# CUDA_C_Programming_Guide

高性能考点总结 王若琪

## 《课件-01-CUDA-C-Basics》

**Page 2 WHAT IS CUDA?**
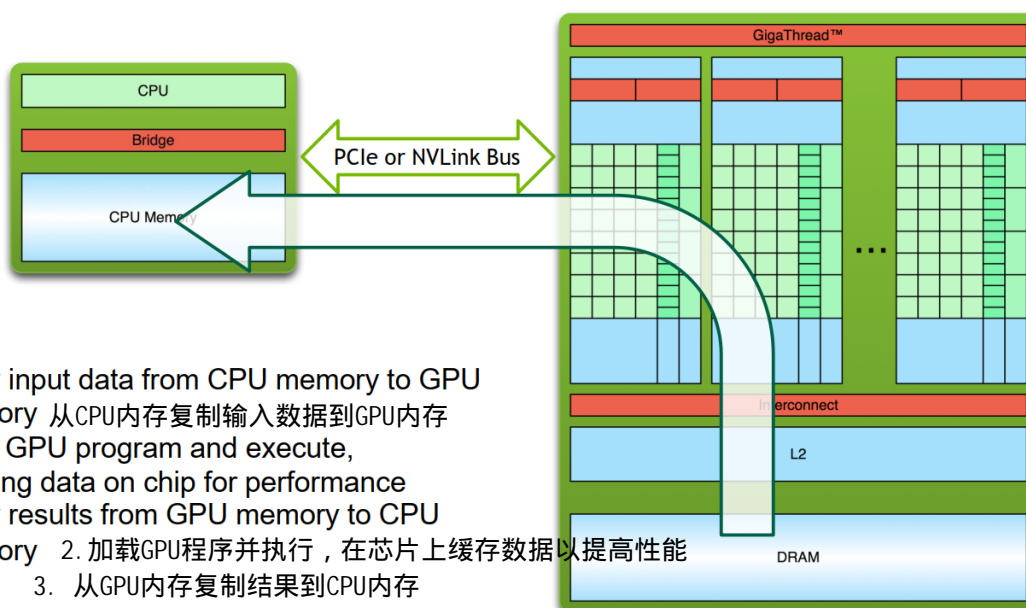
# WHAT IS CUDA?

- ▶ CUDA Architecture    CUDA
    - ▶ Expose GPU parallelism for general-purpose computing
        GPU
    - ▶ Expose/Enable performance

- ▶ CUDA C++
    - ▶ Based on industry-standard C++
        C++
    - ▶ Set of extensions to enable heterogeneous programming

    - ▶ Straightforward APIs to manage devices, memory etc.
        api
- ▶ This session introduces CUDA C++
    CUDA C++
    - ▶ Other languages/bindings available: Fortran, Python, Matlab, etc.
        /    :Fortran, Python, Matlab

**Page 8 SIMPLE PROCESSING FLOW**

# SIMPLE PROCESSING FLOW



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

**Page 10-11 GPU KERNELS: DEVICE CODE**

# GPU KERNELS: DEVICE CODE

```
__global__ void mykernel(void) {

}
```

- CUDA C++ keyword __global__ indicates a function that:
  CUDA c++ keyword global         :
  - Runs on the device

    - Is called from host code (can also be called from other device code)
          (                           )
- **nvcc** separates source code into host and device components
    nvcc
    - Device functions (e.g. **mykernel()**) processed by NVIDIA compiler
              (  mykernel())  NVIDIA
    - Host functions (e.g. **main()**) processed by standard host compiler:
                              (    main()):
        - **gcc, cl.exe**

# GPU KERNELS: DEVICE CODE

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

- Triple angle brackets mark a call to *device* code

    - Also called a "kernel launch"

    - We'll return to the parameters (1,1) in a moment
                              $(1, 1)$
    - The parameters inside the triple angle brackets are the CUDA kernel **execution configuration**
                              CUDA
- That's all that is required to execute a function on the GPU!
            GPU                                    !

