Introduction to CUDA Programming

15-418/618: Parallel Computer Architecture and Programming Recitation 2, February 2, 2018,

Roadmap

- Review
- Possibly Helpful Tips
- Performance Optimization Example

Vocabulary

- CPU: A central processor unit, i.e. a normal processor
- **GPU:** A graphics processing unit, i.e. what we are learning about
- Host: The "normal computer" to which the GPU is connected
 - Of especial note are the CPU(s) and memory
- Device: The GPU and its memory
- CUDA: Compute Unified Device Architecture. nVidia's framework for utilizing their GPUs for general purpose programming
 - OpenCL: Open Computing Language. The "generic version"

Vocabulary, cont, cont

- Global memory: Device memory shared across the various blocks
 - CUDAMalloc(), CUDAMemcpy(), CUDAFree()
- Shared memory: Memory shared only by threads within the associated block (not across blocks)
 - __shared__

Vocabulary, cont

- Kernel: The work, written as a function, to be parallelized across the GPU's cores.
- **Thread**: An abstraction for the work associated with an instance of the kernel.
- **Thread Block**: A partition of threads and associated work that will be dispatched to a *Streaming Media (SM)* processor, basically a GPU.
- **Block**: See *Thread Block*
- Grid: Set of all blocks

Vocabulary, cont, cont, cont

- CUDA Core: A single graphics processor core. Within the CUDA architecture, these aren't one-offs, but architected into Streaming Multiprocessors (SMs).
- Streaming Multiprocessor (SM): A collection of CUDA Cores architected together to form a single GPU. Threads within a thread block concurrently execute on an SM.
- Warp: A division of a block created within the SM to assign work to cores. Warps aren't schedule until a core is available for each thread within the warp.

Syntax, Etc.

- nvcc: nVidia C compiler. Can compiler host and device code.
- __shared___: Qualifier to declare a variable in shared (per thread block) memory
- __global__: Qualifier to place a function into device memory, for execution onto the device, but enabling it to be callable from the host.
- cudaMalloc(), cudaMemcpy(), cudaFree()
 - Allocates, Frees, and copies to/from device memory.
 - cudaMemcpyHostToDevice/cudaMemcpyDeviceToHost flag sets direction of copy
- __syncthreads()__
 - barrier to ensure all threads get there before any continue.

Syntax, Etc, cont

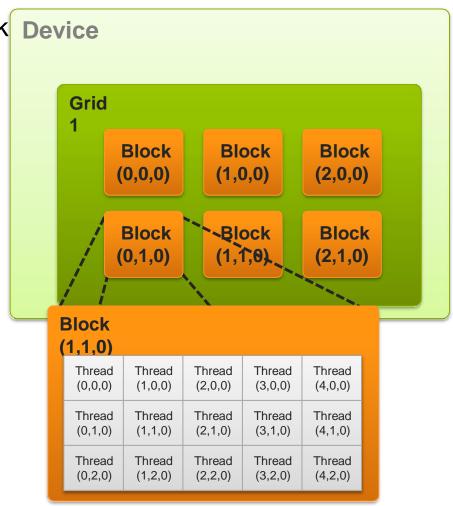
- someGlobalFunctionName<<<N,M>>>
 - "Kernel Launch"
 - N thread blocks
 - M threads per thread block
- blockidx: block index within the (up to 3D) grid
 - threadIdx.x is 1D index
- threadIdx: thread index within the (up to 3D) thread block
 - threadIdx.x is 1D index
- int index = threadIdx.x + blockIdx.x * M;
 - Global thread index, given M threads per block
- blockDim, gridDim
 - 3D block and grid dimensions
 - E.g., blockDim.x, gridDim.x, etc

A Picture Worth Some Number of Words

- A kernel is launched as a grid of block of threads
 - blockIdx and threadIdx are 3D
 - We showed only one dimension (x)
- Built-in variables:
 - threadIdx
 - blockIdx
 - blockDim
 - gridDim

CUDA C/C++ Basics

Cyril Zeller, NVIDIA Corporation Supercomputing 2011 Tutorial



Roadmap

- Review
- Possibly Helpful Tips
- Performance Optimization Example

Did it fly? Wrap All CUDA Library Calls

```
Definitions in file reduce.cu:
// Support for CUDA error checking
// Wrapper for CUDA functions
#define CHK(ans) gpuAssert((ans), __FILE__, __LINE__);
// Checker
inline void gpuAssert (CUDAError t code, const char *file, int line)
  if (code != CUDASuccess) {
    fprintf(stderr, "GPUassert: %s %s %s\n",
    CUDAGetErrorString(code), file, line);
// Cannot wrap kernel launches. Instead, insert this after each
// kernel launch.
#define POSTKERNEL CHK(CUDAPeekAtLastError())
```

Your Old Friend. Better. Than. Ever. cuda-gdb: Can get data off of the device

- break reduce.cu:90
 - Set breakpoint corresponding to line 90 of file reduce.cu.
- print ((@global float *) srcVecDevice)[1]
 - Print contents of array in device memory
- CUDA thread 2
 - Shift focus to specified thread number
- info locals
 - Prints values of all currently-active local variables
 - CUDA info threads
 - Prints status of threads (in current block)

Some Advice

- Don't wire down constants
- Don't assume special properties of N
 - Multiple of block size, power of 2, ...
- Use function or macro to do rounding-up division
- Write checker code
 - Overall functionality
 - Individual steps on device
 - Must transfer data back to host to check
- Avoid printf within kernel functions
 - Only use on small examples or too much unordered output.
- Get the algorithm & abstract implementation and benchmark right before attempting low-level optimizations
 - Exploiting the various memory categories on device
 - Exploiting properties specific to block level

Some More Advice

- Even more so than with C programs, out-of-bounds memory writes in CUDA lead to bizarre and erratic beahvior.
 - Write bounds checking code that gets invoked when program is run in DEBUG mode
- It's possible to put printf's in kernel code, but don't rely on them
 - Often nothing gets printed, or values printed are incorrect.
 - Just because nothing prints, it doesn't mean that part of the code wasn't reached.
- Write host code that duplicates functionality of different parts of CUDA code
 - In debug mode, transfer results back to host memory and check values against this code

Why Is printf() weird?

