target

#pragma omp teams num_teams(2)
      num_threads(8)
{
    //-- creates a "league" of teams
    //-- only local barriers permitted
#pragma omp distribute
for (int i=0; i<N; i++) {


}


}
```

Only target directive makes
it as accelerator region

# teams and distribute Loop Example

```
float dotprod_teams(float B[], float C[], int N, int num_blocks,
  int block_threads)
{
    float sum = 0;
    int i, i0;
    #pragma omp target map(to: B[0:N], C[0:N])
    #pragma omp teams num_teams(num_blocks) thread_limit(block_threads)
      reduction(+:sum)
    #pragma omp distribute
    for (i0=0; i0<N; i0 += num_blocks)
      #pragma omp parallel for reduction(+:sum)
      for (i=i0; i< min(i0+num_blocks,N); i++)
          sum += B[i] * C[i];
    return sum;
}
```

Double-nested loops are mapped to the two levels of thread hierarchy (league and team)

# Jacobi Example: The Impact of Compiler Transformation to Performance

```
#pragma omp target data device (gpu0) map(to:n, m, omega, ax, ay, b, \
    f[0:n][0:m])  map(tofrom:u[0:n][0:m]) map(alloc:uold[0:n][0:m])

while ((k<=mits)&&(error>tol))
{
// a loop copying u[][] to uold[][] is o
 …

#pragma omp target  device(gpu(
    uold[0:n][0:m])  map(tofrom:u
#pragma omp parallel for private
for (i=1;i<(n-1);i++)
  for (j=1;j<(m-1);j++)
  {
    resid = (ax*(uold[i-1][j] + uold[i+1
       + ay*(uold[i][j-1] + uold[i][j+1])+
    u[i][j] = uold[i][j] - omega * resid;
    error = error + resid*resid ;
  } // the rest code omitted  ...
}
```



Jacobi Execution Time (s)

Legend:
- first version
- target-data
- Loop collapse using linearization with static-even scheduling
- Loop collapse using 2-D mapping (16x16 block)
- Loop collapse using 2-D mapping (8x32 block)
- Loop collapse using linearization with round-robin scheduling

Matrix size (float): 128x128, 256x256, 512x512, 1024x1024, 2048x2048

# Mapping Nested Loops to GPUs

- Need to achieve coalesced memory access on GPUs

```
#pragma acc loop gang(2) vector(2)
for ( i = x1; i < X1; i++ ) {
#pragma acc loop  gang(3) vector(4)
for ( j = y1; j < Y1; j++ ) {...... }
}
```





Fig. 9: Double nested loop mapping.     Fig. 10: Triple nested loop mapping.

15

# Compiler vs Hand-Written

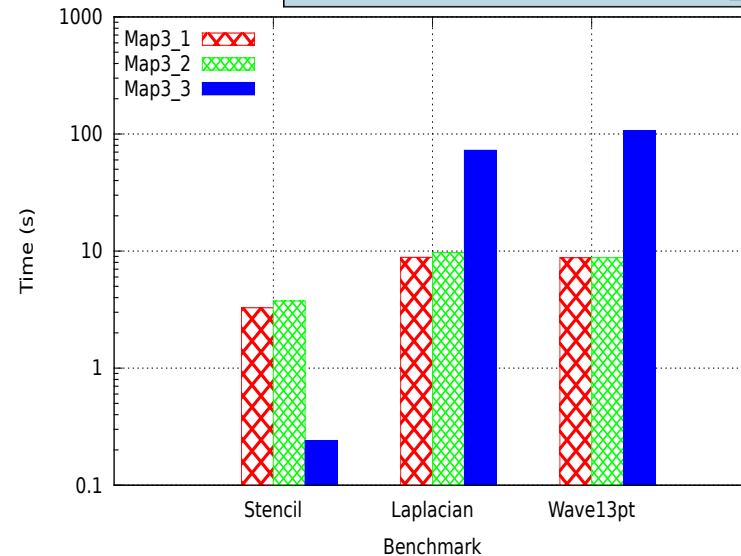| Applications | Domains | OpenACC Directive Combinations | Lines of Code Added vs Serial | | Speedup Over | | |
| --- | --- | --- | --- | --- | --- | --- | --- |
| | | | OpenMP | OpenACC | Seq | OpenMP | CUDA |
| Needleman-Wunsch | Bioinformatics | data copy, copyin<br>kernels present<br>loop gang, vector, private | 6 | 5 | 2.98 | 1.28 | 0.24 |
| Stencil | Cellular Automation | data copyin, copy, deviceptr<br>kernels present<br>loop collapse, independent | 1 | 3 | 40.55 | 15.87 | 0.92 |
| Computational Fluid Dyanmics (CFD) | Fluid Mechanics | data copyin, copy, deviceptr<br>data present, deviceptr<br>kernels deviceptr<br>kernels loop, gang, vector, private<br>loop gang, vector<br>acc_malloc(), acc_free() | 8 | 46 | 35.86 | 4.59 | 0.38 |
| 2D Heat (grid size 4096*4096) | Heat Conduction | data copyin, copy, deviceptr<br>kernels present<br>loop collapse, independent | 1 | 3 | 99.52 | 28.63 | 0.90 |
| Clever (10Ovals) | Data Mining | data copyin<br>kernels present, create, copyin, copy<br>loop independent | 10 | 3 | 4.25 | 1.22 | 0.60 |
| FeldKemp (FDK) | Image Processing | kernels copyin, copyout<br>loop, collapse, independent | 1 | 2 | 48.30 | 6.51 | 0.75 |

NAS Parallel Benchmarks for GPGPUs using a Directive-based Programming Model, Rengan Xu, Xiaonan Tian, Sunita Chandrasekaran, Yonghong Yan and Barbara Chapman 27th

International Workshop on Languages and Compilers for Parallel Computing (LCPC2014)

# AXPY Example with OpenMP: Multiple device