**Page 12-16 RUNNING CODE IN PARALLEL,   VECTOR ADDITION ON THE DEVICE**

# MEMORY MANAGEMENT

- Host and device memory are separate entities
- Device pointers point to GPU memory
  - Typically passed to device code
  - Typically not dereferenced in host code
- Host pointers point to CPU memory
  - Typically not passed to device code
  - Typically not dereferenced in device code

- Simple CUDA API for handling device memory
  - cudaMalloc(), cudaFree(), cudaMemcpy()
  - Similar to the C equivalents malloc(), free(), memcpy()

# RUNNING CODE IN PARALLEL

- GPU computing is about massive parallelism

    - So how do we run code in parallel on the device?

$$add<<< 1, 1 >>>();$$

$$add<<< N, 1 >>>();$$

- Instead of executing `add()` once, execute N times in parallel

# 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 all 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
- Built-in variables like `blockIdx.x` are zero-indexed (C/C++ style), 0..**N**-1, where **N** is from the kernel execution configuration indicated at the kernel launch

# VECTOR ADDITION ON THE DEVICE

```c
#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

```c
    // 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);

    // 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;
}
```
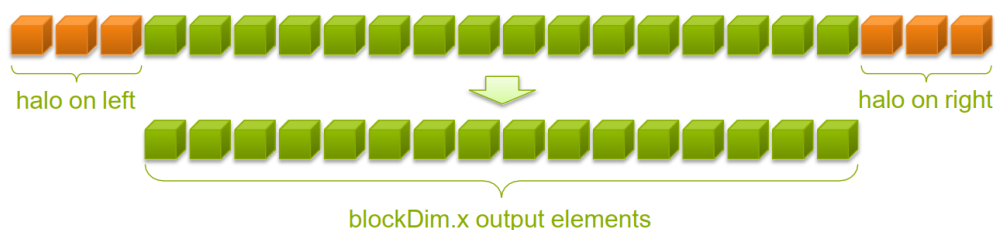
《课件-02-CUDA-Shared-Memory》

**Page 6-13 SHARING DATA BETWEEN THREADS**

# SHARING DATA BETWEEN THREADS

- ▸ Terminology: within a block, threads share data via shared memory
  :

- ▸ Extremely fast on-chip memory, user-managed

- ▸ Declare using `__shared__`, allocated per block

- ▸ Data is not visible to threads in other blocks

# IMPLEMENTING WITH SHARED MEMORY

- ▸ Cache data in shared memory
  - ▸ Read (`blockDim.x + 2 * radius`) input elements from global memory to shared memory
  - ▸ Compute `blockDim.x` output elements
  - ▸ Write `blockDim.x` output elements to global memory
- ▸ Each block needs a halo of `radius` elements at each boundary



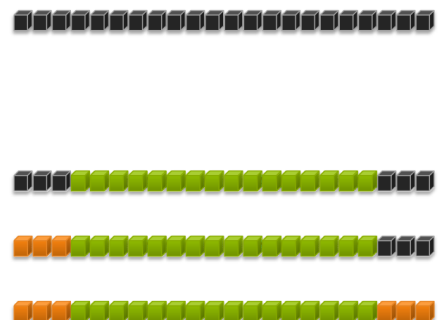halo on left                    halo on right

blockDim.x output elements

# STENCIL KERNEL

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

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

# STENCIL KERNEL

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


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

# DATA RACE!

- The stencil example will not work...

- Suppose thread 15 reads the halo before thread 0 has fetched

```
    temp[lindex] = in[gindex];
      if (threadIdx.x < RADIUS) {          Store at temp[18]
          temp[lindex - RADIUS] = in[gindex - RADIUS];
          temp[lindex + BLOCK_SIZE] = in[gindex + BLOCK_SIZE];
      }                                    Skipped, threadIdx > RADIUS
      int result = 0;
      result += temp[lindex + 1];          Load from temp[19]
```

# __SYNCTHREADS()

- **`void __syncthreads();`**

- Synchronizes all threads within a block

    - Used to prevent RAW / WAR / WAW hazards

- All threads must reach the barrier

    - In conditional code, the condition must be uniform across the block

# STENCIL KERNEL

```
__global void stencil_1d(int *in, int *out) {
    __shared int temp[BLOCK_SIZE + 2 * RADIUS];
    int gindex = threadIdx.x + blockIdx.x * blockDim.x;
    int lindex = threadIdx.x + radius;