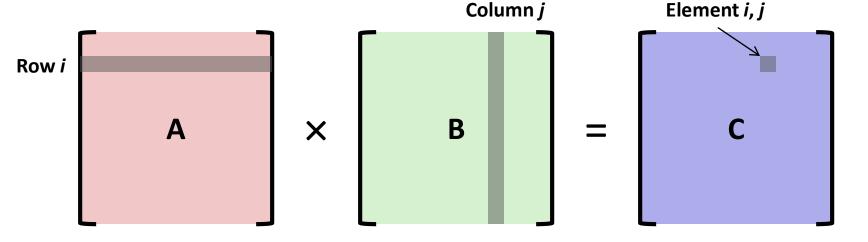
- printf() output is stored in a circular buffer of a fixed size.
 - If the buffer fills, old output will be overwritten. The buffer's size defaults to 1MB and can be configured with CUDADeviceSetLimit(CUDALimitPrintfFifoSize, size_t size).

This buffer is flushed only for

- the start of a kernel launch
- synchronization (e.g. CUDADeviceSynchronize())
- blocking memory copies (e.g. CUDAMemcpy(...))
- module load/unload
- context destruction
- Note: The list above does not include program exit.
 - If the call to CUDADeviceSynchronize() was removed from the example program above, the we would see no output
- Concurrency serialized upon output

Credit: Steven Fackler, Former 418/618 TA, SCS'13

Application Example: N × N Matrix Multiplication



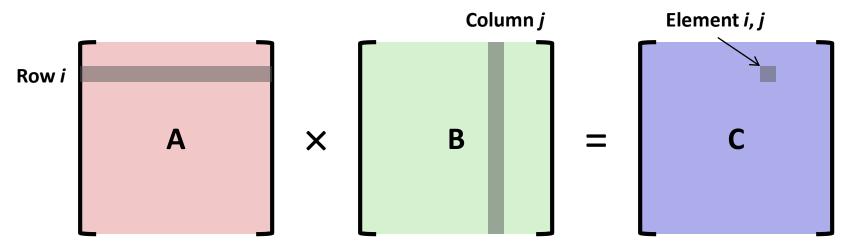
$$c_{i,j} = \mathop{\mathring{\mathbf{a}}}_{k=0}^{N-1} a_{i,k} \times b_{k,j}$$

Complexity

- N³ multiplications
- N³ additions
- Assume row-major access

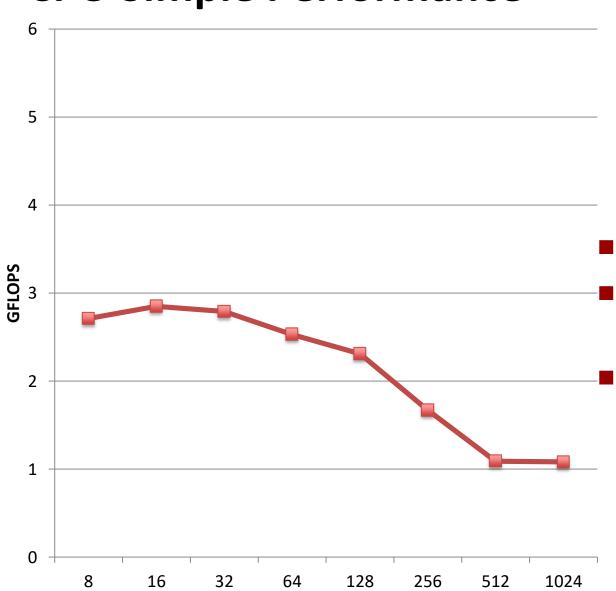
#define RM(r, c, width) ((r) * (width) + (c))

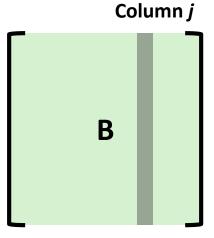
Matrix Multiplication: Simple CPU Implementation



```
void multMatrixSimple(int N, float *matA, float *matB, float *matC) {
    for (int i = 0; i < N; i++)
        for (int j = 0; j < N; j++) {
            float sum = 0.0;
            for (int k = 0; k < N; k++)
                 sum += matA[RM(i,k,N)] * matB[RM(k,j,N)];
            matC[RM(i,j,N)] = sum;
        }
}</pre>
```

CPU Simple Performance



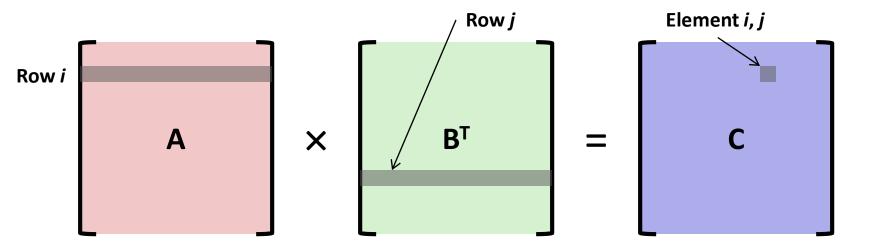


Measured in GFLOPS

Drops off for large values of N

B has bad access pattern

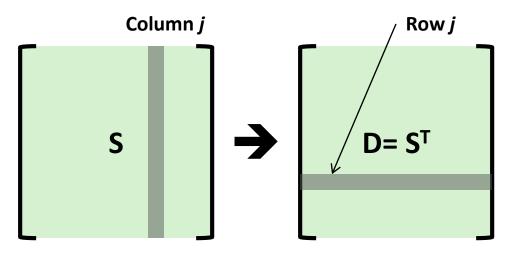
Optimization #1: Pretranspose



- Transposed version of B has better access pattern
 - Transpose once
 - Use each element N times

$$c_{i,j} = \mathop{igcap}\limits_{k=0}^{N-1} a_{i,k} imes b_{j,k}^T$$

Transposing a Matrix



Column-major ordering of elements

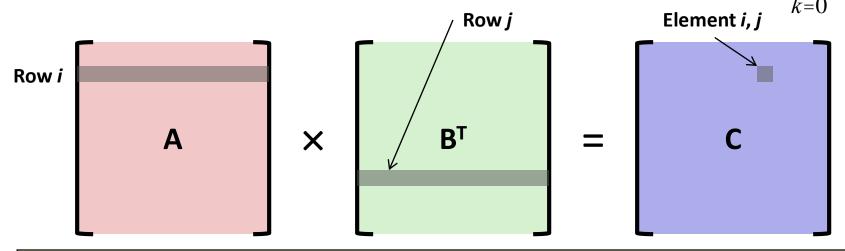
```
#define CM(r, c, height) ((c) * (height) + (r))
```

