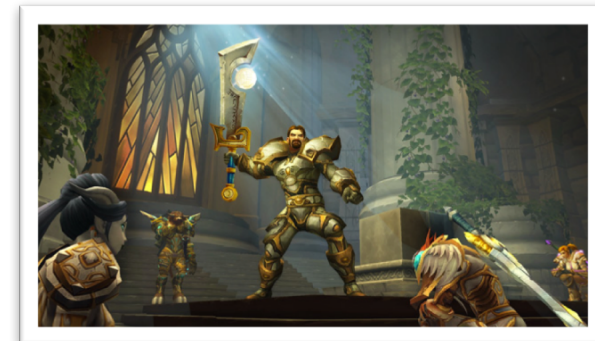
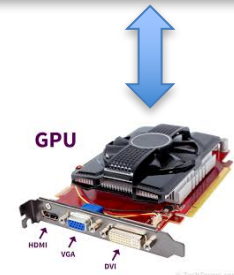
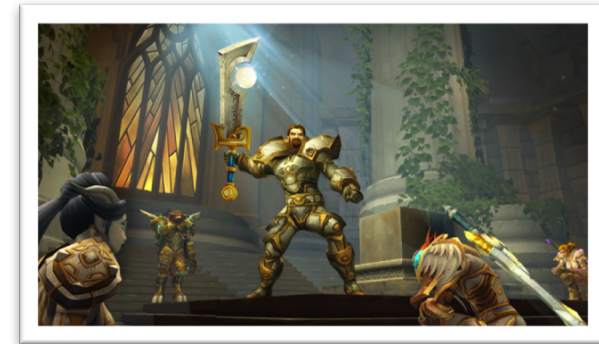
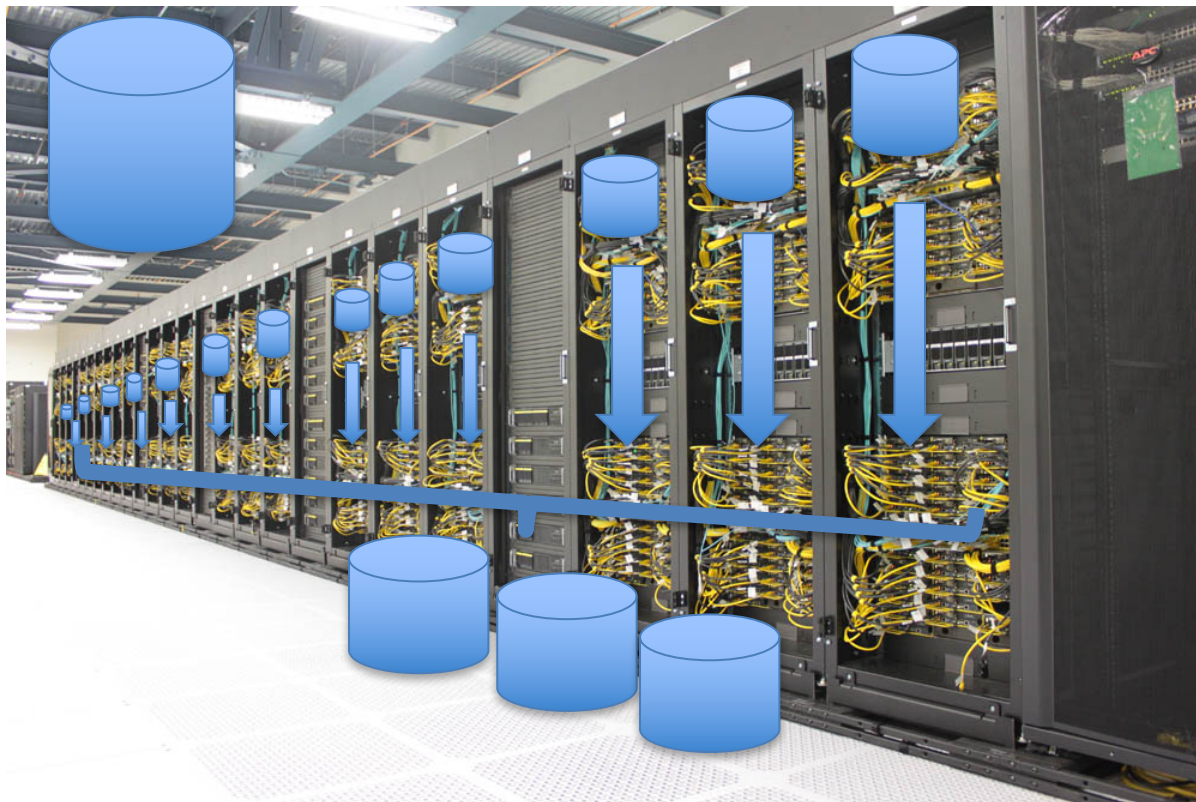


Big ML and GPUs



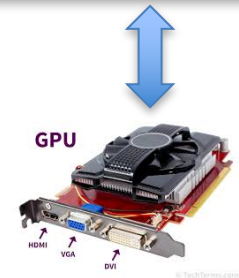
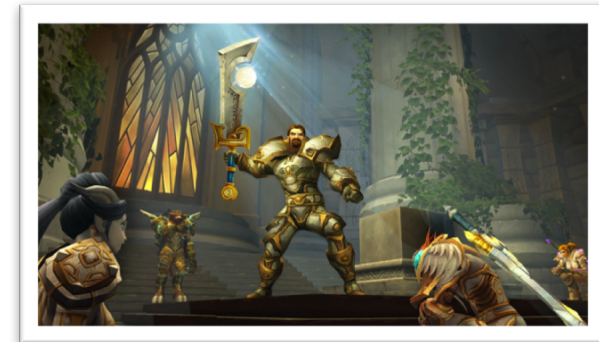
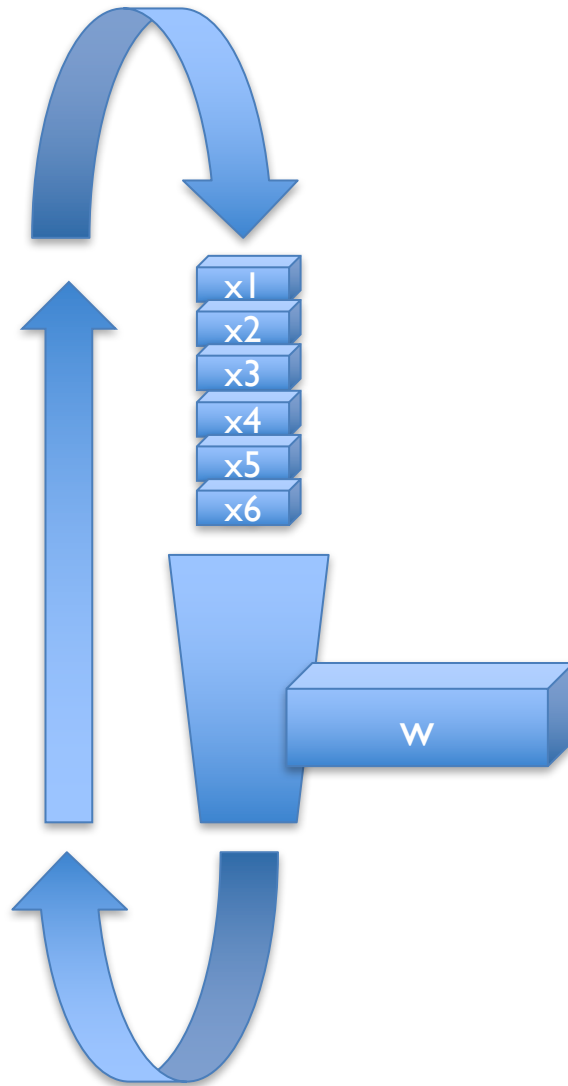
Parallel computing with map-reduce:

- Stream-and-sort in parallel
- Enormous datasets
- Tasks are i/o bound
- Many unreliable processors
 - which are basically commodity PCs
- Parallelize with mapreduce
 - loosely coupled, heavy-weight jobs
 - communicate via network/disk
- Don't iterate (typically)



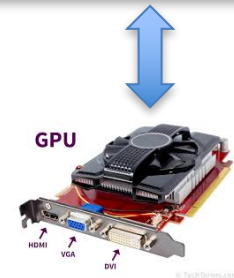
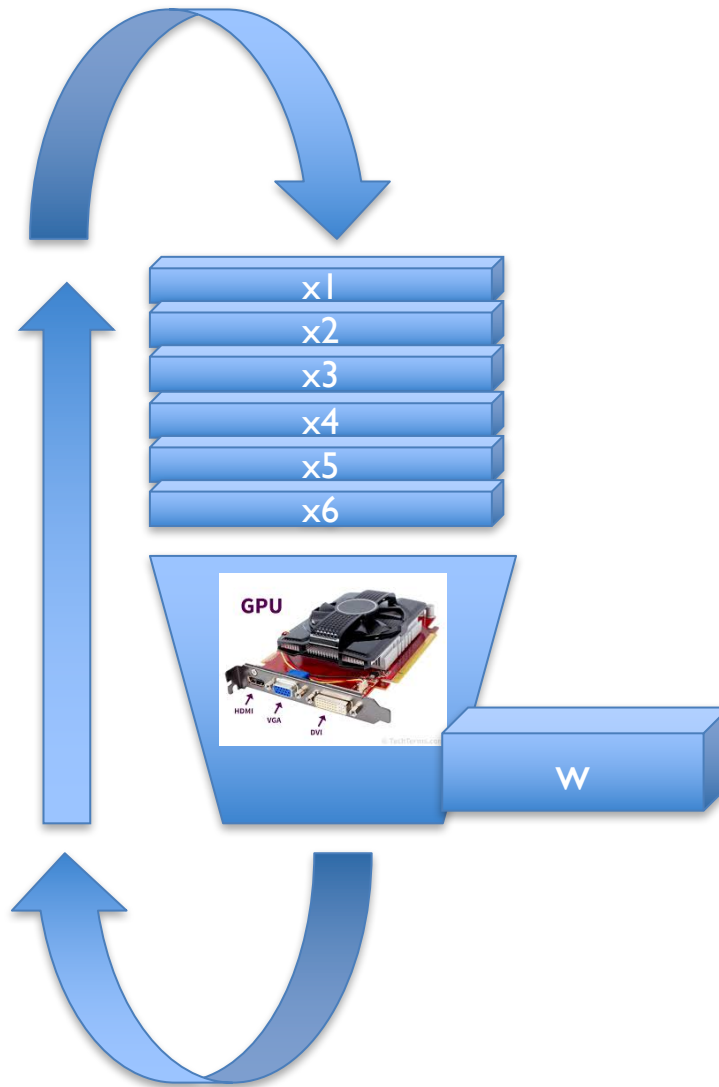
Streaming SGD:

- Iterative
- Sequential
- Fast
- Scale up by bounding memory
- You can handle very large datasets ... but slowly



Streaming SGD:

- Iterative
- Sequential
- Fast
- Scale up by bounding memory
- You can handle very large datasets ... but slowly
- You can speed it up by making the tasks in the stream bigger and doing them in **parallel**
- A GPU is a good way of doing that

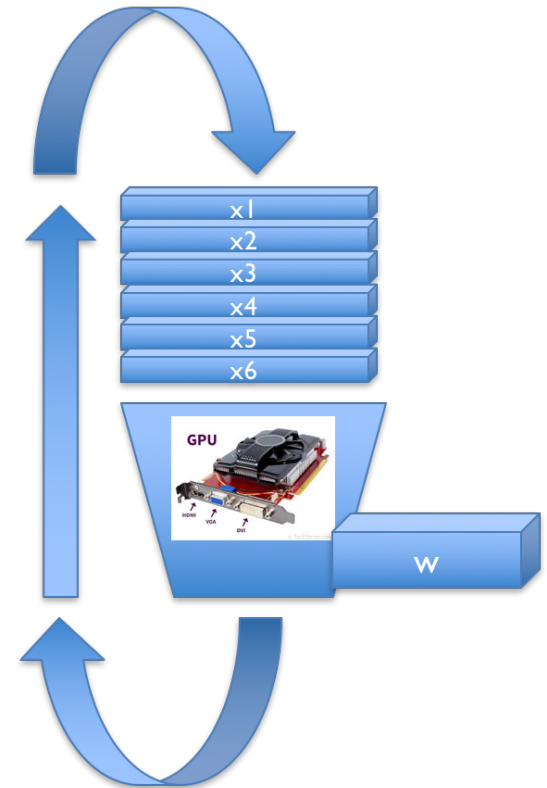
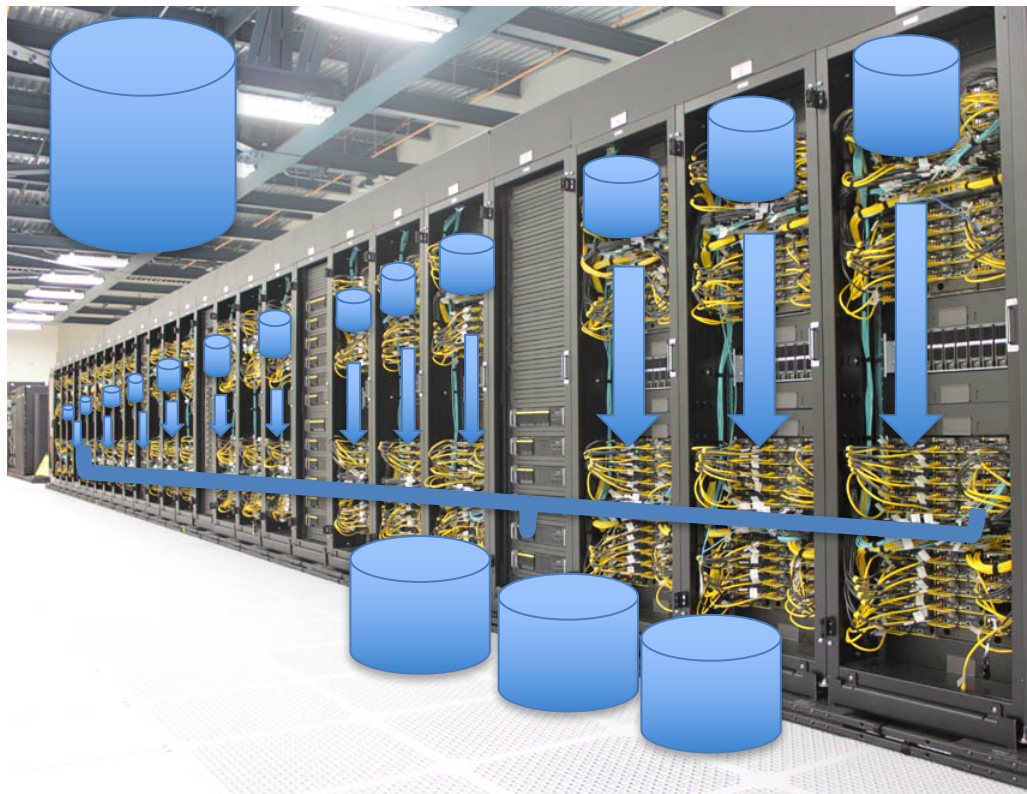


Parallel computing with map-reduce:

- Stream-and-sort in parallel
- Enormous datasets
- Tasks are i/o bound
- Many unreliable processors
 - which are basically commodity PCs
- Parallelize with mapreduce
 - loosely coupled, heavy-weight jobs
 - communicate via network/disk
- Don't iterate (typically)

Parallel ML computing with GPUS:

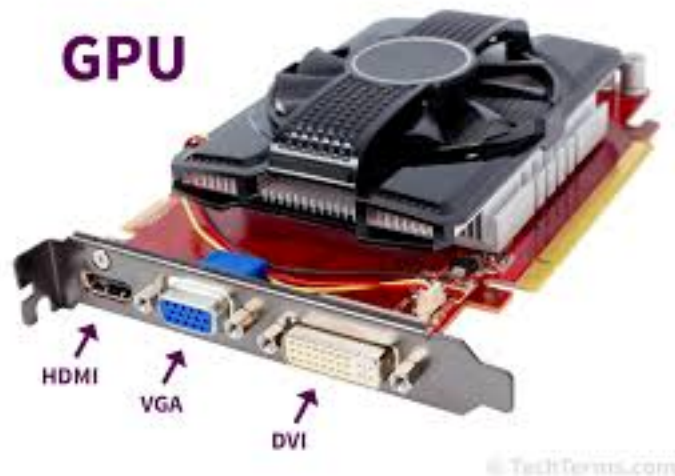
- Iterative streaming ML in parallel
- Big-but-not-too-big datasets
- Tasks are compute bound
- Many fast-but-simple processors
- Replace streaming operations with medium-sized computations that can be done in parallel
- Usually iterate many times



WHAT ARE GPUS?