    // Read input elements into shared memory
    temp[lindex] = in[gindex];
    if (threadIdx.x < RADIUS) {
        temp[lindex - RADIUS] = in[gindex - RADIUS];
        temp[lindex + BLOCK_SIZE] = in[gindex + BLOCK_SIZE];
    }
    // Synchronize (ensure all the data is available)
    __syncthreads();
```

# STENCIL KERNEL

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

    // Store the result
    out[gindex] = result;
}
```
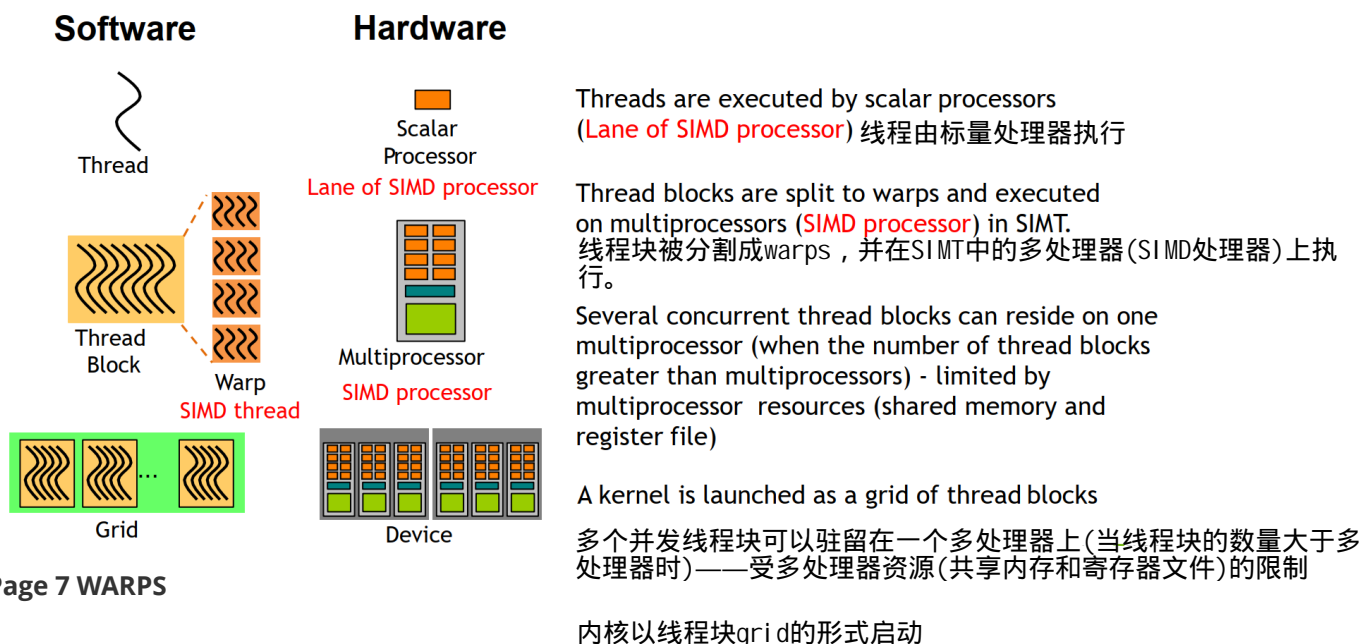
《课件-03-CUDA-Fundamental-Optimization-Part-1》

**Page 6 EXECUTION MODEL**

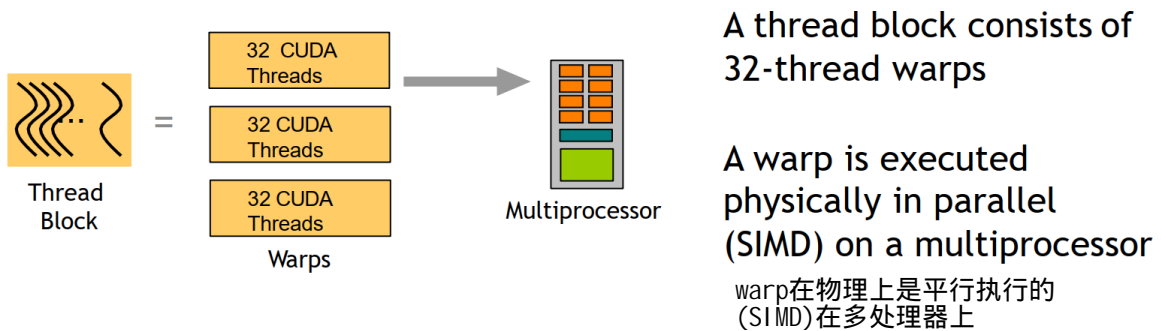# EXECUTION MODEL

**Software**        **Hardware**

Thread

Scalar Processor

Lane of SIMD processor

Threads are executed by scalar processors (Lane of SIMD processor)