Transposing converts from row-major to column-major order

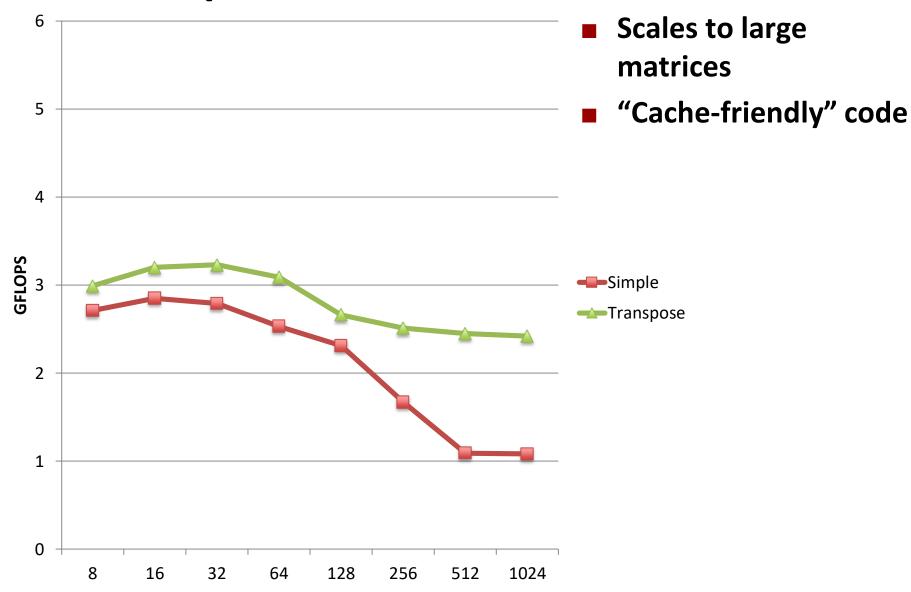
```
void transposeMatrix(int N, float *matS, float *matD) {
    for (int i = 0; i < N; i++)
        for (int j = 0; j < N; j++)
        matD[CM(i,j,N)] = matS[RM(i,j,N)];
}</pre>
```

Matrix Multiplication: Pretranspose Implementation $C_{i,j}$

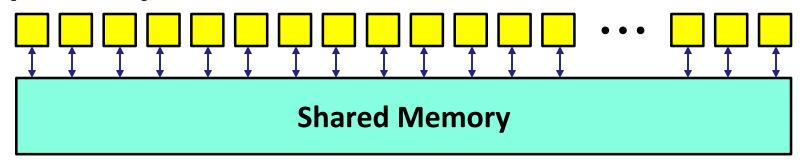




Pretranspose Performance



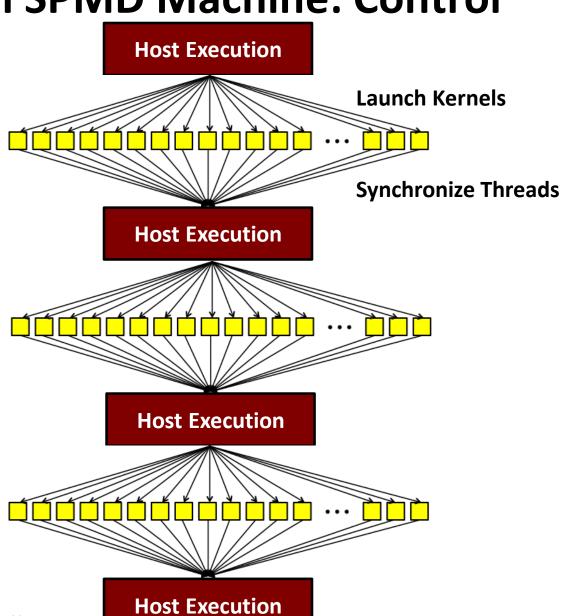
Abstract Single Program Multiple Data (SPMD) Model



- M Processors, all executing same code
 - Called "kernels"
 - M based on problem size
- Share common global memory
 - And also have private memory for local variables
 - Make no assumptions about effect of memory access conflicts
- No synchronization primitives
- Called threads, but not at all like pthreads
 - Very simple & lightweight
- All execute the same program
 15-418/618: Lightly edited from Prof. Bryant, 15-418/618, Spring 2017

Interacting with SPMD Machine: Control

- Overall execution managed by code executing on host
- Launch set of kernels
 - Number & kernel function can vary with each launch
- Wait until all completed
 - Explicit or implicit synchronization
- Repeat as necessary



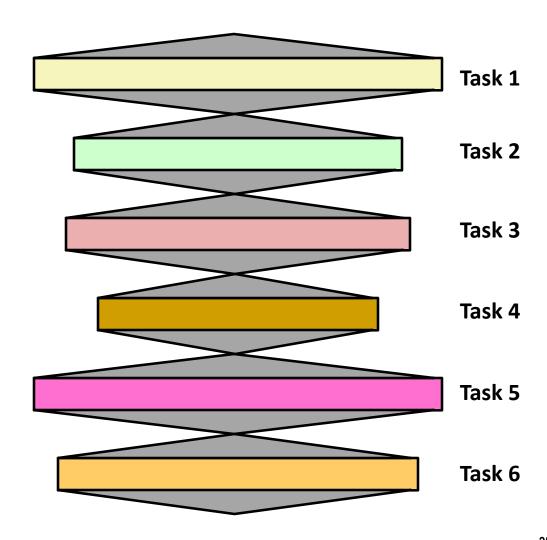
Structure of SPMD Program

Concept

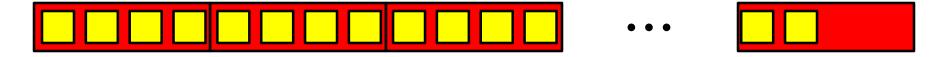
- Partition computation into sequence of tasks
- Perform each task over all data with single operation

Performance Limitations

- Synchronization requires waiting for slowest task
- No locality of data
- No locality of synchronization