```c
void axpy_omp_mdev(REAL* x, REAL* y, int n, REAL
  int ndev = omp_get_num_devices();
  #pragma omp parallel num_threads(ndev)
  {
    int devid = omp_get_thread_num();
    int start, size, remnant;
    remnant = n % ndev; size = n / ndev;
    if (devid < remnant) {
      size++; start = size*devid;
    } else start = size*devid+remnant;
    #pragma omp target device (devid) \
        map(tofrom: y[start:size]) \
        map(to: x[start:size],a,size)
    #pragma omp parallel for shared(x, y, si
      for (int i = 0; i < size; ++i)
        y[i] += a * x[i];
  }
}
```



- Parallel region
  - One CPU thread per device
- Manually partition array x and y
- Each thread offload subregion of x and y
- Chunk the loop in alignment with the data partition

1. Copy input data from CPU memory to GPU memory
2. Load GPU program and execute, caching data on chip for performance
3. Copy results from GPU memory to CPU memory

# Hybrid OpenMP (HOMP) for Multiple Accelerators

```c
1  /* align computation with data using ALIGN(x)*/
2  void axpy_homp_v1(REAL* x, REAL* y, int n, REAL a) {
3    #pragma omp parallel target device (*) \
4        map(tofrom: y[0:n] distribute(BLOCK)) \
5        map(to: x[0:n] distribute(BLOCK),a,n)
6    #pragma omp parallel for distribute(ALIGN(x))
7    for (int i = 0; i < n; ++i)
8      y[i] += a * x[i];
9  }
10
11 /* align data with computation using ALIGN*/
12 void axpy_homp_v2(REAL* x, REAL* y, int n, REAL a) {
13   #pragma omp parallel target device (*) \
14       map(tofrom: y[0:n] distribute(ALIGN(loop))) \
15       map(to: x[0:n] distribute(ALIGN(loop)),a,n)
16   #pragma omp parallel for distribute(AUTO)
17 loop:  for (int i = 0; i < n; ++i)
18      y[i] += a * x[i];
19 }
```

HOMP: Automated Distribution of Parallel Loops and Data in Highly Parallel Accelerator-Based Systems, Yonghong Yan, Jiawen Liu, and Kirk W. Cameron, The IEEE International Parallel & Distributed Processing Symposium (IPDPS) 2017

# Three Challenges to Implement HOMP

1. Load balance when distributing loop iterations across computational different devices (CPU, GPU, and MIC)
   - We developed 7 algorithms of loop distribution and the runtime select algorithms based on computation/data intensity

2. Only copy the associated data to the device that are needed for the loop chunks assigned to that device
   - Runtime support for ALIGN interface to move or share data between memory spaces

1. Select devices for computations for the optimal performance because more devices ≠ better performance