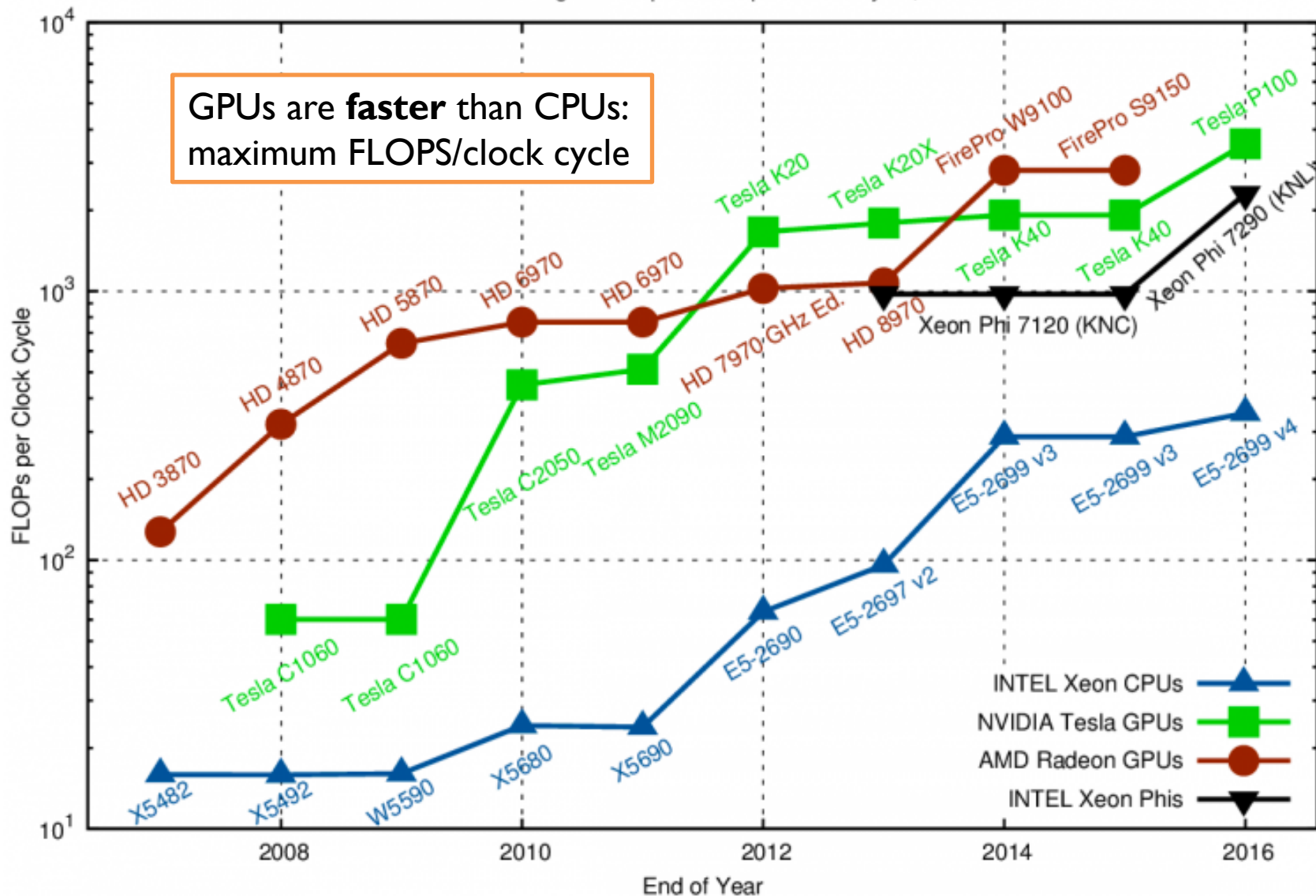
What is a GPU?

A **graphics processing unit (GPU)** is a specialized [electronic circuit](#) designed to rapidly manipulate and alter [memory](#) to accelerate the creation of [images](#) in a [frame buffer](#) intended for output to a [display device](#). [wikipedia]



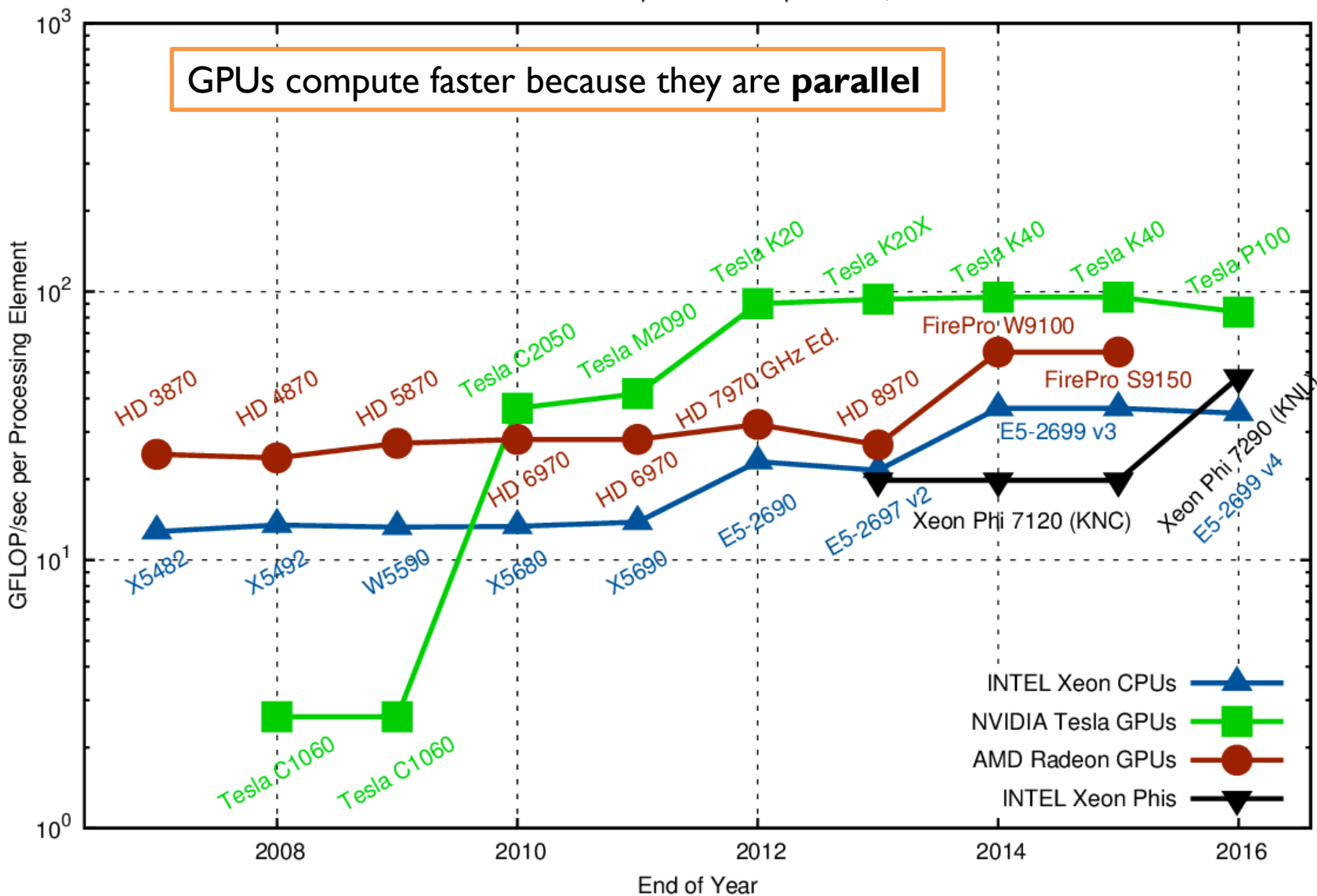
The term GPU was popularized by [Nvidia](#) in 1999, who marketed the [GeForce 256](#) as "the world's first ...Graphics Processing Unit." It was presented as a "single-chip processor with integrated [transform, lighting, triangle setup/clipping](#), and rendering engines".^[3]

Theoretical Peak Floating Point Operations per Clock Cycle, Double Precision



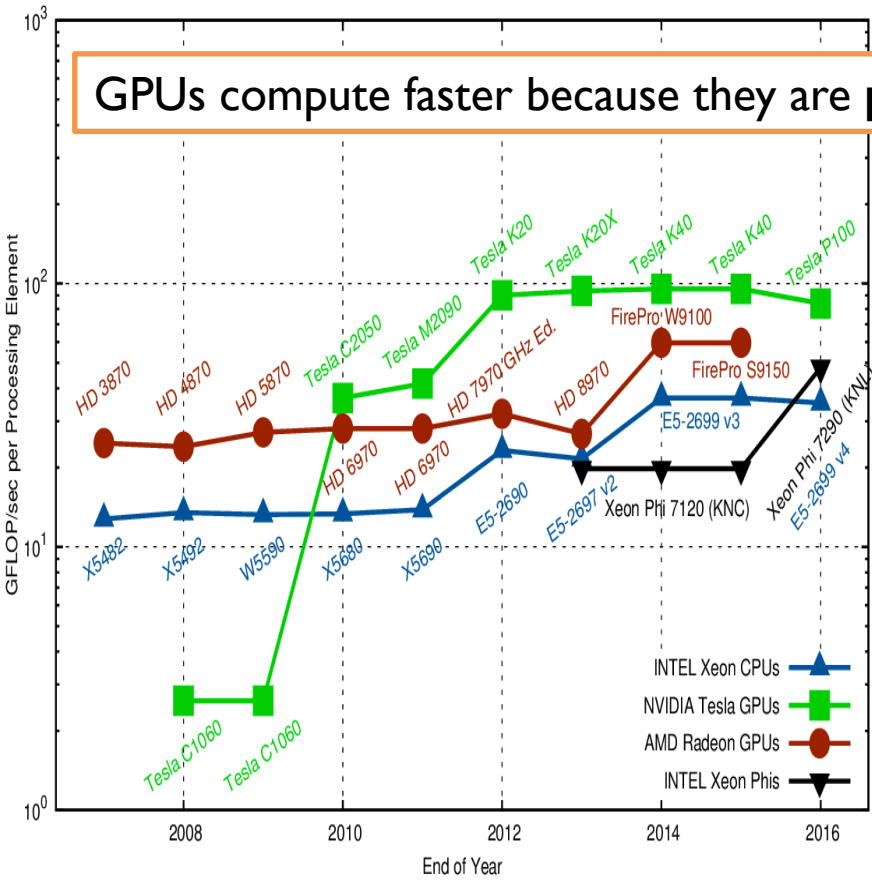
<https://www.karlrupp.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time/>

Theoretical Peak Performance per Core/Multiprocessor, Double Precision

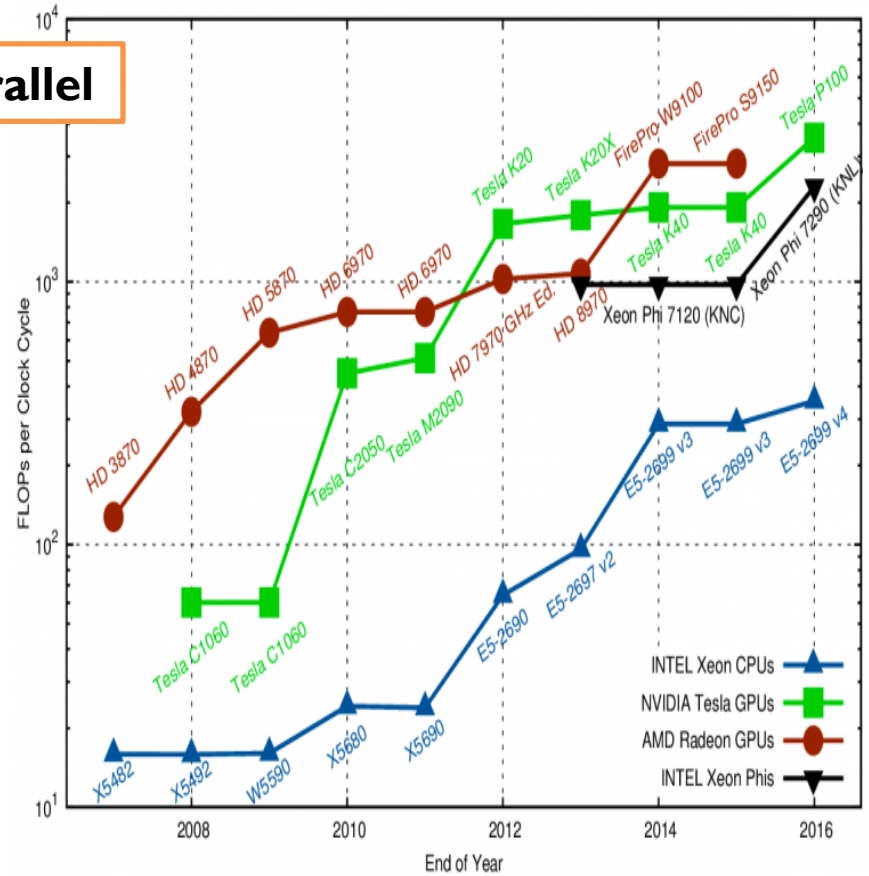


<https://www.karlrupp.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time/>