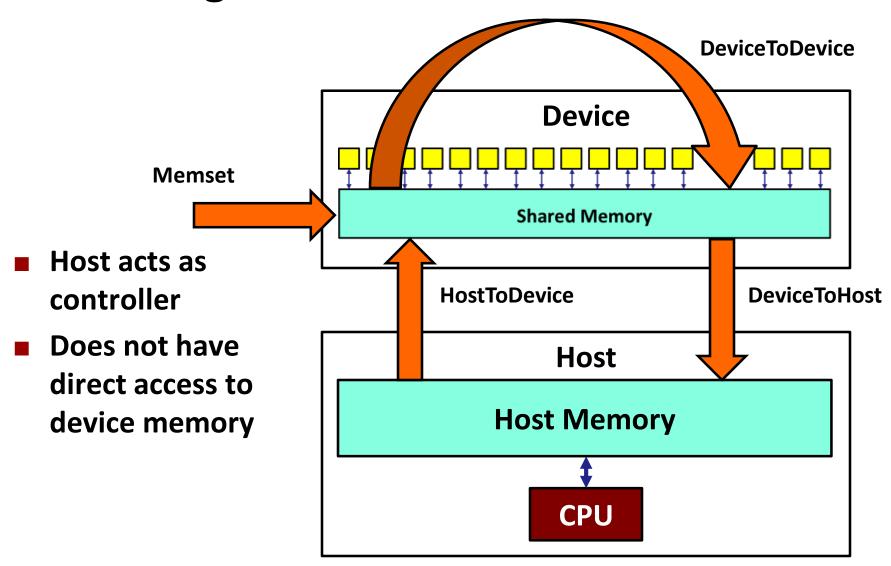
Block/Thread Notation



- Idea (One-dimensional version)
 - Executing threads grouped into blocks
 - Each contains same number of threads
 - Host program specifies block size (blockDim.x)
 - Host program makes sure there are enough blocks to generate N threads
 - Some threads in last block should not get used

```
__global___ void
inplaceReduceKernel(int length, int nlength, float *data) {
    int idx = blockIdx.x * blockDim.x + threadIdx.x;
    if (idx < nlength) {
        . . . .
    }
}
```

Interacting with SPMD Machine: Data



CUDA Program

- CUDA file (.cu) contains mix of device code & host code
 - It's up to you to understand which is which!

Device Code

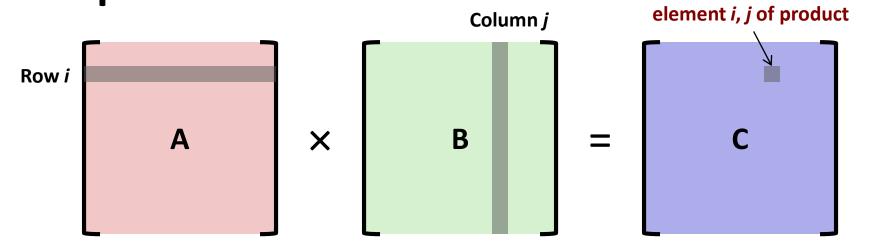
- Kernels (__global___)
 - Code for threads
 - Must only reference device memory
- Device functions (__device__)
 - Called by kernels
 - Only reference device memory
 - Do not generate new threads

```
__global__ void
inplaceReduceKernel(int length, int nlength, float *data) {
   int idx = blockIdx.x * blockDim.x + threadIdx.x;
   if (idx < nlength) {
        . . . .
   }
}</pre>
```

CUDA Program (cont.)

- CUDA file (.cu) contains mix of device code & host code
 - It's up to you to understand which is which!
- Host Code
 - Conventional C/C++
 - Can only reference host memory
 - But, can have pointers to device memory
 - Manages the launching of threads
 - Manages movement of data between host & device memories

Matrix Multiplication: Simple CUDA Implementation Each thread computes



```
__global___ void
CUDASimpleKernel(int N, float *dmatA, float *dmatB, float *dmatC) {
   int i = blockIdx.y * blockDim.y + threadIdx.y;
   int j = blockIdx.x * blockDim.x + threadIdx.x;
   if (i >= N || j >= N)
        return;
   float sum = 0.0;
   for (int k = 0; k < N; k++) {
        sum += dmatA[RM(i,k,N)] * dmatB[RM(k,j,N)];
   }
   dmatC[RM(i,j,N)] = sum;
}</pre>
```

Host Code Example

LBLK = 32 32 X 32 = 1024 threads / block

Launch kernels to perform vector product

Useful stuff

• Compute $\hat{\mathfrak{g}}_n / d\hat{\mathfrak{g}}$

```
// Integer division, rounding up
static inline int updiv(int n, int d) {
   return (n+d-1)/d;
}
```

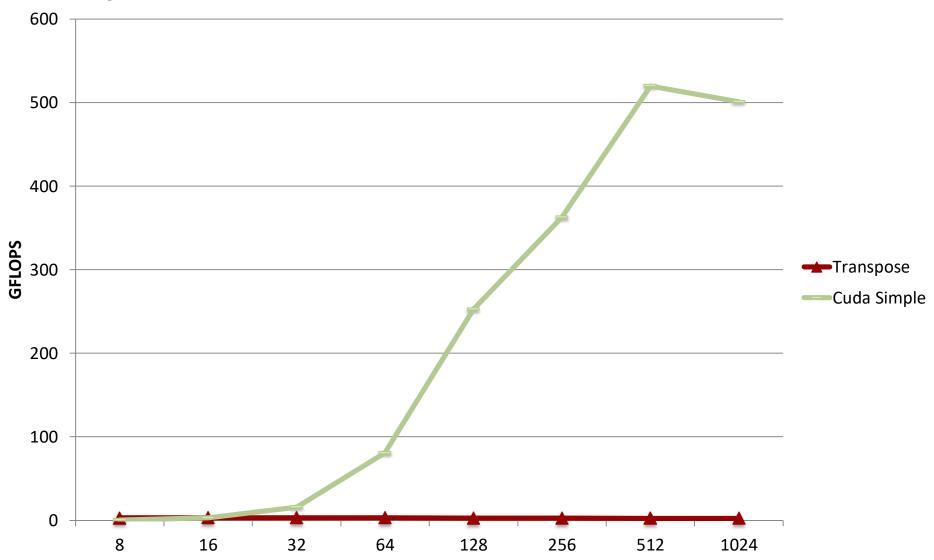
- Setting number of threads per block:
 - Should be multiple of 32
 - Max value = 1024

Host Code Example (cont).

```
void CUDAMultiply(int N, float *aData, float *bData, float *cData) {
    float *aDevData, *bDevData, *cDevData
    CUDAMalloc((void **) &aDevData, N*N * sizeof(float));
    CUDAMalloc((void **) &bDevData, N*N * sizeof(float));
    CUDAMalloc((void **) &cDevData, N*N * sizeof(float));
    CUDAMemcpy (aDevData, aData, N*N * sizeof(float),
               CUDAMemcpyHostToDevice);
    CUDAMemcpy (bDevData, bData, N*N * sizeof(float),
               CUDAMemcpyHostToDevice);
    CUDAMultMatrixSimple(N, aDevData, bDevData, tDevData);
    CUDAMemcpy (cData, cDevData, N*N * sizeof(float),
               CUDAMemcpyDeviceToHost);
    CUDAFree (aDevData); CUDAFree (bDevData); CUDAFree (cDevData);
```

Observe: Host can hold pointers to device memory, but cannot read or write device memory locations