   - CUTOFF ratio to select device

# Offloading Execution Time (ms) on 2 CPUs + 4 GPUs + 2 MICs and using CUTOFF_RATIO

**Execution Time (ms) on 2 CPUs + 4 GPUs + 2 MICs**

The algorithm that delivers the best performance without CUTOFF

The algorithm with 15% CUTOFF_RATIO that delivers the best performance and its speedup against the best algorithm that does not use CUTOFF

## sum-300M

| Algorithm | Value |
|---|---|
| MODEL_PROFILE_AUTO(15% CUTOFF) | 100.93 — 2.09 |
| MODEL_PROFILE_AUTO | 291.86 |
| SCHED_PROFILE_AUTO | 407.61 |
| MODEL_2_AUTO | 545.66 |
| MODEL_1_AUTO | 261.99 |
| SCHED_GUIDED | 779.83 |
| SCHED_DYNAMIC | 211.11 |
| BLOCK | 752.22 |

## stencil2d-256

| Algorithm | Value |
|---|---|
| MODEL_1_AUTO(15% CUTOFF) | 432.17 — 3.43 |
| MODEL_2_AUTO | 1504.45 |
| MODEL_1_AUTO | 1482.90 |
| BLOCK | 5054.33 |

## matvec-48k

| Algorithm | Value |
|---|---|
| SCHED_PROFILE_AUTO(15% CUTOFF) | 709.63 — 0.56 |
| MODEL_PROFILE_AUTO | 555.21 |
| SCHED_PROFILE_AUTO | 793.04 |
| MODEL_2_AUTO | 953.50 |
| MODEL_1_AUTO | 1060.35 |
| SCHED_DYNAMIC | 400.75 |
| BLOCK | 1459.78 |

## matul-6144

| Algorithm | Value |
|---|---|
| MODEL_2_AUTO(15% CUTOFF) | 1422.96 — 2.68 |
| MODEL_PROFILE_AUTO | 6532.80 |
| SCHED_PROFILE_AUTO | 10393.55 |
| MODEL_2_AUTO | 3544.63 |
| MODEL_1_AUTO | 3809.59 |
| SCHED_GUIDED | 21587.16 |
| SCHED_DYNAMIC | 3664.08 |
| BLOCK | 20989.67 |

## bm2d-256

| Algorithm | Value |
|---|---|
| MODEL_1_AUTO(15% CUTOFF) | 272.80 — 1.01 |
| MODEL_PROFILE_AUTO | 1723.84 |
| SCHED_PROFILE_AUTO | 1695.45 |
| MODEL_2_AUTO | 274.32 |
| MODEL_1_AUTO | 277.19 |
| SCHED_GUIDED | 885.80 |
| SCHED_DYNAMIC | 3508.77 |

## axpy-10B

| Algorithm | Value |
|---|---|
| MODEL_PROFILE_AUTO (15% CUTOFF) | 306.26 — 1.35 |
| MODEL_PROFILE_AUTO | 514.63 |
| SCHED_PROFILE_AUTO | 767.02 |
| MODEL_2_AUTO | 1017.77 |
| MODEL_1_AUTO | 653.72 |
| SCHED_DYNAMIC | 412.29 |
| BLOCK | 1423.10 |

TOTAL OFF TIME(ms): 0.00, 200.00, 400.00, 600.00, 800.00, 1000.00, 1200.00, 1400.00, 1600.00, 1800.00, 2000.00

# Speedup From CUTOFF

- Apply 15% CUTOFF ratio to modeling and profiling
  - Only those devices who may compute more than 15% of total iterations will be used
    - Thinking of 8 devices (1/8 = 12.5%)

| Benchmarks | Devices used | CUTOFF Speedup |
|------------|--------------|----------------|
| axpy-10B | 2 CPU + 4 GPUs | 1.35 |
| bm2d-256 | 2 CPU + 4 GPUs | 1.01 |
| matul-6144 | 4 GPUs | 2.68 |
| matvec-48k | **4 GPUs** | 0.56 |
| stencil2d-256 | 4 GPUs | 3.43 |
| sum-300M | 2 CPUs + 4 GPUs | 2.09 |