Thread Block

Warp

SIMD thread

Multiprocessor

SIMD processor

Thread blocks are split to warps and executed on multiprocessors (SIMD processor) in SIMT.
warps        SIMT        (SIMD        )

Several concurrent thread blocks can reside on one multiprocessor (when the number of thread blocks greater than multiprocessors) - limited by multiprocessor resources (shared memory and register file)

Grid

Device

A kernel is launched as a grid of thread blocks

(        )——        (        )

grid

**Page 7 WARPS**

# WARPS

| | | |
|---|---|---|
| Thread Block | 32 CUDA Threads | Multiprocessor |
| | 32 CUDA Threads | |
| | 32 CUDA Threads | |
| | Warps | |

A thread block consists of 32-thread warps

A warp is executed physically in parallel (SIMD) on a multiprocessor

warp
(SIMD)

《课件-04-CUDA-Fundamental-Optimization-Part-2》

**Page 8-17 GPU MEM OPERATIONS**

# GPU MEM OPERATIONS
GPU MEM

- Loads:

  - Caching

    - Default mode

    - Attempts to hit in L1, then L2, then GMEM          L1          L2          GMEM

    - Load granularity is 128-byte line          128

- Stores:

  - Invalidate L1, write-back for L2 L1          L2

# GPU MEM OPERATIONS

- Loads:

  - Non-caching          cache

    - Compile with -Xptxas -dlcm=cg option to nvcc
      -Xptxas -dlcm =cg                    nvcc
    - Attempts to hit in L2, then GMEM
          L2                    GMEM
      Do not hit in L1, invalidate the line if it's in L1 already
          L1                    L1
    - Load granularity is 32-bytes (segment)
          32     ( )

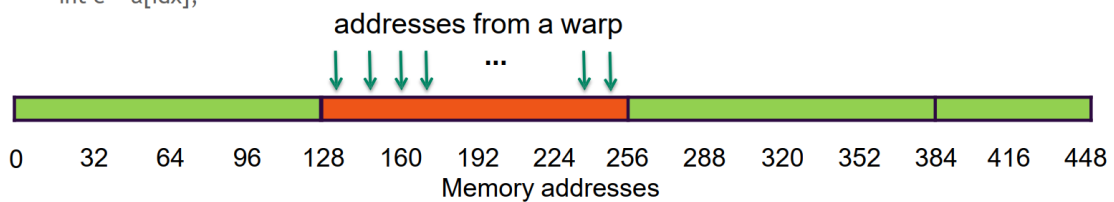We won't spend much time with non-caching loads in this training session

# LOAD OPERATION

- ▸ Memory operations are issued per warp (32 threads)
  warp(32     )
  - ▸ Just like all other instructions

- ▸ Operation:
  - ▸ Threads in a warp provide memory addresses
    wrap
  - ▸ Determine which lines/segments are needed
    /
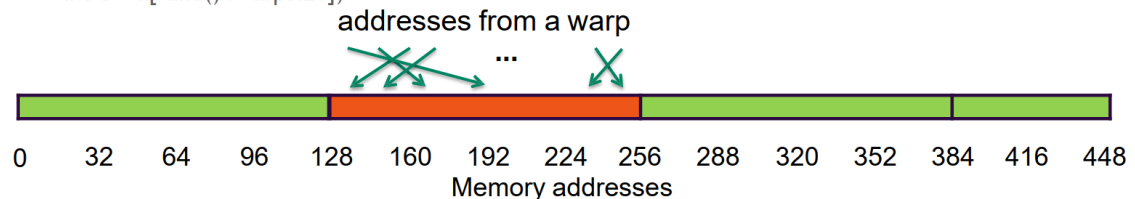  - ▸ Request the needed lines/segments
    /

# CACHING LOAD

- ▸ Warp requests 32 aligned, consecutive 4-byte words
  wrap     32              4
- ▸ Addresses fall within 1 cache-line

  - ▸ Warp needs 128 bytes

  - ▸ 128 bytes move across the bus on a miss
    128
  - ▸ Bus utilization: 100%
                : 100%

  - ▸ int c = a[idx];

addresses from a warp

...

| 0 | 32 | 64 | 96 | 128 | 160 | 192 | 224 | 256 | 288 | 320 | 352 | 384 | 416 | 448 |