Simple CUDA Performance



Inverted Indexing Accessing Pattern

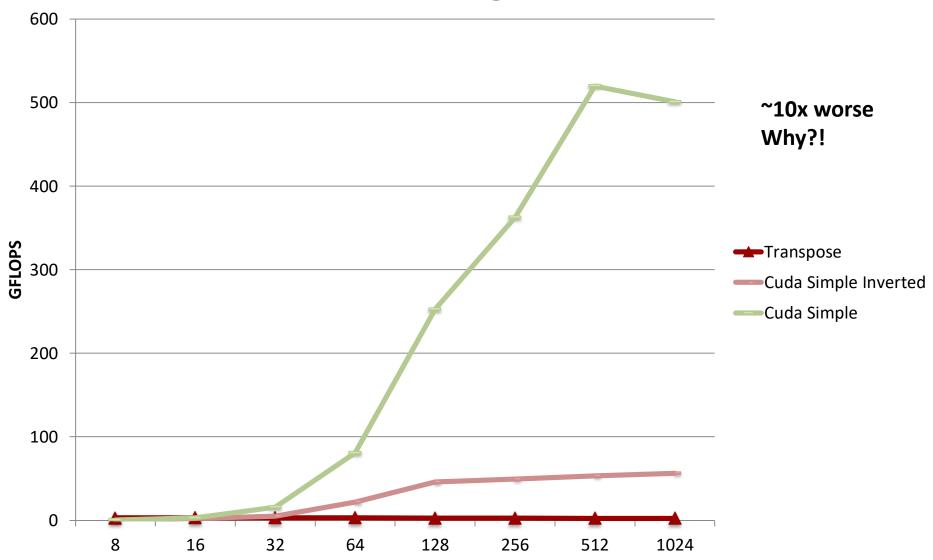
Regular

```
__global___ void
CUDASimpleKernel(int N, float *dmatA, float *dmatB, float *dmatC) {
  int i = blockIdx.y * blockDim.y + threadIdx.y;
  int j = blockIdx.x * blockDim.x + threadIdx.x;
  . . . .
}
```

Inverted

```
__global__ void
CUDASimpleKernelOld(int N, float *dmatA, float *dmatB, float *dmatC) {
  int i = blockIdx.x * blockDim.x + threadIdx.x;
  int j = blockIdx.y * blockDim.y + threadIdx.y;
  . . . .
}
```

CUDA Inverted Indexing Performance



What's the Difference?

Regular

```
int i = blockIdx.y * blockDim.y + threadIdx.y;
  int j = blockIdx.x * blockDim.x + threadIdx.x;
Inverted
  int i = blockIdx.x * blockDim.x + threadIdx.x;
  int j = blockIdx.y * blockDim.y + threadIdx.y
```

- CUDA threads numbered within block in row-major order
 - X = column number, Y = row number
- Threads with same value of Y map to single warp.
- Threads with same value of Y and consecutive values of X map to consecutive positions in single warp
- When single warp accesses consecutive memory locations, do block read or write
- When single warp accesses separated memory locations, 15-418/618: Lightly edited momeroi. Sry gather, (read) or scatter (write)

Impact on Memory Referencing: Regular

```
int i = blockIdx.y * blockDim.y + threadIdx.y;
int j = blockIdx.x * blockDim.x + threadIdx.x;
```

```
Read A = dmatA[RM(i,k,N)]; Threads in warp reference single location

Read B = dmatB[RM(k,j,N)]; Threads in warp do block read

Write B dmatC[RM(i,j,N)] Threads in block do block write
```

Threads within warp have:

- same value of k
- same value of i
- consecutive values of j
- Warp reads & writes match memory organization

Impact on Memory Referencing: Inverted

```
int i = blockIdx.x * blockDim.x + threadIdx.x;
int j = blockIdx.y * blockDim.y + threadIdx.y;
```

```
Read B = dmatB[RM(k,j,N)]; Thread

Write B dmatC[RM(i,j,N)] = Thread
```

= dmatA[RM(i,k,N)];

Threads in warp do gather

Threads in warp reference single location

Threads in block do scatter

Threads within warp have:

- same value of k
- consecutive values of i
- same value of j
- Warp reads/writes does not match memory organization

Read A

Relation to Hardware



Optimizing memory instruction performance

- Load faster than gather
- Store faster than scatter

Avoiding memory conflicts

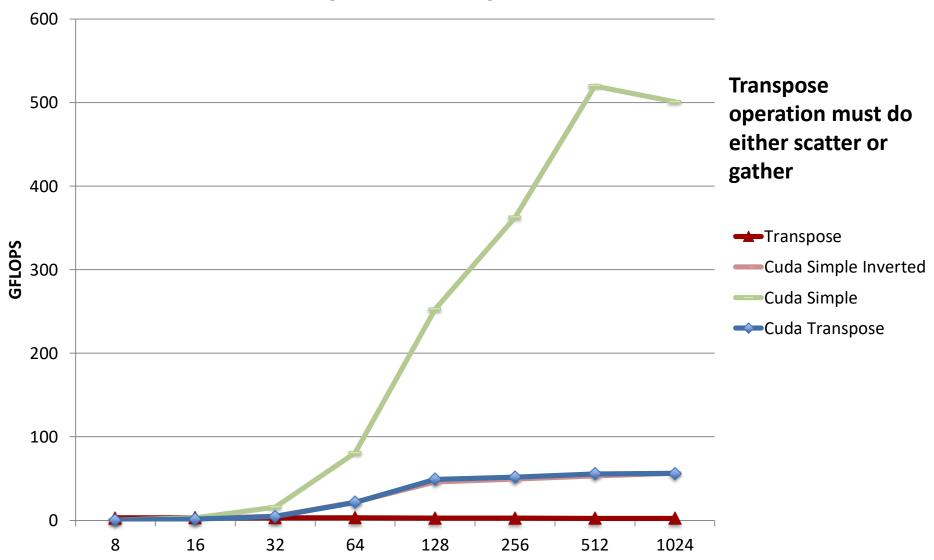
Inverted code has multiple warps competing for same block of memory

Pretransposing with CUDA

```
/* Transpose matrix */
    __global___ void
CUDATransposeKernel(int N, const float *dmatS, float *dmatD) {
    int i = blockIdx.y * blockDim.y + threadIdx.y;
    int j = blockIdx.x * blockDim.x + threadIdx.x;
    if (i >= N || j >= N)
        return;
    dmatD[CM(i,j,N)] = dmatS[RM(i,j,N)];
}
```

```
__global___ void
CUDATransposedKernel(int N, float *dmatA, float *dmatB, float *dmatC) {
  int i = blockIdx.y * blockDim.y + threadIdx.y;
  int j = blockIdx.x * blockDim.x + threadIdx.x;
  if (i >= N || j >= N)
            return;
  float sum = 0.0;
  for (int k = 0; k < N; k++) {
        sum += dmatA[RM(i,k,N)] * dmatB[RM(j,k,N)];
  }
  dmatC[RM(i,j,N)] = sum;
}</pre>
```