Theoretical Peak Performance per Core/Multiprocessor, Double Precision

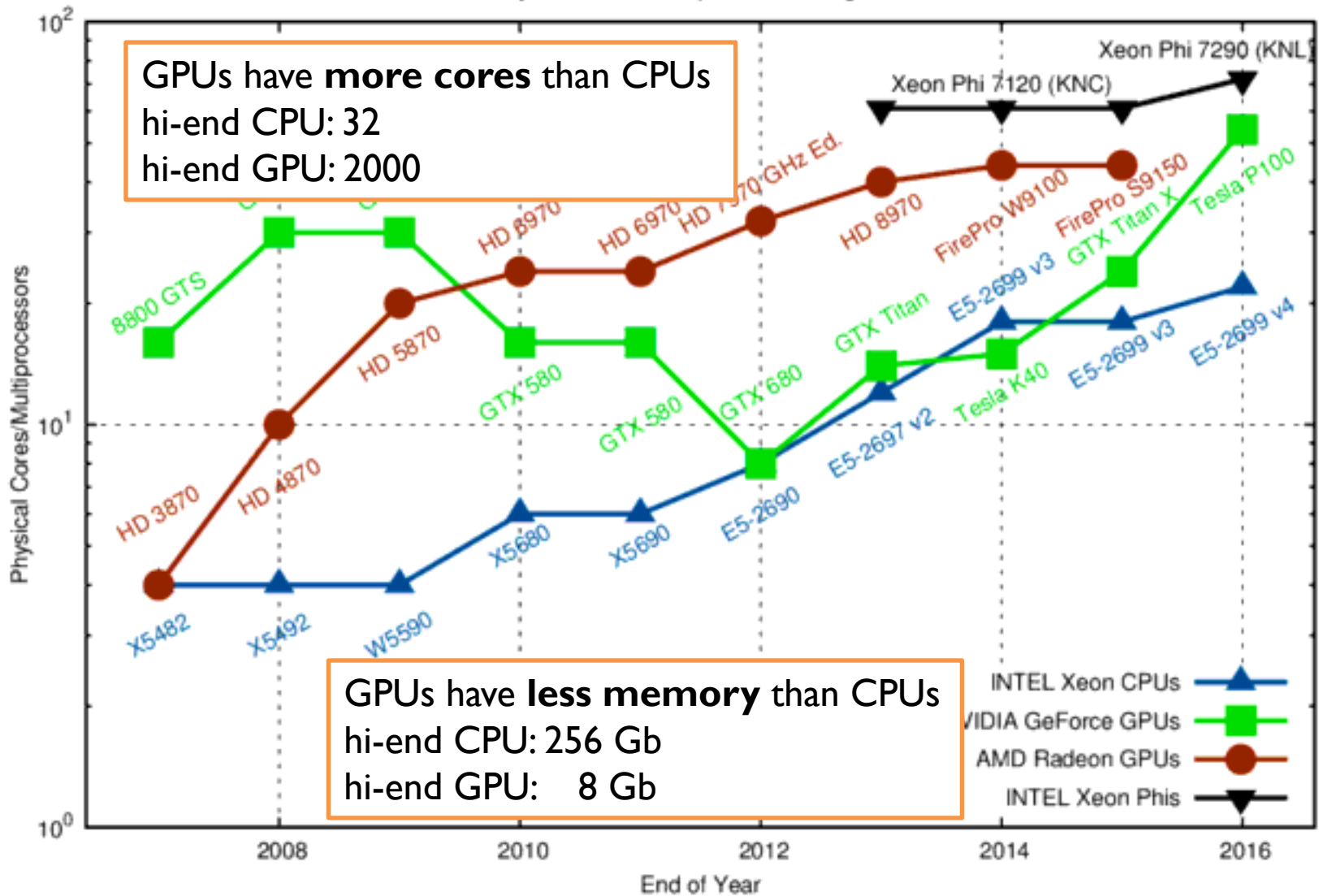


Theoretical Peak Floating Point Operations per Clock Cycle, Double Precision



<https://www.karlsruh.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time/>

Number of Physical Cores/Multiprocessors, High-End Hardware

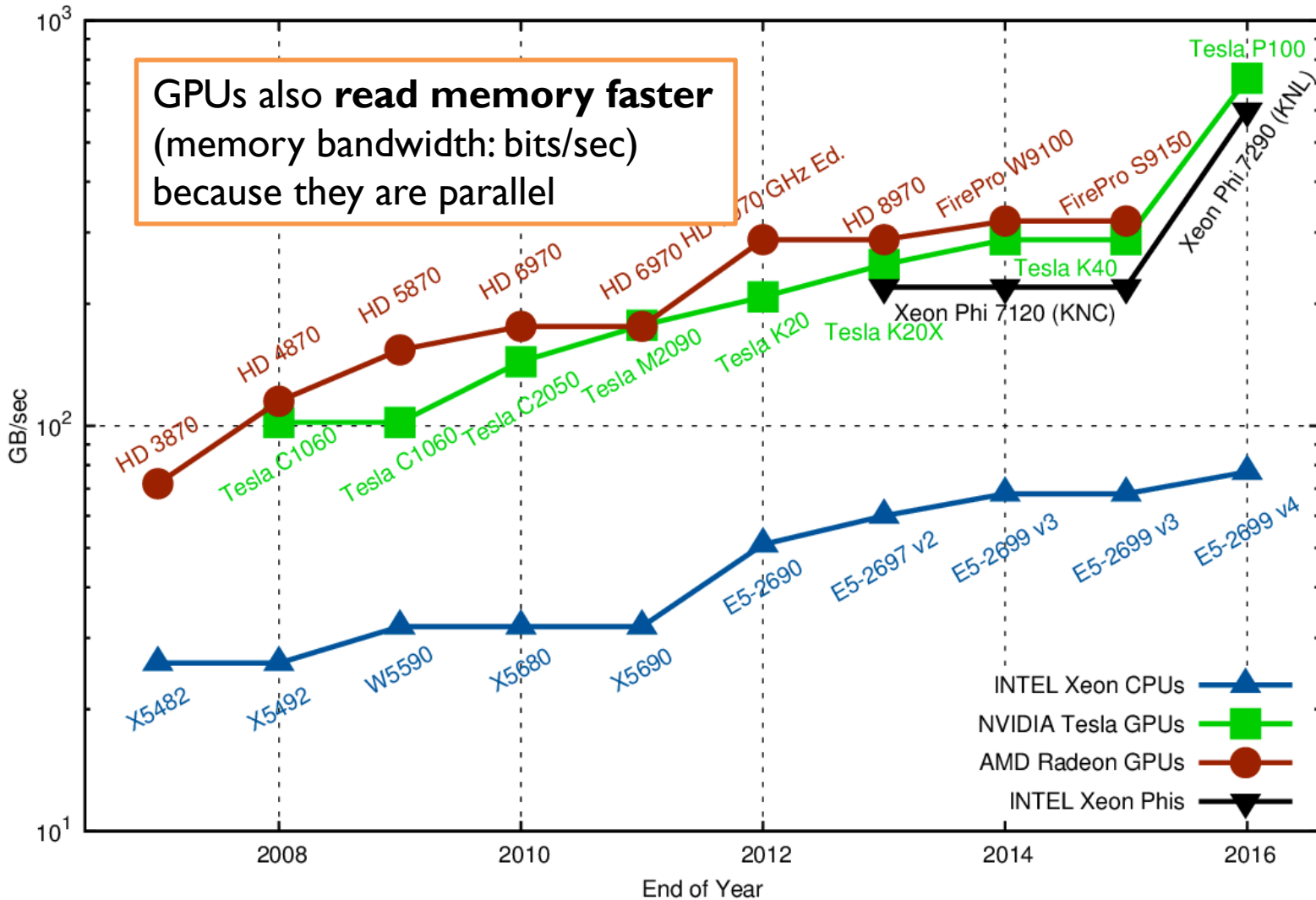


GPUs have **more cores** than CPUs
 hi-end CPU: 32
 hi-end GPU: 2000

GPUs have **less memory** than CPUs
 hi-end CPU: 256 Gb
 hi-end GPU: 8 Gb

<https://www.karlrupp.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time/>

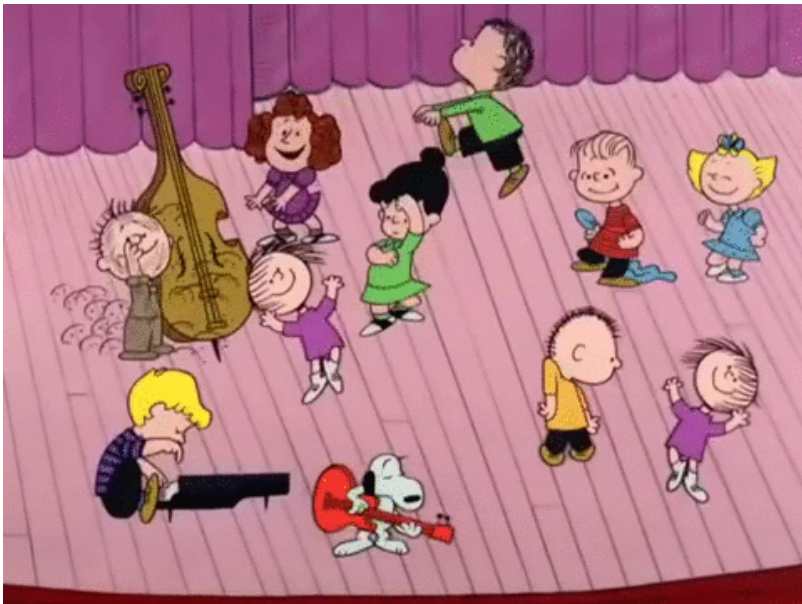
Theoretical Peak Memory Bandwidth Comparison

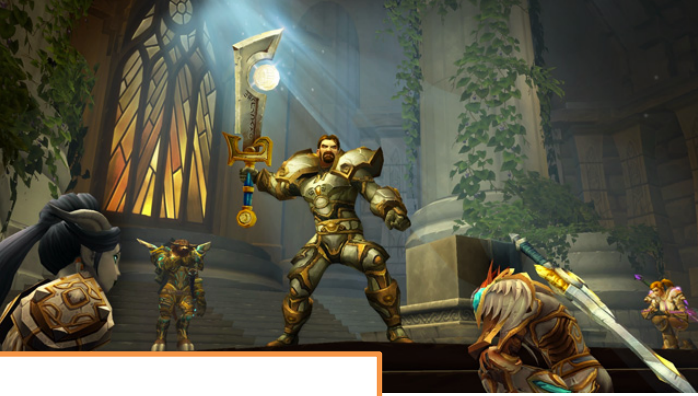


<https://www.karlrupp.net/2013/06/cpu-gpu-and-mic-hardware-characteristics-over-time/>

Summary of GPUs vs CPUs

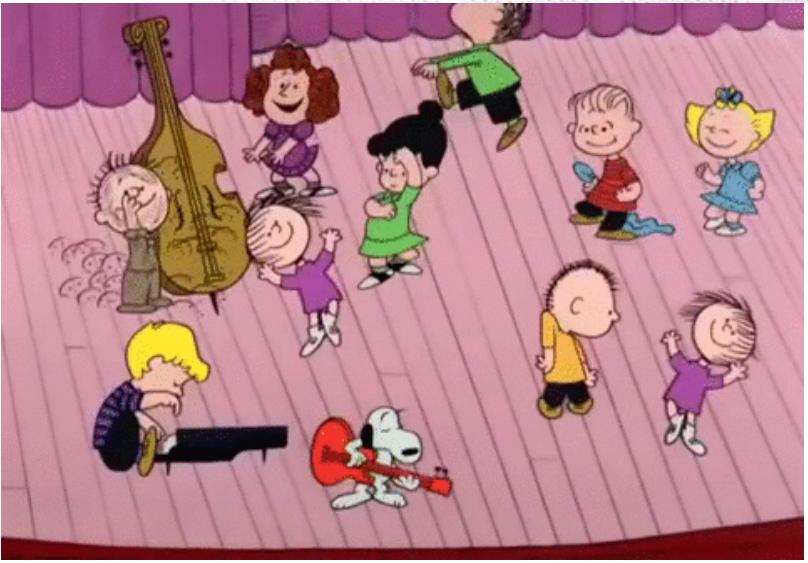
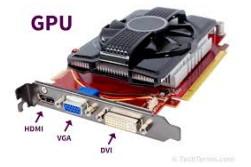
- less total memory
- more cores and more parallelism
- Multicore CPUs are mostly multiple-instruction multiple-data (MIMD)
- GPUs are mostly single-instruction multiple-data (SIMD)





Summary of GPUs vs CPU clusters

- way way less total memory, and memory is RAM not disk
- SIMD vs MIMD
- Expensive and reliable vs cheap and fault-tolerant



Art, Science and GPU's
Adam Savage & Jamie Hyneman
Explain Parallel Processing



▶ ▶ 🔊 0:03 / 1:33



<https://www.youtube.com/watch?v=-P28LKWTzrl>

HOW DO YOU USE A GPU?

http://developer.download.nvidia.com/compute/developertrainingmaterials/presentations/cuda_language/Introduction_to_CUDA_C.pptx

Using GPUs for ML

- Programming parallel machines is complicated
- To use the parallelism of a GPU in an ML algorithm we almost always use **matrix algebra** as an abstraction layer – i.e. *vectorize*

```
def logisticRegression(....):  
    ....  
    for X,Y in ....  
        Z = W.matrix_multiply(X)  
        P = logistic(Z)  
        W = W + learning_rate * (P - Y) * X  
        .....
```

Using GPUs for ML



```
def logisticRegression(....):  
    ....  
    for X,Y in ....  
        Z = W.matrix_multiply(X)  
        P = logistic(Z)  
        W = W + learning_rate * (P - Y) * X  
    .....
```