CUDA Pretranspose Implementations



Thinking About CUDA

NVIDIA GTX 680 (2012)

NVIDIA Kepler GK104 architecture SMX unit (one "core")



GPU Hierarchy

Block Level

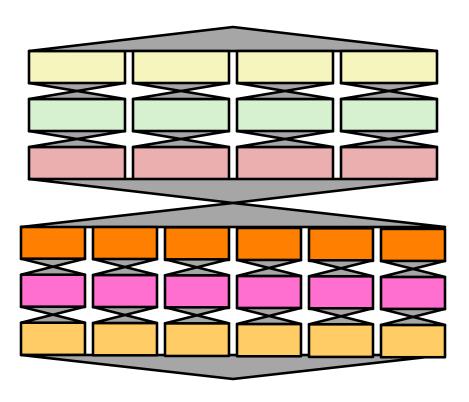
- Programmer partitions problem into blocks of K threads each
 - 32 ≤ K ≤ 1024
 - Multiple of 32
- Within block, have access to fast shared memory
- Within block, can synchronize with syncthreads ()

Warp Level

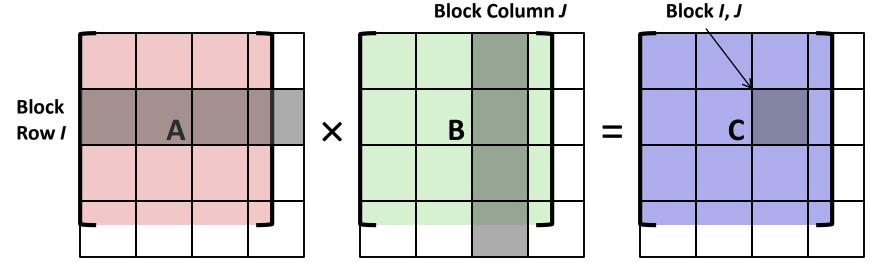
- Each block implemented as set of warps
 - 32 threads each
- Implemented using "SIMT" processor
 - Single-instruction, multiple threads
 - Guarantees stay synchronized

Programming with Blocks

- Localize computation within blocks
- Each performs sequence of tasks
- Each uses shared memory and local synchronization



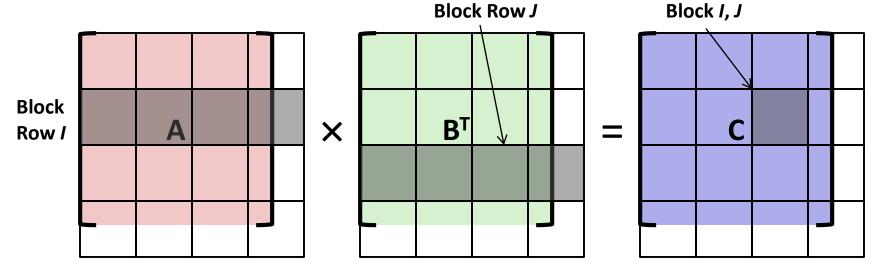
MM Optimization #2: Partitioning into Blocks



$$C_{I,J} = \mathop{\bigcirc}\limits_{K=0}^{N_b-1} A_{I,K} \times B_{K,J}$$

- Generate results on block-byblock basis
- Localizes access to A and B
 - N need not be multiple of block size

CPU-based Blocked Implementation



$$C_{I,J} = \mathop{ \mathring{\mathbf{S}}}_{K=0}^{N_b-1} A_{I,K} \times B_{K,J}$$

Use pretranspose

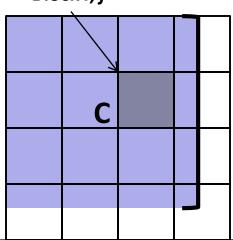
Required for performance

Structure

- Outer loops index over blocks
- Inner loops compute product for single block
- Block size SBLK = 8

Blocked Multiplication Implementation: Outer Loops

Block i, j



$$C_{I,J} = \mathop{ \stackrel{N_b-1}{\overset{}{\circ}}}_{K=0} A_{I,K} \times B_{K,J}$$

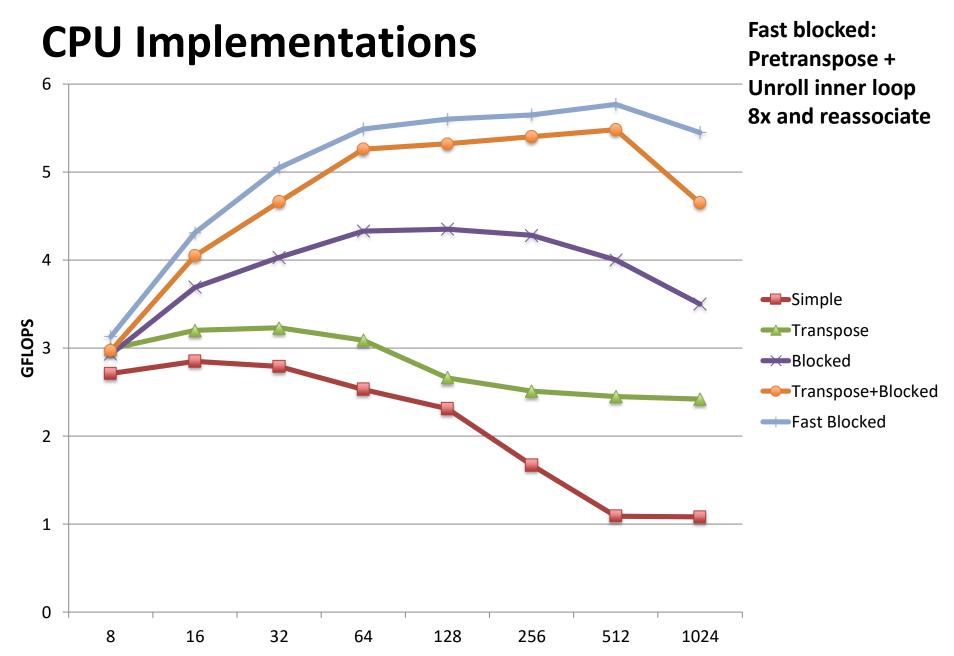
Look at actual code to see how it handles cases where N is not multiple of block size

Blocked Multiplication Implementation: Inner Loops

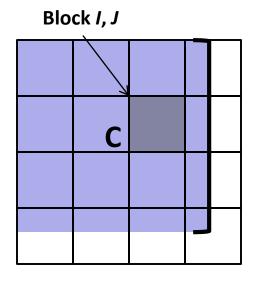
$$c_{i+bi,j+bj} = \mathop{\text{a}}_{bk=0}^{b-1} a_{i+bi,k+bk} \times b_{j+bj,k+bk}^{T}$$

- i, j, k provide starting indices of blocks
- bi, bj, bk provide offsets within blocks

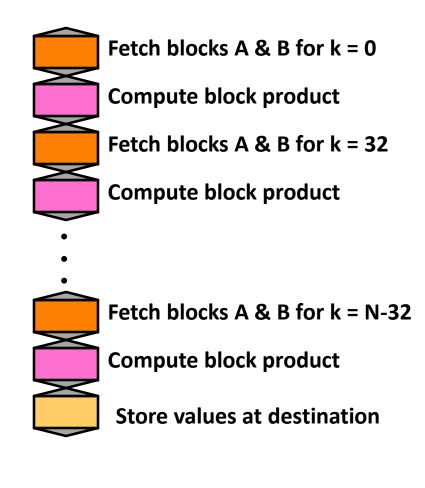
```
for (int bi = 0; bi < SBLK; bi++)
  for (int bj = 0; bj < SBLK; bj++) {
    float sum = 0.0;
    for (int bk =0; bk < SBLK; bk++)
        sum += matA[RM(i+bi,k+bk,N)] * tranB[RM(j+bj,k+bk,N)];
    matC[RM(i+bi,j+bj,N)] += sum;
}</pre>
```



Blocking with CUDA



- Block size LBLK = 32
- Use one CUDA block for each block of destination matrix
- Enough CUDA blocks to cover C
- Each thread in block accumulates single destination value



CUDA Block Kernel Structure

```
global void
CUDABlockKernel(int N, float *dmatA, float *dmatB, float *dmatC) {
    int i = blockIdx.x * blockDim.x + threadIdx.x;
    int j = blockIdx.y * blockDim.y + threadIdx.y;
    int bi = threadIdx.x;
    int bj = threadIdx.y;
    float sum = 0.0; // Accumulate result for C[i][j]
   // Shared space for two submatrices of A and B
    shared float subA[LBLK*LBLK];
    shared float subB[LBLK*LBLK];
   Loop over values of k
    if (i < N \&\& j < N)
       dmatC[RM(i,j,N)] = sum;
```

- Block size LBLK = 32
 - blockDim.x = blockDim.y = 32
- i, j index into source and destination arrays
- bi, bj index local arrays

CUDA Block Loop Structure

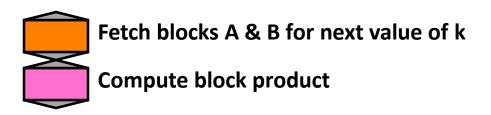
```
for (int k = 0; k < N; k+= LBLK) {
    Fetch elements bi, bj for local arrays subA and subB

    // Wait until entire block gets filled
    __syncthreads();

    Compute contribution to element i, j of output

    // Wait until all products computed
    syncthreads();
}</pre>
```

- Within loop, each thread plays two distinct roles
 - Fetch elements from source arrays into shared memory
 - Compute one element of subblock product



Fetching Blocks

```
if (i < N && k+bj < N) {
    subA[RM(bi,bj,LBLK)] = dmatA[RM(i,k+bj,N)];
} else {
    subA[RM(bi,bj,LBLK)] = 0.0;
}
if (j < N && k+bi < N) {
    subB[RM(bi,bj,LBLK)] = dmatB[RM(k+bi,j,N)];
} else {
    subB[RM(bi,bj,LBLK)] = 0.0;
}</pre>
```



Fetch blocks A & B for next value of k

- k is multiple of LBLK
 - Coarse-grained
- Fetch element i, k+bj from A to get subA[bi,bj]
- Fetch element k+bi, j from B to get subB[bi,bj]
- Set to 0 if out of range

Computing Block Product

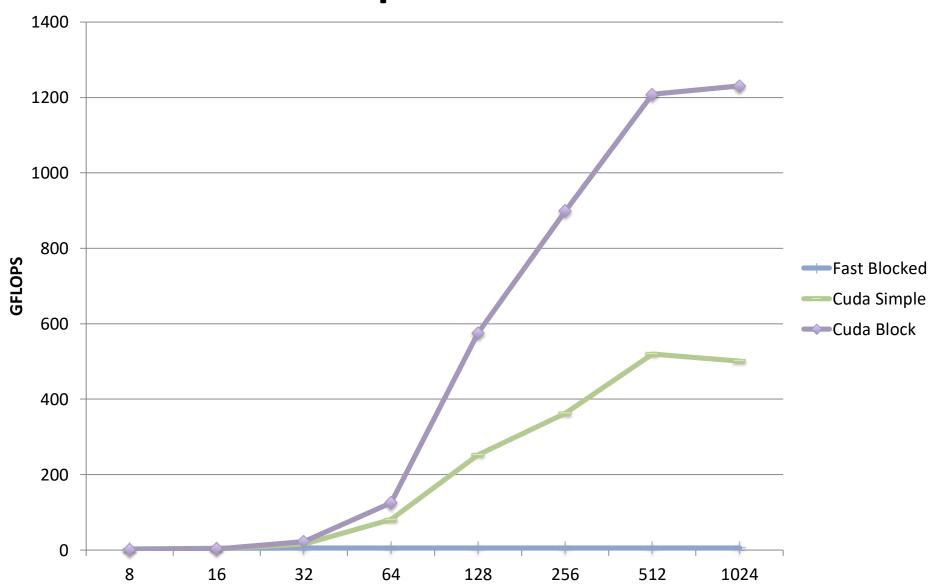
```
for (int bk = 0; bk < LBLK; bk++)
sum += subA[RM(bi,bk,LBLK)] * subB[RM(bk,bj,LBLK)];</pre>
```