Using GPUs for ML

```
#include <iostream>
#include <algorithm>

using namespace std;

#define N 1024
#define RADIUS 3
#define BLOCK_SIZE 16

__global__ void stencil_1d(int *in, int *out) {
    __shared__ int temp[BLOCK_SIZE * 2 * RADIUS];
    int gindex = threadIdx.x + blockIdx.x * blockDim.x;
    int index = threadIdx.x + RADIUS;

    // Read input elements into shared memory
    temp[index] = in[gindex];
    if (threadIdx.x < RADIUS) {
        temp[index - RADIUS] = in[gindex - RADIUS];
        temp[index + BLOCK_SIZE] = in[gindex + BLOCK_SIZE];
    }

    // Synchronize (ensure all the data is available)
    __syncthreads();

    // Apply the stencil
    int result = 0;
    for (int offset = -RADIUS; offset <= RADIUS; offset++)
        result += temp[index + offset];

    // Store the result
    out[gindex] = result;
}

void fill_in(int *x, int n) {
    fill_n(x, n, 1);
}

int main(void) {
    int *in, *out; // host copies of a, b, c
    int *d_in, *d_out; // device copies of a, b, c
    int size = (N + 2*RADIUS) * sizeof(int);

    // Alloc space for host copies and setup values
    in = (int *)malloc(size); fill_in(in, N + 2*RADIUS);
    out = (int *)malloc(size); fill_in(out, N + 2*RADIUS);

    // Alloc space for device copies
    cudaMalloc((void **)&d_in, size);
    cudaMalloc((void **)&d_out, size);

    // Copy to device
    cudaMemcpy(d_in, in, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_out, out, size, cudaMemcpyHostToDevice);

    // Launch stencil_1d() kernel on GPU
    stencil_1d<<<(N/BLOCK_SIZE, BLOCK_SIZE)>>>(d_in + RADIUS,
    d_out + RADIUS);

    // Copy result back to host
    cudaMemcpy(out, d_out, size, cudaMemcpyDeviceToHost);

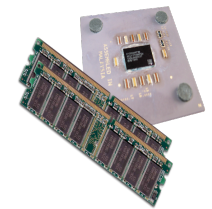
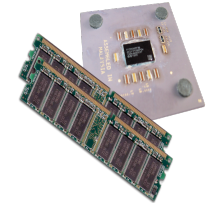
    // Cleanup
    free(in); free(out);
    cudaFree(d_in); cudaFree(d_out);
    return 0;
}
```

parallel fn

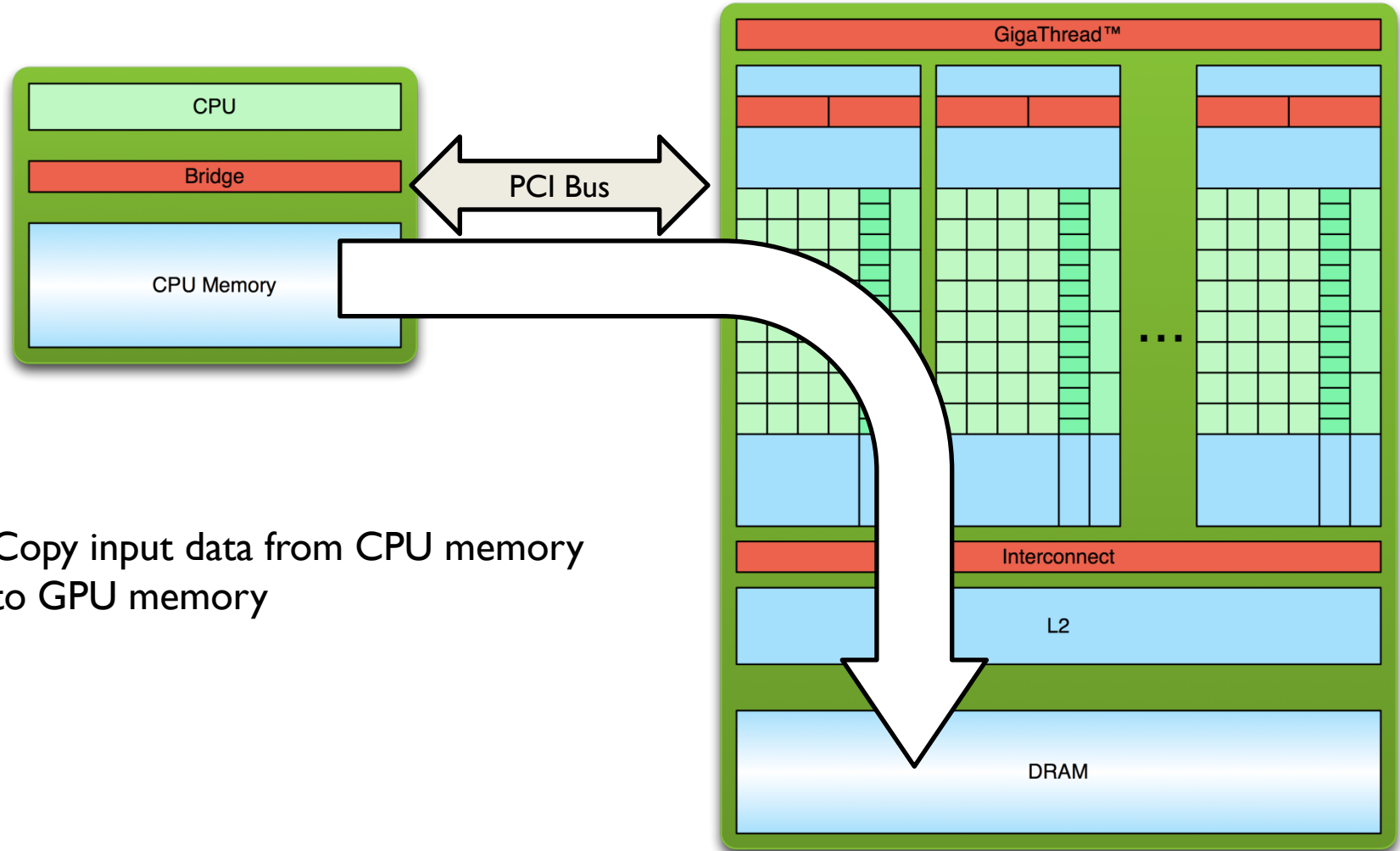
serial code

parallel code

serial code

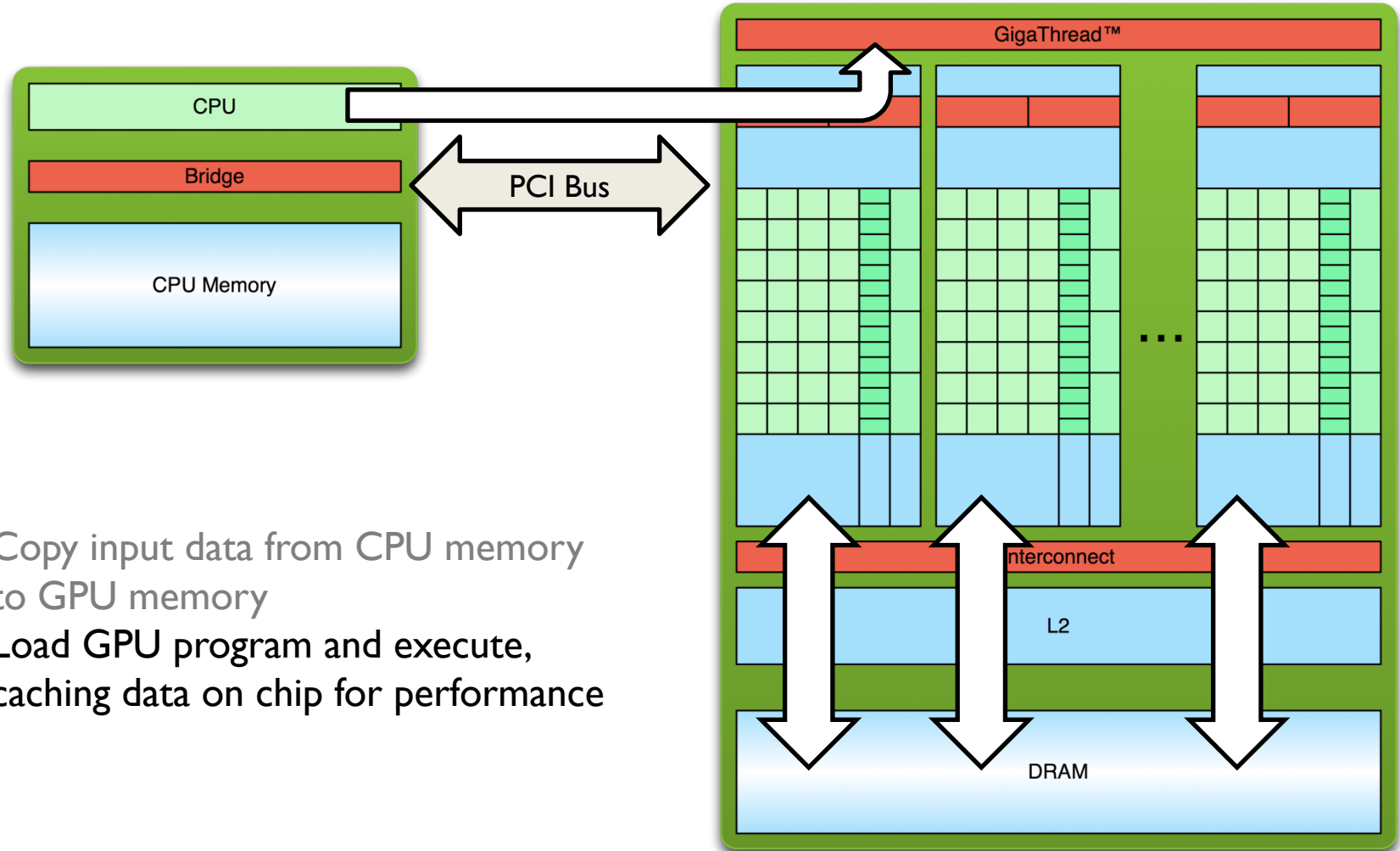


Simple Processing Flow



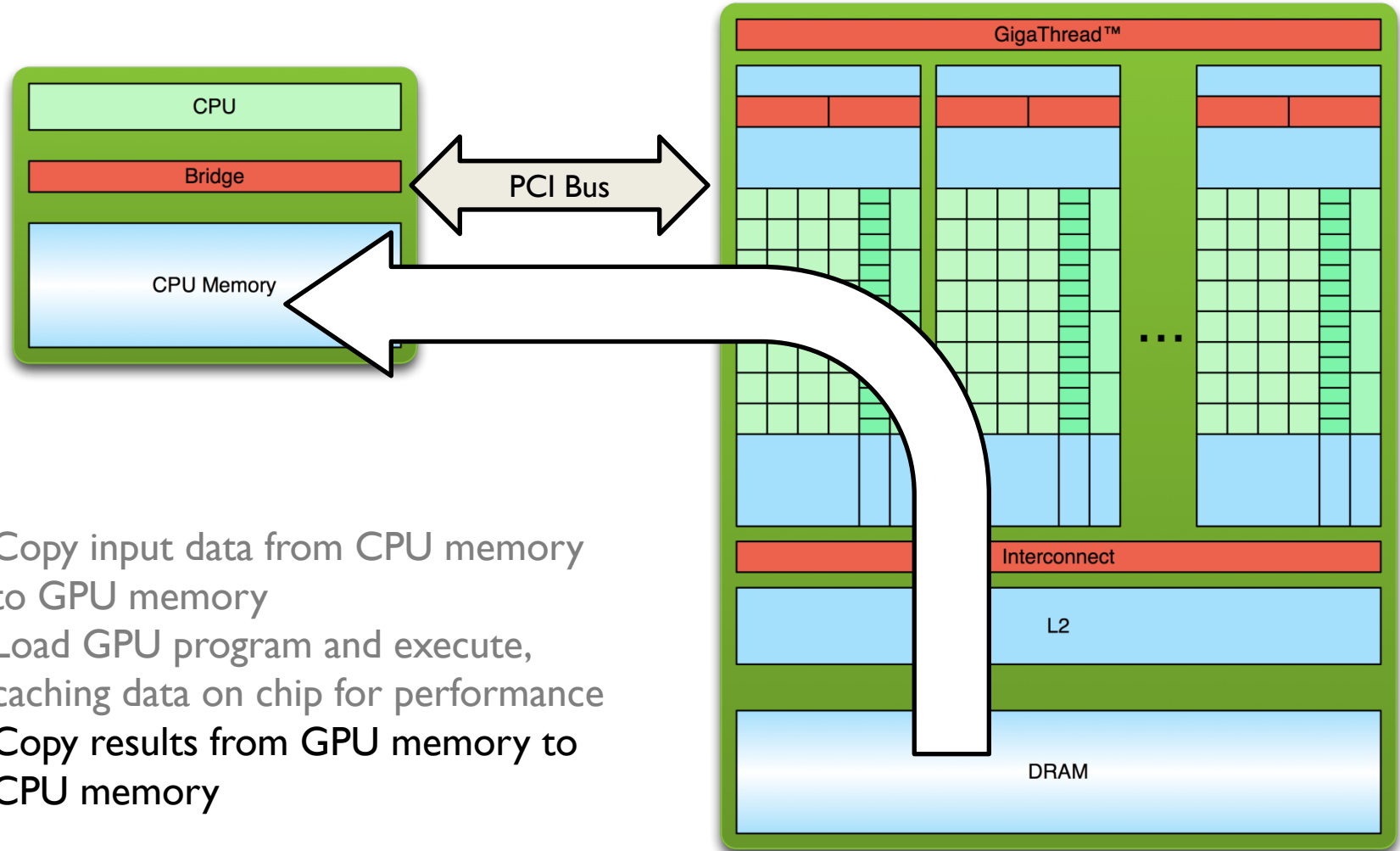
- I. Copy input data from CPU memory to GPU memory

Simple Processing Flow



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

Simple Processing Flow



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

SOME EXAMPLE CODE

Hello World!

```
int main(void) {  
    printf("Hello World!\n");  
    return 0;  
}
```

- Standard C that runs on the host (the CPU)
- NVIDIA compiler (nvcc)
- We can also write code for the device (the GPU)

Output:

```
$ nvcc  
hello_world.  
cu  
$ a.out  
Hello World!  
$
```


Hello World! with Device Code

```
__global__ void mykernel(void) {  
}
```

```
int main(void) {  
    mykernel<<<1,1>>>();  
    printf("Hello World!\n");  
    return 0;  
}
```

- Two new things here...

Hello World! with Device Code

```
__global__ void mykernel(void) {  
}
```

- CUDA C/C++ keyword `__global__` indicates a function that:
 - Runs on the device
 - Is called from host code
- `nvcc` separates source code into host and device components
 - Device functions (e.g. `mykernel()`) processed by NVIDIA compiler
 - Host functions (e.g. `main()`) processed by standard host compiler
 - `gcc, cl.exe`

Hello World! with Device COde

```
mykernel<<<1,1>>>();
```

- Triple angle brackets mark a call from *host* code to *device* code
 - Also called a “kernel launch”
 - We’ll get to the parameters (1,1) soon

Hello World! with Device Code

```
__global__ void mykernel(void) {  
}
```

```
int main(void) {  
    mykernel<<<1,1>>>();  
    printf("Hello World!\n");  
    return 0;  
}
```

- `mykernel()` does nothing at all in this example ... so let's fix that.

Output:

```
$ nvcc  
hello.cu  
$ a.out  
Hello World!  
$
```

Addition on the Device

- A simple kernel to add two integers (coming up: adding two arrays)

```
__global__ void add(int *a, int *b, int *c) {  
    *c = *a + *b;  
}
```

- As before `__global__` is a CUDA C/C++ keyword meaning
 - `add()` will execute on the device
 - `add()` will be called from the host

Addition on the Device

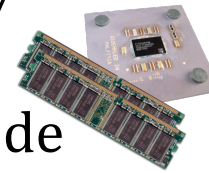
- Note that we use pointers for the variables

```
__global__ void add(int *a, int *b, int *c) {  
    *c = *a + *b;  
}
```

- `add()` runs on the device, so `a`, `b` and `c` must point to device memory
- We need to allocate memory on the GPU

Memory Management

- Host and device memory are separate entities
 - *Device* pointers point to GPU memory
 - May be passed to/from host code
 - May *not* be dereferenced in host code
 - *Host* pointers point to CPU memory
 - May be passed to/from device code
 - May *not* be dereferenced in device code
- Simple CUDA API for handling device memory
 - `cudaMalloc()`, `cudaFree()`, `cudaMemcpy()`
 - Similar to the C equivalents `malloc()`, `free()`, `memcpy()`



Addition on the Device: `add()`

- Returning to our `add()` kernel

```
__global__ void add(int *a, int *b, int *c) {  
    *c = *a + *b;  
}
```

- Let's take a look at `main()`...

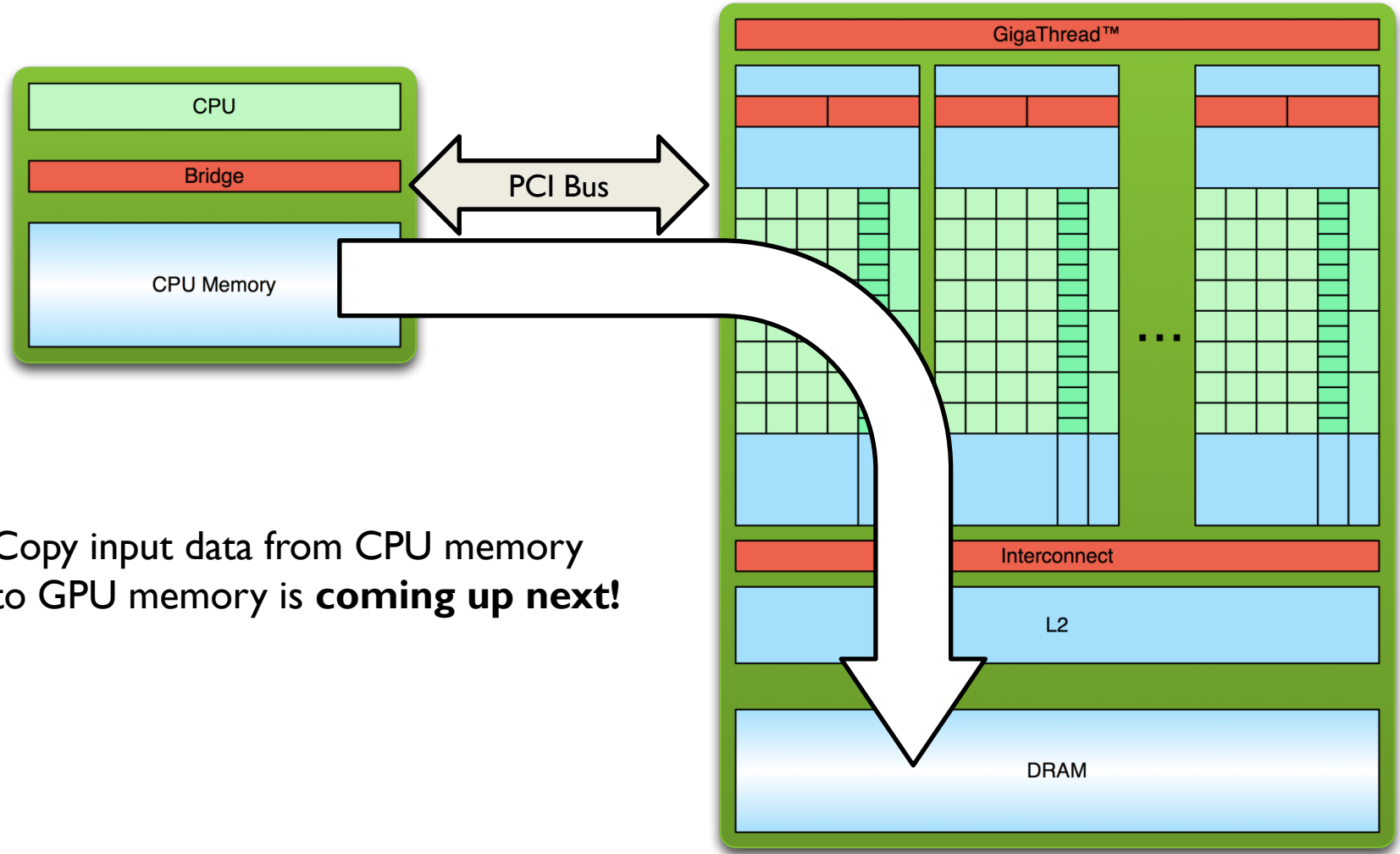
Addition on the Device: `main()`

```
int main(void) {
    int a, b, c;           // host copies of a, b, c
    int *d_a, *d_b, *d_c; // device copies of a, b, c
    int size = sizeof(int);

    // Allocate space for device copies of a, b, c
    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    // Setup input values
    a = 2;
    b = 7;
```

We're getting ready to do this...



- I. Copy input data from CPU memory to GPU memory is **coming up next!**

Addition on the Device: `main()`

```
// Copy inputs to device
cudaMemcpy(d_a, &a, size, cudaMemcpyHostToDevice);
cudaMemcpy(d_b, &b, size, cudaMemcpyHostToDevice);

// Launch add() kernel on GPU
add<<<1,1>>>(d_a, d_b, d_c);

// Copy result back to host
cudaMemcpy(&c, d_c, size, cudaMemcpyDeviceToHost);

// Cleanup
cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
return 0;
}
```

Next: Vector Addition on the Device

- With `add()` running in parallel we can do vector addition
- Terminology: each parallel invocation of `add()` is referred to as a **block**
 - The set of blocks is referred to as a **grid**
 - Each invocation can refer to its block index using `blockIdx.x`

```
__global__ void add(int *a, int *b, int *c) {  
    c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];  
}
```

- By using `blockIdx.x` to index into the array, each block handles a different index

Vector Addition on the Device

```
__global__ void add(int *a, int *b, int *c) {  
    c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];  
}
```

- On the device, each block can execute in parallel:

Block 0

```
c[0] = a[0] + b[0];
```

Block 1

```
c[1] = a[1] + b[1];
```

Block 2

```
c[2] = a[2] + b[2];
```

Block 3

```
c[3] = a[3] + b[3];
```

Vector Addition on the Device: `add()`

- Returning to our parallelized `add()` kernel

```
__global__ void add(int *a, int *b, int *c) {  
    c[blockIdx.x] = a[blockIdx.x] + b[blockIdx.x];  
}
```

- Let's take a look at `main()`...

Vector Addition on the Device: `main()`

```
#define N 512
int main(void) {
    int *a *b *c // host copies of a, b, c
    int *d_a, *d_b, *d_c; // device copies of a, b, c
    int size = N * sizeof(int);

    // Alloc space for device copies of a, b, c
    cudaMalloc((void **)&d_a, size);
    cudaMalloc((void **)&d_b, size);
    cudaMalloc((void **)&d_c, size);

    // Alloc space for host copies of a, b, c and setup input values
    a = (int *)malloc(size); random_ints(a, N);
    b = (int *)malloc(size); random_ints(b, N);
    c = (int *)malloc(size);
```

Vector Addition on the Device: `main()`

```
// Copy inputs to device
```

```
cudaMemcpy(d_a, a, size, cudaMemcpyHostToDevice);
```

```
cudaMemcpy(d_b, b, size, cudaMemcpyHostToDevice);
```

```
// Launch add() kernel on GPU with N blocks
```

```
add<<<N,1>>>(d_a, d_b, d_c);
```

a little more magic...

```
// Copy result back to host
```

```
cudaMemcpy(c, d_c, size, cudaMemcpyDeviceToHost);
```

```
// Cleanup
```

```
free(a); free(b); free(c);
```

```
cudaFree(d_a); cudaFree(d_b); cudaFree(d_c);
```

```
return 0;
```

```
}
```


Coordinating Host & Device

- Kernel launches are **asynchronous**
 - Control returns to the CPU immediately
- CPU needs to synchronize before consuming the results

`cudaMemcpy ()`

Blocks the CPU until the copy is complete
Copy begins when all preceding CUDA calls have completed

`cudaMemcpyAsync ()`

Asynchronous, does not block the CPU

`cudaDeviceSynchronize ()`

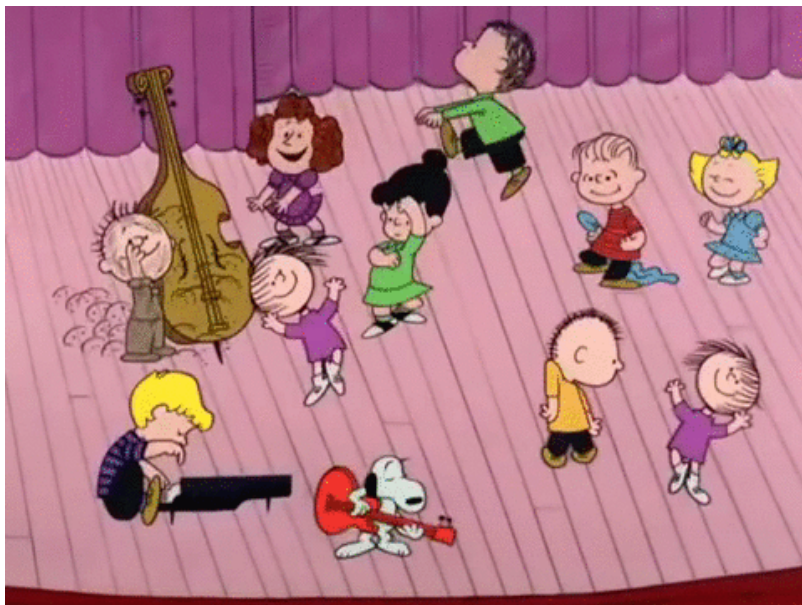
Blocks the CPU until all preceding CUDA calls have completed

MORE DETAILS ON GPU PROGRAMMING

http://developer.download.nvidia.com/compute/developertrainingmaterials/presentations/cuda_language/Introduction_to_CUDA_C.pptx

Summary of GPUs vs CPUs

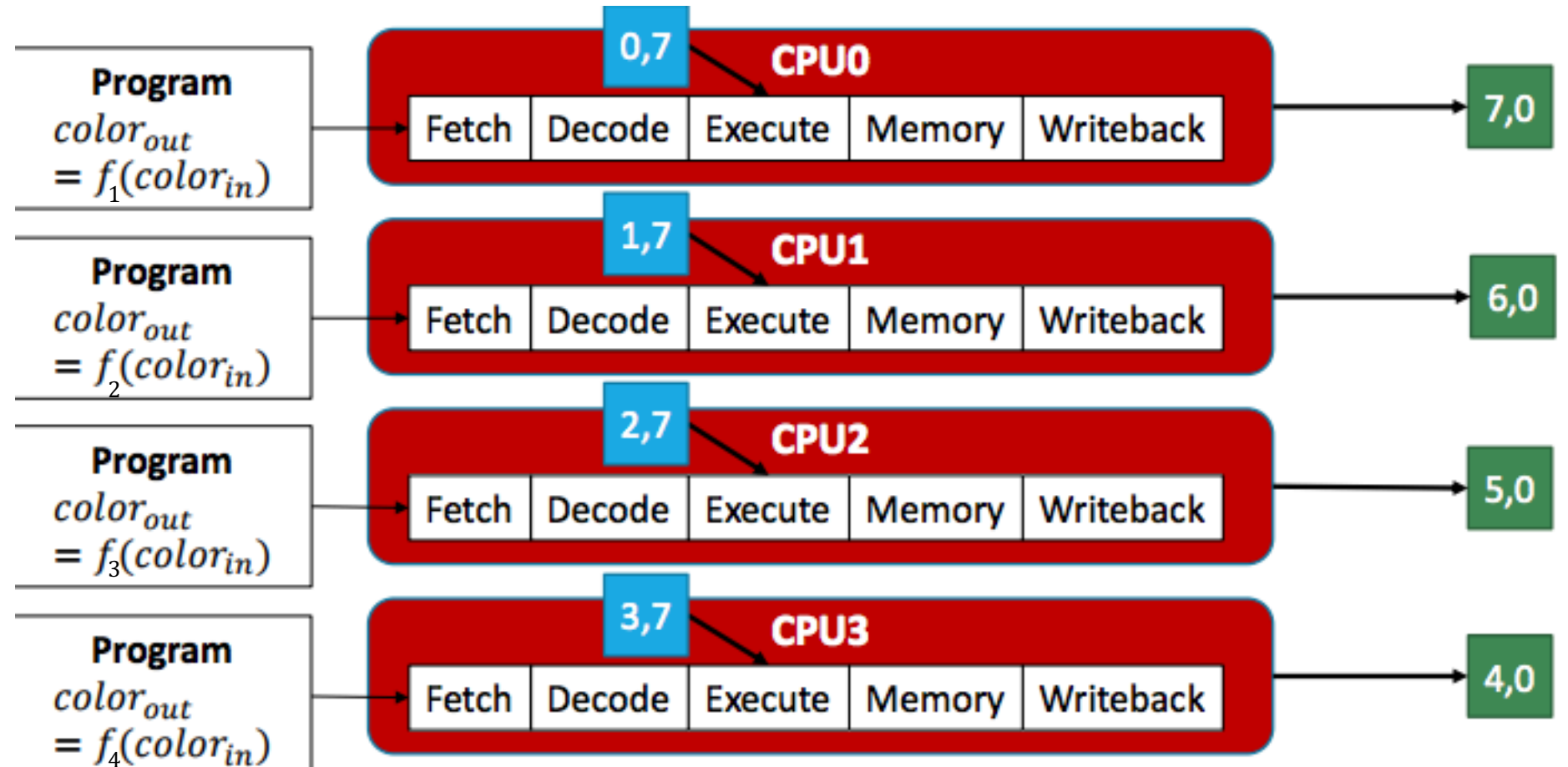
- less total memory
- more cores and more parallelism
- Multicore CPUs are *mostly* multiple-instruction multiple-data (MIMD)
- GPUs are *mostly* single-instruction multiple-data (SIMD)