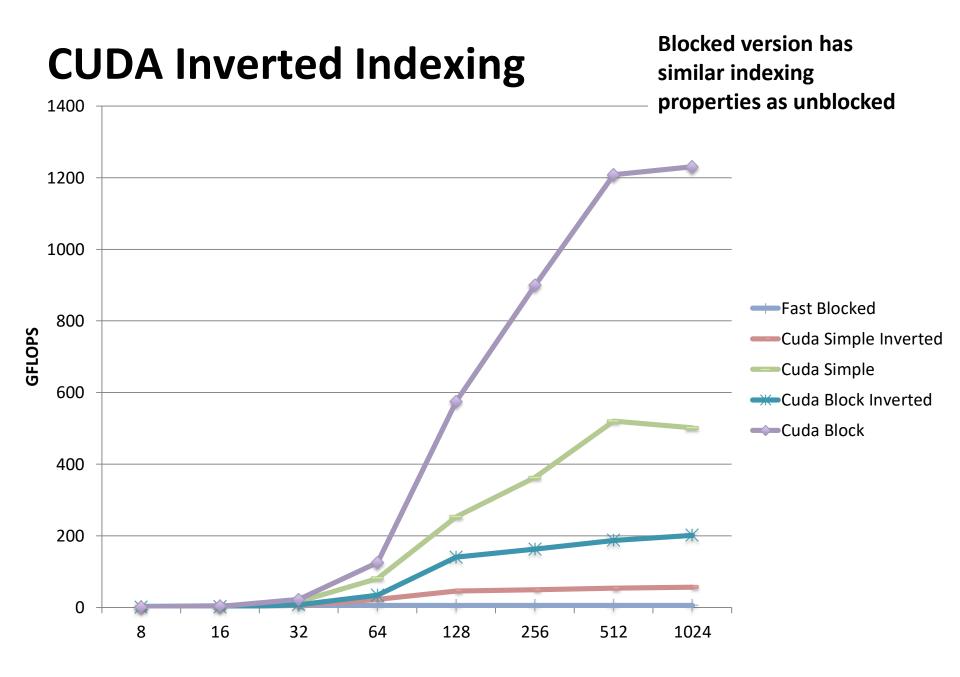
 Each thread in block accumulates single destination value



$$c_{bi,bj} = \mathop{\text{a}}_{bk=0}^{b-1} a_{bi,bk} \times b_{bk,bj}$$

CUDA Blocked Implementations





Warning!

```
for (int k = 0; k < N; k+= LBLK) {
    if (i >= N || j >= N)
        continue; // Skip if out of bounds
    Computation when in-bounds
    // Wait until everyone finished
      syncthreads();
    Compute contribution to element i, j of output
    // Wait until all products computed
      syncthreads();
```

What's wrong with this code?

Observations

Making use of CUDA hierarchy can help

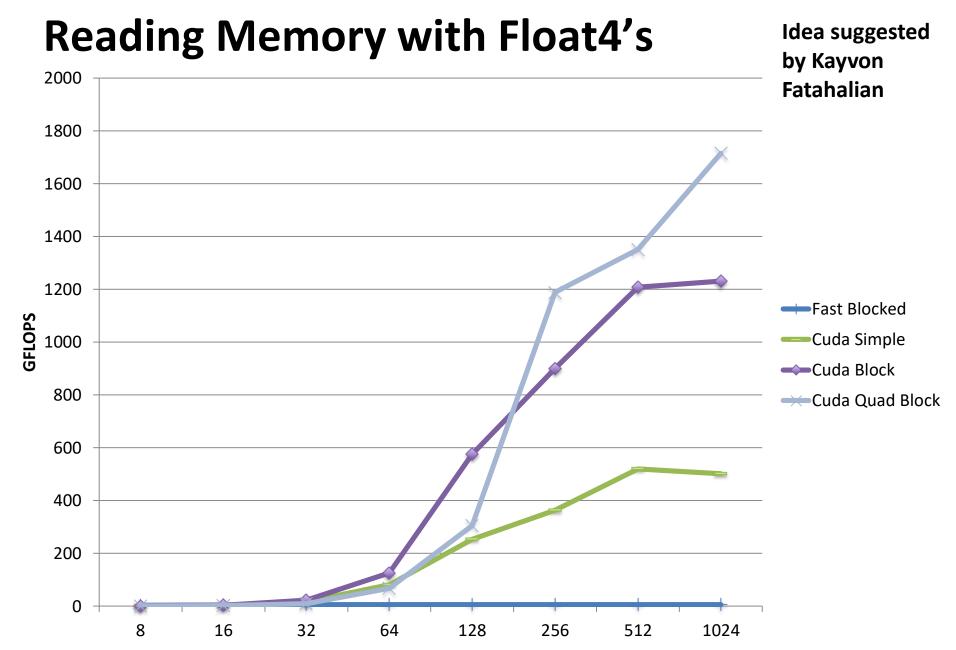
- Lighter weight synchronization
- Shared access to fast memory
- Different blocks can proceed at different rates
 - (Not shown in this example)

Advice

- Implement pure data-parallel version first
- Only exploit hierarchy for performance critical parts
- Watch out for synchronization bugs
- Proper memory referencing more important than these low-level optimizations

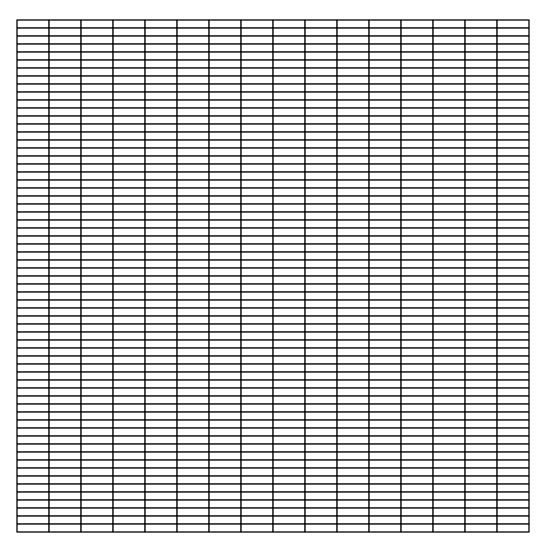






Idea

16 columns of float4's

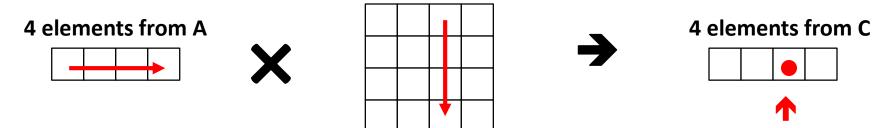


- Thread blocks compute products of 64x64 submatrices
- 1024 threads
- Organize as 64 rows X 16 columns
- Threads read & write memory in chunks of 16 bytes
 - 4 float's each

64 rows

Added Inner Step of Computation

4 X 4 elements from B



Each thread loops 16 times

- Within loop, compute product:
 - 1x4 portion of A
 - 4x4 of B
 - Add sum to 1x4 portion of C
 - 16 multiplies, 16 adds

Why so fast?

Makes maximum use of memory bus capability