Blocks, Grids, Threads, Warps

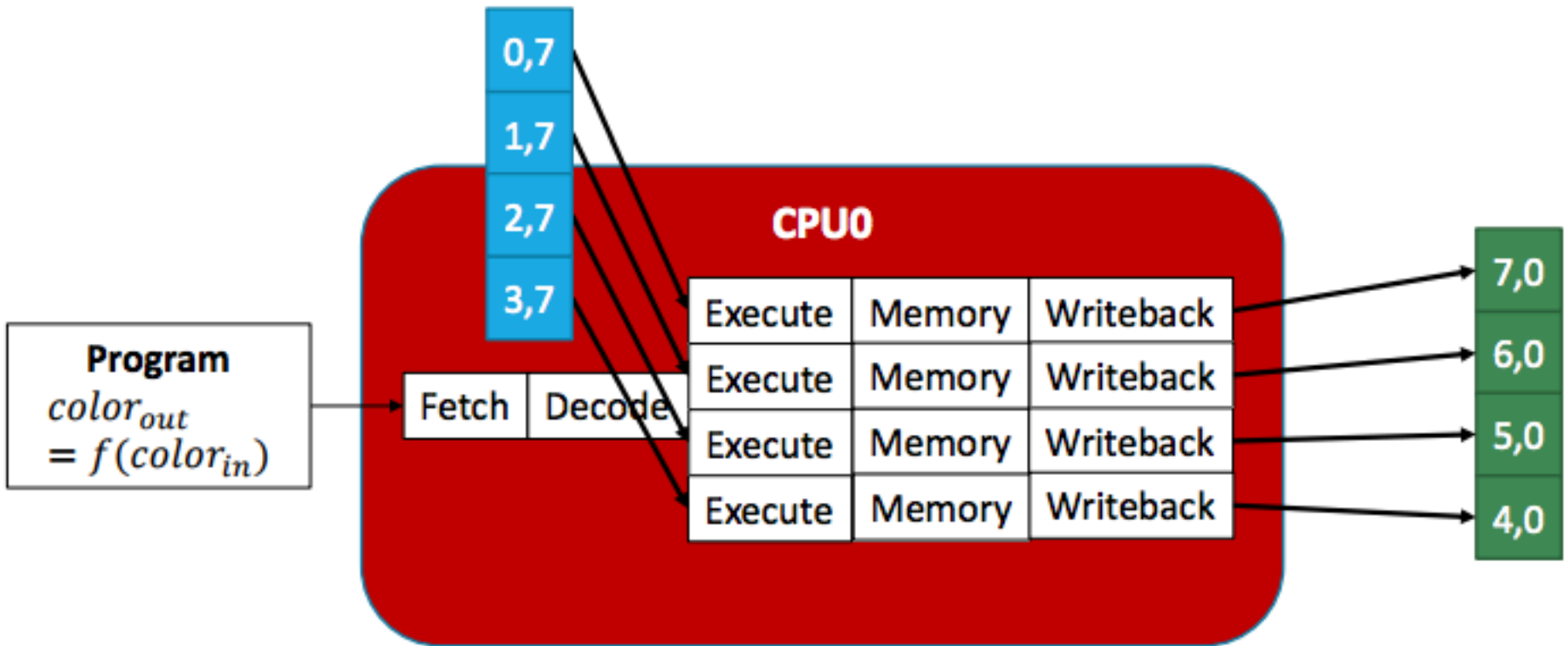
- Recall blocks are the things that work in parallel, and blocks are arranged in **grids**
- $c[\text{blockIdx}.x] = a[\text{blockIdx}.x] + b[\text{blockIdx}.x];$
- That would be SIMD (single instruction multiple data)
- It's actually more complicated than that....

MIMD

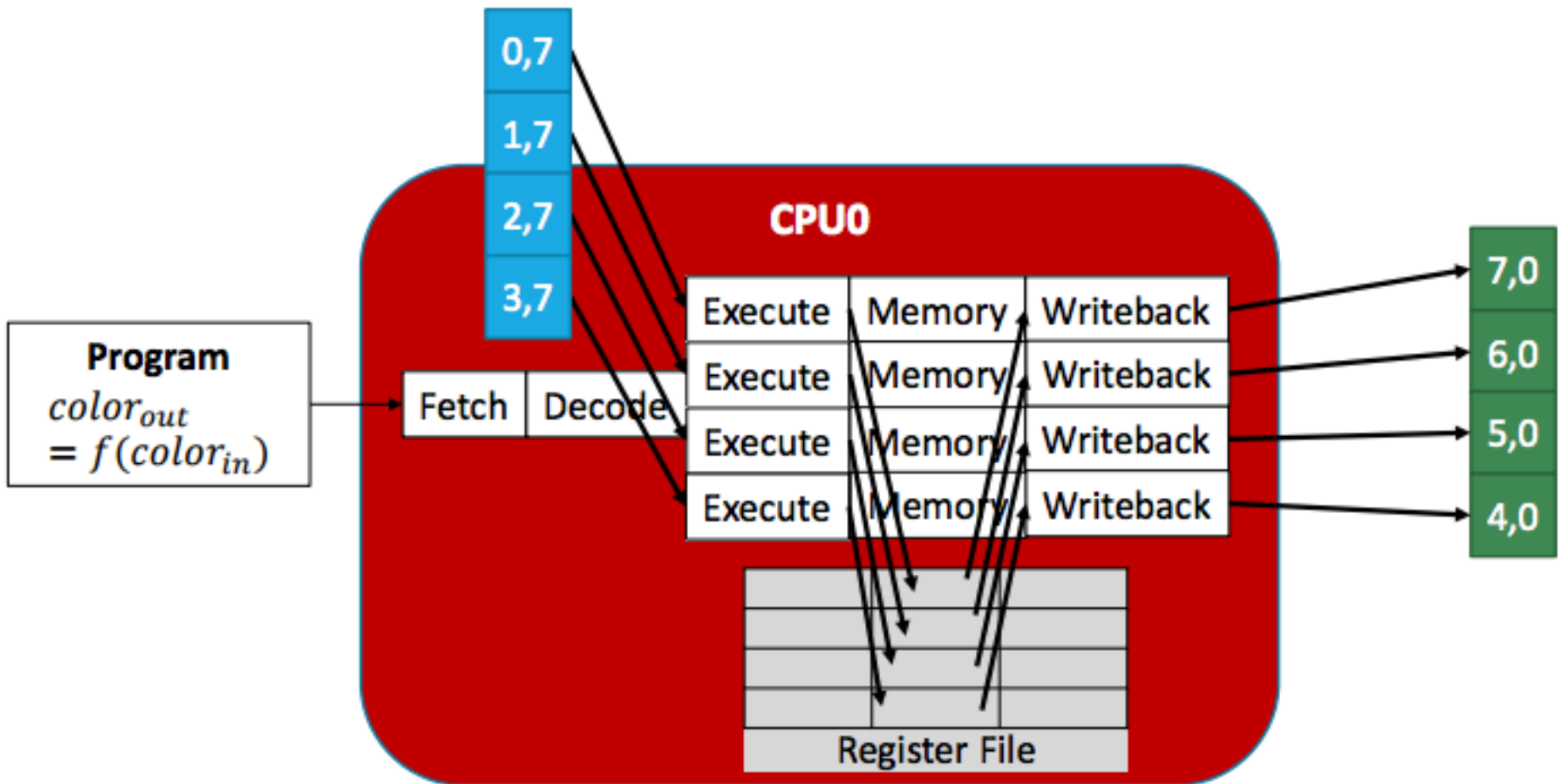


<https://courses.cs.washington.edu/courses/cse471/13sp/lectures/GPUsStudents.pdf>

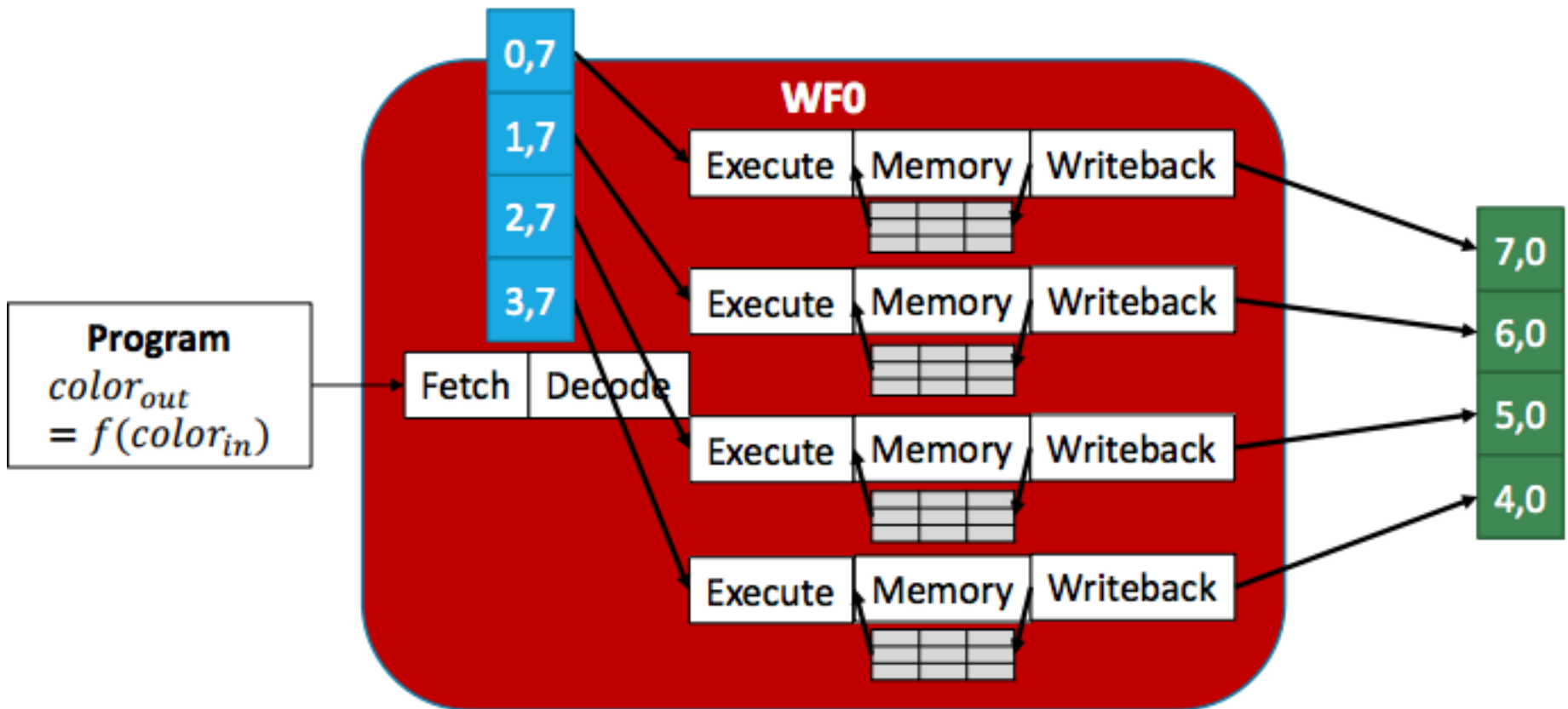
SIMD



SIMD: Zooming in

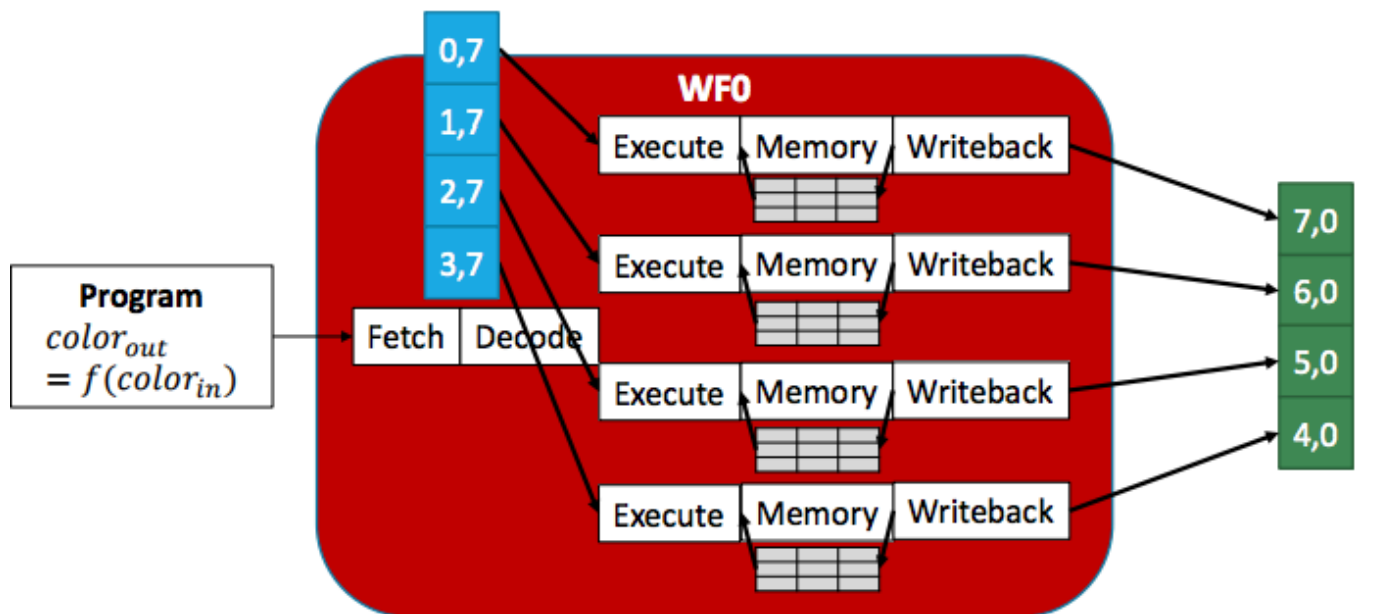


SIMT: Single Instruction Multiple Threads



SIMT: Single Instruction Multiple Threads

A **thread** can access its own **block id** and also **thread id**.
Blocks and threads are in a **grid**, which is 2D or 3D (there's a .x and a .y part)



Comparison

MIMD/SPMD



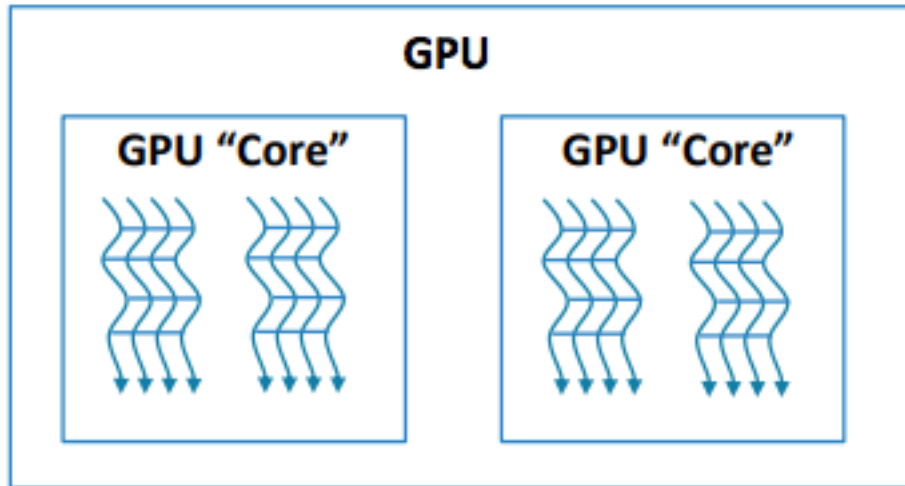
SIMD/Vector



SIMT



What's in a GPU?



Threads (SIMT, synchronous threads) are grouped into **cores** (which are decoupled, like a MIMD machine)

MIMD/SPMD



SIMD/Vector



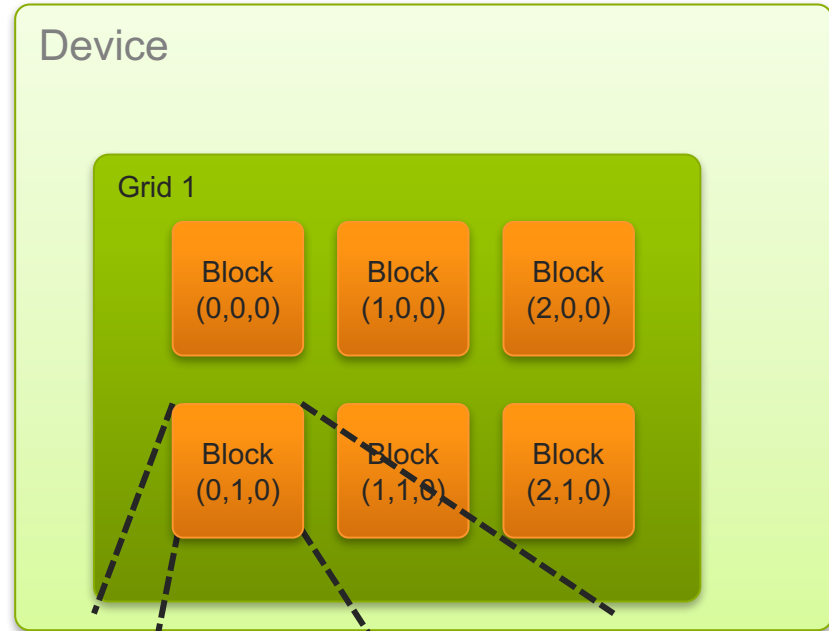
SIMT



IDs and Dimensions

– A kernel is launched as a grid of blocks of threads

- `blockIdx` and `threadIdx` are 3D
- We showed only one dimension (x)



Block (1,1,0)

Thread (0,0,0)	Thread (1,0,0)	Thread (2,0,0)	Thread (3,0,0)	Thread (4,0,0)
Thread (0,1,0)	Thread (1,1,0)	Thread (2,1,0)	Thread (3,1,0)	Thread (4,1,0)
Thread (0,2,0)	Thread (1,2,0)	Thread (2,2,0)	Thread (3,2,0)	Thread (4,2,0)

• Built-in variables:

- `threadIdx`
- `blockIdx`
- `blockDim`
- `gridDim`

Thread and block parallelism

```
tx = cuda.threadIdx.x
ty = cuda.threadIdx.y
bx = cuda.blockIdx.x
by = cuda.blockIdx.y
bw = cuda.blockDim.x
bh = cuda.blockDim.y
x = tx + bx * bw
y = ty + by * bh
array[x, y] = something(x, y)
```

How Do You REALLY Use a GPU?

Recap: Using GPUs for ML

- Programming parallel machines is complicated
- To use the parallelism of a GPU in an ML algorithm we almost always use **matrix algebra** as an abstraction layer – i.e. *vectorize*

```
def logisticRegression(....):  
    ....  
    for X,Y in ....  
        Z = W.matrix_multiply(X)  
        P = logistic(Z)  
        W = W + learning_rate * (P - Y) * X  
        .....
```

Recap: Using GPUs for ML



```
def logisticRegression(....):  
    ....  
    for X,Y in ....  
        Z = W.matrix_multiply(X)  
        P = logistic(Z)  
        W = W + learning_rate * (P - Y) * X  
    .....
```


Recap: Using GPUs for ML

```
#include <iostream>
#include <algorithm>

using namespace std;

#define N 1024
#define RADIUS 3
#define BLOCK_SIZE 16

__global__ void stencil_1d(int *in, int *out) {
    __shared__ int temp[BLOCK_SIZE * 2 * RADIUS];
    int gidindex = threadIdx.x + blockIdx.x * blockDim.x;
    int index = threadIdx.x + RADIUS;

    // Read input elements into shared memory
    temp[gidindex] = in[gidindex];
    if (threadIdx.x < RADIUS) {
        temp[index - RADIUS] = in[index - RADIUS];
        temp[index + BLOCK_SIZE] = in[index + BLOCK_SIZE];
    }

    // Synchronize (ensure all the data is available)
    __syncthreads();

    // Apply the stencil
    int result = 0;
    for (int offset = -RADIUS; offset <= RADIUS; offset++)
        result += temp[index + offset];

    // Store the result
    out[gidindex] = result;
}

void fill_ints(int *x, int n) {
    fill_n(x, n, 1);
}

int main(void) {
    int *in, *out; // host copies of a, b, c
    int *d_in, *d_out; // device copies of a, b, c
    int size = (N + 2*RADIUS) * sizeof(int);

    // Alloc space for host copies and setup values
    in = (int *)malloc(size); fill_ints(in, N + 2*RADIUS);
    out = (int *)malloc(size); fill_ints(out, N + 2*RADIUS);

    // Alloc space for device copies
    cudaMalloc((void **)&d_in, size);
    cudaMalloc((void **)&d_out, size);

    // Copy to device
    cudaMemcpy(d_in, in, size, cudaMemcpyHostToDevice);
    cudaMemcpy(d_out, out, size, cudaMemcpyHostToDevice);

    // Launch stencil_1d() kernel on GPU
    stencil_1d<<<N/BLOCK_SIZE, BLOCK_SIZE>>>(d_in + RADIUS,
    d_out + RADIUS);

    // Copy result back to host
    cudaMemcpy(out, d_out, size, cudaMemcpyDeviceToHost);

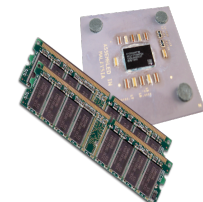
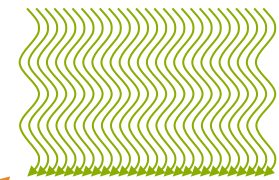
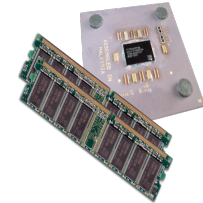
    // Cleanup
    free(in); free(out);
    cudaFree(d_in); cudaFree(d_out);
    return 0;
}
```

parallel fn

serial code

parallel code

serial code



How Do You REALLY Use a GPU?

**Example: Vectorizing Logistic Regression
with Stochastic Gradient Descent**

Vectorizing ML Algorithms

- We want to specify an algorithm so that
 - It's concise
 - It's easy to verify correctness
 - It reveals to a compiler what can be done in parallel on a GPU
 - Matrix-vector computations!

Vectorizing logistic regression

- Conditional probability of an example:

$$P(Y = y|X = \mathbf{x}, \mathbf{w}) = \begin{cases} \frac{1}{1+e^{-\mathbf{x}\cdot\mathbf{w}}} & \text{if } y = 1 \\ 1 - \frac{1}{1+e^{-\mathbf{x}\cdot\mathbf{w}}} & \text{if } y = 0 \end{cases}$$

- Conditional log likelihood of an example:

$$\log P(Y = y|X = \mathbf{x}, \mathbf{w}) = \begin{cases} \log p & \text{if } y = 1 \\ \log(1 - p) & \text{if } y = 0 \end{cases}$$

$$p \equiv \frac{1}{1 + e^{-\mathbf{x}\cdot\mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

$$\frac{\partial}{\partial w^j} \log P(Y = y|X = \mathbf{x}, \mathbf{w}) = (y - p)x^j$$

Vectorizing logistic regression

- Computation we'd like to parallelize:
 - For each \mathbf{x} in the minibatch, compute

$$p \equiv \frac{1}{1 + e^{-\mathbf{x} \cdot \mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

- For each feature j : update w^j using

$$\frac{\partial}{\partial w^j} \log P(Y = y | X = \mathbf{x}, \mathbf{w}) = (y - p)x^j$$

Vectorizing logistic regression

- Computation we'd like to parallelize:
 - For each \mathbf{x} in the minibatch X_{batch} , compute

$$p \equiv \frac{1}{1 + e^{-\mathbf{x} \cdot \mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

Note : $\mathbf{x} \cdot \mathbf{w}$ can be partially computed in parallel...

Vectorizing logistic regression

- Computation we'd like to parallelize:
 - For each \mathbf{x} in the minibatch X_{batch} , compute

$$p \equiv \frac{1}{1 + e^{-\mathbf{x} \cdot \mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

Note: $X_{batch} \cdot \mathbf{w}$
can be
computed in
parallel...

$$X_{batch} = \begin{bmatrix} x_1^1 & \cdots & x_1^J \\ \vdots & \ddots & \vdots \\ x_B^1 & \cdots & x_B^J \end{bmatrix}$$

Vectorizing logistic regression

- Computation we'd like to parallelize:
 - For each \mathbf{x} in the minibatch X_{batch} , compute

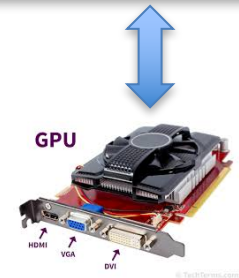
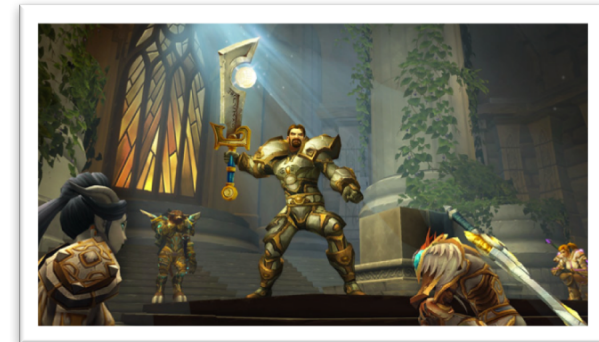
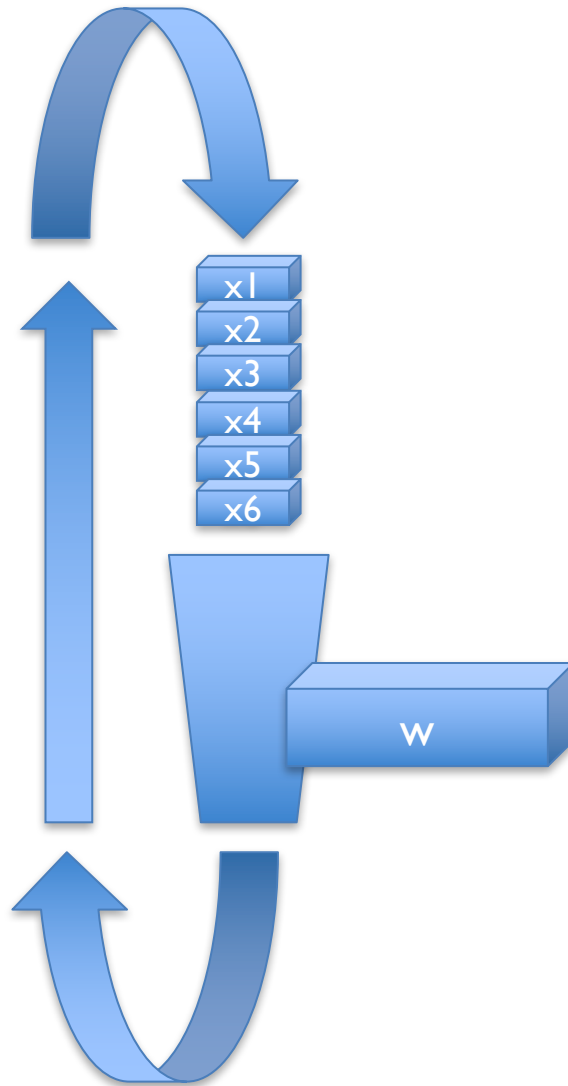
$$p \equiv \frac{1}{1 + e^{-\mathbf{x} \cdot \mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

$$X_{batch} \mathbf{w} = \begin{bmatrix} x_1^1 & \cdots & x_1^J \\ \vdots & \ddots & \vdots \\ x_B^1 & \cdots & x_B^J \end{bmatrix} \begin{bmatrix} w^1 \\ \vdots \\ w^J \end{bmatrix} = \begin{bmatrix} \mathbf{w} \cdot \mathbf{x}_1 \\ \vdots \\ \mathbf{w} \cdot \mathbf{x}_B \end{bmatrix}$$

We're most of the way now...

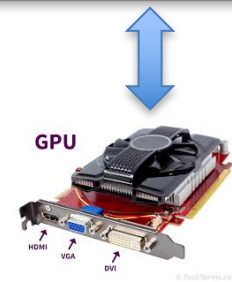
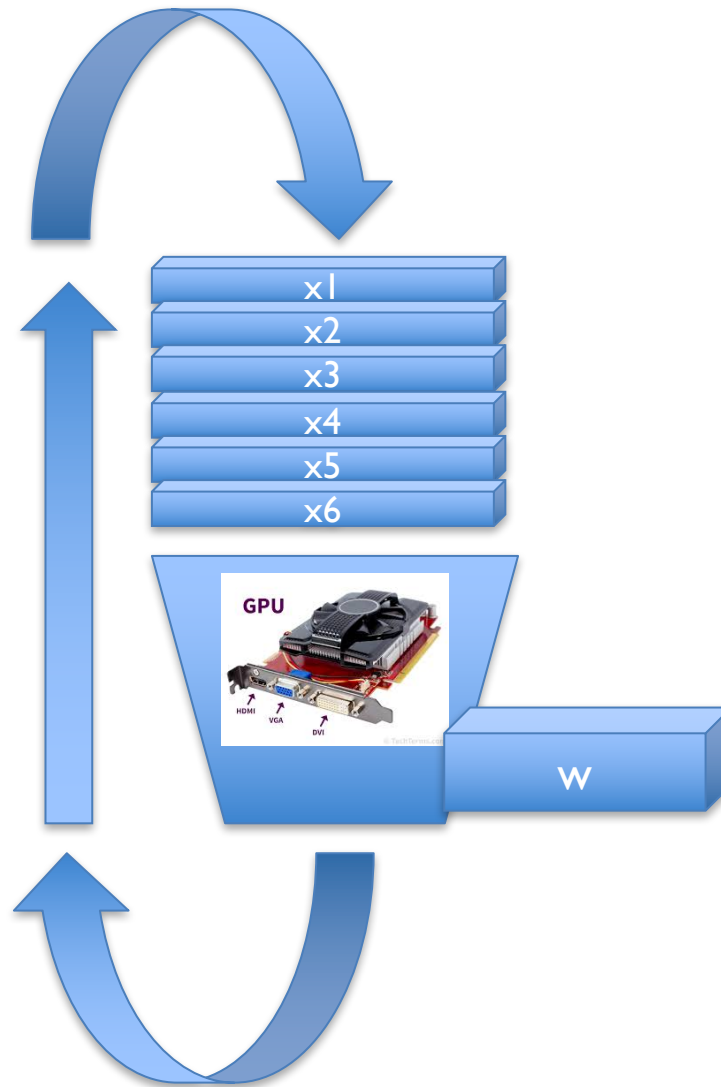
Streaming SGD:

- Iterative
- Sequential
- Fast
- Scale up by bounding memory
- You can handle very large datasets ... but slowly



Streaming SGD:

- Iterative
- Sequential
- Fast
- Scale up by bounding memory
- You can handle very large datasets ... but slowly
- You can speed it up by making the tasks in the stream bigger and doing them in **parallel**
- A GPU is a good way of doing that



Vectorizing logistic regression

- Computation we'd like to parallelize:
 - For each \mathbf{x} in the minibatch X_{batch} , compute

$$p \equiv \frac{1}{1 + e^{-\mathbf{x} \cdot \mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

$$X_{batch} \mathbf{w} = \begin{bmatrix} x_1^1 & \cdots & x_1^J \\ \vdots & \ddots & \vdots \\ x_B^1 & \cdots & x_B^J \end{bmatrix} \begin{bmatrix} w^1 \\ \vdots \\ w^J \end{bmatrix} = \begin{bmatrix} \mathbf{w} \cdot \mathbf{x}_1 \\ \vdots \\ \mathbf{w} \cdot \mathbf{x}_B \end{bmatrix}$$

We're most of the way now...

Vectorizing logistic regression

- Computation we'd like to parallelize:
 - For each \mathbf{x} in the minibatch X_{batch} , compute

$$p \equiv \frac{1}{1 + e^{-\mathbf{x} \cdot \mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

$$\begin{bmatrix} \mathbf{w} \cdot \mathbf{x}_1 \\ \vdots \\ \mathbf{w} \cdot \mathbf{x}_B \end{bmatrix} + 1$$

We can overload the operator for the vector/matrix class so that adding 1 will "work"

Vectorizing logistic regression

```
class Matrix(object):  
    ...  
    def __add__(self, addend):  
        result = Matrix()  
        ....  
        return result
```

```
m = Matrix(...)  
z = m + 1
```

$$\begin{bmatrix} \mathbf{w} \cdot \mathbf{x}_1 \\ \vdots \\ \mathbf{w} \cdot \mathbf{x}_B \end{bmatrix} + 1$$

We can overload the operator for the vector/matrix class so that adding c will create a copy with c added component-wise

Vectorizing logistic regression

```
class Matrix(object):  
    ...  
    def __add__(self, addend):  
        result = Matrix()  
        ....  
        return result
```

```
m = Matrix(...)  
z = (m+1).exp()
```

$$\exp\left(\begin{bmatrix} \mathbf{w} \cdot \mathbf{x}_1 \\ \vdots \\ \mathbf{w} \cdot \mathbf{x}_B \end{bmatrix} + 1\right)$$

....and also define operations like `exp()` to work component-wise...

Vectorizing logistic regression

- Computation we'd like to parallelize:
 - For each \mathbf{x} in the minibatch X_{batch} , compute

$$p \equiv \frac{1}{1 + e^{-\mathbf{x} \cdot \mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

```
Xb = Matrix(...)
```

```
w = Matrix(...)
```

```
def logistic(X): return (X+1).exp().reciprocal()
```

```
p = logistic(Xb.dot(w))      # B rows, 1 column
```

Vectorizing logistic regression

- Computation we'd like to parallelize:
 - For each \mathbf{x} in the minibatch, compute

$$p \equiv \frac{1}{1 + e^{-\mathbf{x} \cdot \mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

$$\frac{\partial}{\partial w^j} \log P(Y = y | X = \mathbf{x}, \mathbf{w}) = (y - p)x^j$$

$$X_{batch} \mathbf{w} = \begin{bmatrix} \mathbf{w} \cdot \mathbf{x}_1 \\ \vdots \\ \mathbf{w} \cdot \mathbf{x}_B \end{bmatrix} \rightarrow \begin{bmatrix} \mathbf{p}_1 \\ \vdots \\ \mathbf{p}_B \end{bmatrix} \quad \begin{bmatrix} \mathbf{y}_1 - \mathbf{p}_1 \\ \vdots \\ \mathbf{y}_B - \mathbf{p}_B \end{bmatrix}$$

Vectorizing logistic regression

- Computation we'd like to parallelize:
 - For each \mathbf{x} in the minibatch, compute

$$p \equiv \frac{1}{1 + e^{-\mathbf{x} \cdot \mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

$$\frac{\partial}{\partial w^j} \log P(Y = y | X = \mathbf{x}, \mathbf{w}) = (y - p)x^j$$

```
def logistic(X): return (X+1).exp().reciprocal()
```

```
p = logistic(Xb.dot(w)) # B rows, 1 column
```

```
grad = Xb.dot(y - p).rowsum() * 1/B
```

```
w = w + grad*rate
```

EXAMPLE: NUMPY PACKAGE

```
1 import numpy as np
2 import numpy.random as random
3 from examples.utils.data_utils import gaussian_cluster_generator as make_data
4 # Predict the class using multinomial logistic regression (softmax regression).
5 def predict(w, x):
6     a = np.exp(np.dot(x, w))
7     a_sum = np.sum(a, axis=1, keepdims=True)
8     prob = a / a_sum
9     return prob
10 # Using gradient descent to fit the correct classes.
11 def train(w, x, loops):
12     for i in range(loops):
13         prob = predict(w, x)
14         loss = -np.sum(label * np.log(prob)) / num_samples
15         if i % 10 == 0:
16             print('Iter {}, training loss {}'.format(i, loss))
17             # gradient descent
18             dy = prob - label
19             dw = np.dot(data.T, dy) / num_samples
20             # update parameters; fixed learning rate of 0.1
21             w -= 0.1 * dw
22
23 # Initialize training data.
24 num_samples = 10000
25 num_features = 500
26 num_classes = 5
27 data, label = make_data(num_samples, num_features, num_classes)
28
29 # Initialize training weight and train
30 weight = random.randn(num_features, num_classes)
31 train(weight, data, 100)
```

Binary to softmax logistic regression

$$p \equiv \frac{1}{1 + e^{-\mathbf{x} \cdot \mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

$$X_{batch} \mathbf{w} = \begin{bmatrix} x_1^1 & \cdots & x_1^J \\ \vdots & \ddots & \vdots \\ x_B^1 & \cdots & x_B^J \end{bmatrix} \begin{bmatrix} w^1 \\ \vdots \\ w^J \end{bmatrix} = \begin{bmatrix} \mathbf{w} \cdot \mathbf{x}_1 \\ \vdots \\ \mathbf{w} \cdot \mathbf{x}_B \end{bmatrix}$$

Binary to softmax logistic regression

$$p \equiv \frac{1}{1 + e^{-\mathbf{x} \cdot \mathbf{w}}} = \frac{1}{1 + \exp(-\sum_j x^j w^j)}$$

$$p^y \equiv \frac{\exp(\mathbf{x} \cdot \mathbf{w}^y)}{\sum_{y'} \exp(\mathbf{x} \cdot \mathbf{w}^{y'})}$$

$$XW = \begin{bmatrix} x_1^1 & \cdots & x_1^J \\ \vdots & \ddots & \vdots \\ x_B^1 & \cdots & x_B^J \end{bmatrix} \begin{bmatrix} w^1 \\ \vdots \\ w^J \end{bmatrix} = \begin{bmatrix} \mathbf{w} \cdot \mathbf{x}_1 \\ \vdots \\ \mathbf{w} \cdot \mathbf{x}_B \end{bmatrix}$$

$$XW = \begin{bmatrix} x_1^1 & \cdots & x_1^J \\ \vdots & \ddots & \vdots \\ x_B^1 & \cdots & x_B^J \end{bmatrix} \begin{bmatrix} w_1^{y1} & \cdots & w_1^{yK} \\ \vdots & \ddots & \vdots \\ w_J^{y1} & \cdots & w_J^{yK} \end{bmatrix} = \begin{bmatrix} \mathbf{w}^{y1} \cdot \mathbf{x}_1 & \cdots & \mathbf{w}^{yK} \cdot \mathbf{x}_1 \\ \vdots & \ddots & \vdots \\ \mathbf{w}^{y1} \cdot \mathbf{x}_B & \cdots & \mathbf{w}^{yK} \cdot \mathbf{x}_B \end{bmatrix}$$

```
1 import numpy as np
2 import numpy.random as random
3 from examples.utils.data_utils import gaussian_cluster
4 # Predict the class using multinomial logistic regression
5 def predict(w, x):
6     a = np.exp(np.dot(x, w))
7     a_sum = np.sum(a, axis=1, keepdims=True)
8     prob = a / a_sum
9     return prob
```

Matrix multiply,; then
exponentiate
component-wise

prob will have B rows
and K columns, and each
row will sum to 1

Sum the columns to get
the denominator;
keepdim=True means...

... that this line will work
correctly even though 'a'
and 'a_sum' have
different shapes

$$p^y \equiv \frac{\exp(\mathbf{x} \cdot \mathbf{w}^y)}{\sum_{y'} \exp(\mathbf{x} \cdot \mathbf{w}^{y'})}$$

$$XW = \begin{bmatrix} x_1^1 & \dots & x_1^J \\ \vdots & \ddots & \vdots \\ x_B^1 & \dots & x_B^J \end{bmatrix} \begin{bmatrix} w_1^{y1} & \dots & w_1^{yK} \\ \vdots & \ddots & \vdots \\ w_J^{y1} & \dots & w_J^{yK} \end{bmatrix} = \begin{bmatrix} \mathbf{w}^{y1} \cdot \mathbf{x}_1 & \dots & \mathbf{w}^{yK} \cdot \mathbf{x}_1 \\ \vdots & \ddots & \vdots \\ \mathbf{w}^{y1} \cdot \mathbf{x}_B & \dots & \mathbf{w}^{yK} \cdot \mathbf{x}_B \end{bmatrix}$$

```
1 import numpy as np
2 import numpy.random as random
3 from examples.utils.data_utils import gaussian_cluster_generator as make_data
4 # Predict the class using multinomial logistic regression (softmax regression).
5 def predict(w, x):
6     a = np.exp(np.dot(x, w))
7     a_sum = np.sum(a, axis=1, keepdims=True)
8     prob = a / a_sum
9     return prob
10
11 # Using gradient descent to fit the correct classes.
12 def train(w, x, loops):
13     for i in range(loops):
14         prob = predict(w, x)
15         loss = -np.sum(label * np.log(prob)) / num_samples
16         if i % 10 == 0:
17             print('Iter {}, training loss {}'.format(i, loss))
18             # gradient descent
19             dy = prob - label
20             dw = np.dot(data.T, dy) / num_samples
21             # update parameters; fixed learning rate of 0.1
22             w -= 0.1 * dw
23
24 # Initialize training data.
25 num_samples = 10000
26 num_features = 500
27 num_classes = 5
28 data, label = make_data(num_samples, num_features, num_classes)
29
30 # Initialize training weight and train
31 weight = random.randn(num_features, num_classes)
32 train(weight, data, 100)
```

```
1 import numpy as np
2 import numpy.random as random
3 from examples.utils.data_utils import gaussian_cluster
```

```
4 # Predict the class using multinomial logistic regr
```

$$x.T dy = \begin{bmatrix} x_1^1 & \dots & x_B^1 \\ \vdots & \ddots & \vdots \\ x_1^J & \dots & x_B^J \end{bmatrix} \cdot \begin{bmatrix} dy_{x_1}^{y1} & \dots & dy_{x_1}^{yK} \\ \vdots & \ddots & \vdots \\ dy_{x_B}^{y1} & \dots & dy_{x_B}^{yK} \end{bmatrix}$$

Error on each example x in batch and each class y

```
10 def train(w, x, loops):
11     for i in range(loops):
12         prob = predict(w, x)
13         loss = -np.sum(label * np.log(prob)) / num_sam
```

python bug: should be $x.T$ (transpose)

```
14 # gradient descent
15 dy = prob - label
16 dw = np.dot(data.T, dy) / num_samples
17 # update parameters; fixed Learning rate
18 w -= 0.1 * dw
```

The gradient step!

```
19 # Initialize training data.
```

$$\frac{\partial}{\partial w^j} \log P(Y = y | X = \mathbf{x}, \mathbf{w}) = \underline{(y - p)x^j}$$

```
25 weight = random.randn(num_features, num_classes)
26 train(weight, data, 100)
```


So this will run in parallel on a GPU?

http://minpy.readthedocs.io/en/latest/get-started/logistic_regression.html

```
1 import numpy as np
2 import numpy.random as random
3 from examples.utils.data_utils import gaussian_cluster_generator as make_data
4 # Predict the class using multinomial logistic regression (softmax regression).
5 def predict(w, x):
6     a = np.exp(np.dot(x, w))
7     a_sum = np.sum(a, axis=1, keepdims=True)
8     prob = a / a_sum
9     return prob
10
11 # Using gradient descent to fit the correct classes.
12 def train(w, x, loops):
13     for i in range(loops):
14         prob = predict(w, x)
15         loss = -np.sum(label * np.log(prob)) / num_samples
16         if i % 10 == 0:
17             print('Iter {}, training loss {}'.format(i, loss))
18         # gradient descent
19         dy = prob - label
20         dw = np.dot(data.T, dy) / num_samples
21         # update parameters; fixed learning rate of 0.1
22         w -= 0.1 * dw
23
24 # Initialize training data.
25 num_samples = 10000
26 num_features = 500
27 num_classes = 5
28 data, label = make_data(num_samples, num_features, num_classes)
29
30 # Initialize training weight and train
31 weight = random.randn(num_features, num_classes)
32 train(weight, data, 100)
```

No: we're not there yet....

So this will run in parallel on a GPU? Not yet....

http://minpy.readthedocs.io/en/latest/get-started/logistic_regression.html

```
1 import numpy as np
2 import numpy.random as random
3 from examples.utils.data_utils import gaussian_cluster_generator as make_data
4
5 # Predict the class using multinomial logistic regression (softmax regression).
6 def predict(w, x):
7     a = np.exp(np.dot(x, w))
8     a_sum = np.sum(a, axis=1, keepdims=True)
9     prob = a / a_sum
10    return prob
11
12 # Using gradient descent to fit the correct classes.
13 def train(w, x, loops):
14     for i in range(loops):
15         prob = predict(w, x)
16         loss = -np.sum(label * np.log(prob)) / num_samples
17         if i % 10 == 0:
18             print('Iter {}, training loss {}'.format(i, loss))
19         # gradient descent
20         dy = prob - label
21         dw = np.dot(data.T, dy) / num_samples
22         # update parameters; fixed learning rate
23         w -= 0.1 * dw
24
25 # Initialize training data.
26 num_samples = 10000
27 num_features = 500
28 num_classes = 5
29 data, label = make_data(num_samples, num_features, num_classes)
30
31 # Initialize training weight and train
32 weight = random.randn(num_features, num_classes)
33 train(weight, data, 100)
34
35 ~~
```

Option 1: switch from numpy (old package) to cupy (new GPU-oriented package)

Option 2: switch to a package that will compile to a GPU and also compute the gradients for you (like Theano, Tensorflow, Torch, ...)