Memory addresses

# CACHING LOAD

- ▸ Warp requests 32 aligned, permuted 4-byte words

- ▸ Addresses fall within 1 cache-line

  - ▸ Warp needs 128 bytes

  - ▸ 128 bytes move across the bus on a miss

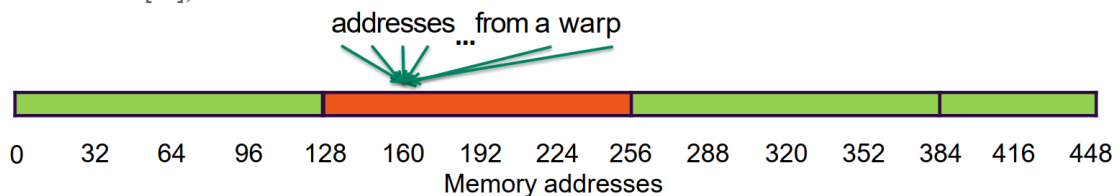  - ▸ Bus utilization: 100%

  - ▸ int c = a[rand()%warpSize];

addresses from a warp

...

| 0 | 32 | 64 | 96 | 128 | 160 | 192 | 224 | 256 | 288 | 320 | 352 | 384 | 416 | 448 |

Memory addresses

# CACHING LOAD

- Warp requests 32 misaligned, consecutive 4-byte words
- Addresses fall within 2 cache-lines
  - Warp needs 128 bytes
  - 256 bytes move across the bus on misses
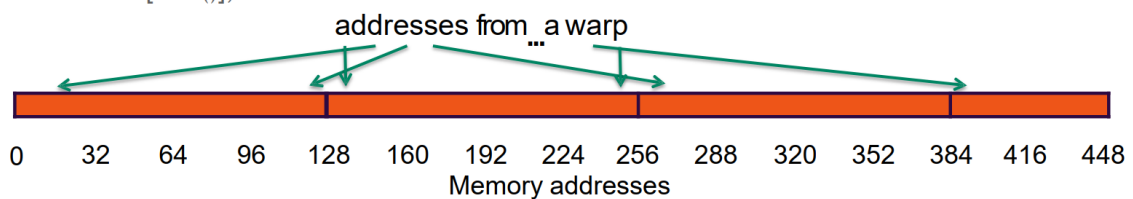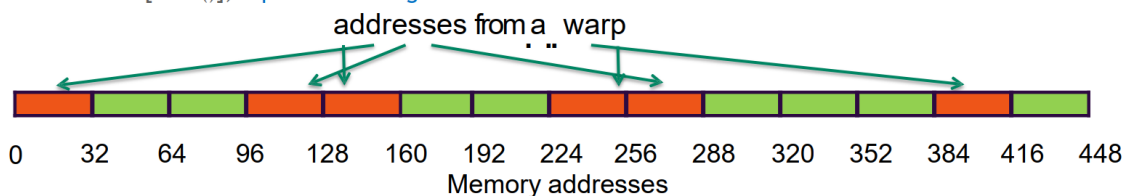  - Bus utilization: 50%
  - int c = a[idx-2];

addresses from a warp
...

| 0 | 32 | 64 | 96 | 128 | 160 | 192 | 224 | 256 | 288 | 320 | 352 | 384 | 416 | 448 |

Memory addresses

# CACHING LOAD

- All threads in a warp request the same 4-byte word
- Addresses fall within a single cache-line
  - Warp needs 4 bytes
  - 128 bytes move across the bus on a miss
  - Bus utilization: 3.125%
  - int c = a[40];

addresses  from a warp
...

| 0 | 32 | 64 | 96 | 128 | 160 | 192 | 224 | 256 | 288 | 320 | 352 | 384 | 416 | 448 |

Memory addresses

# CACHING LOAD

- Warp requests 32 scattered 4-byte words
- Addresses fall within N cache-lines
  - Warp needs 128 bytes
  - N*128 bytes move across the bus on a miss
  - Bus utilization: 128 / (N*128) (3.125% worst case N=32)
  - int c = a[rand()];

addresses from  a warp
...

| 0 | 32 | 64 | 96 | 128 | 160 | 192 | 224 | 256 | 288 | 320 | 352 | 384 | 416 | 448 |

Memory addresses

# NON-CACHING LOAD

▶ Warp requests 32 scattered 4-byte words

▶ Addresses fall within N segments

 ▶ Warp needs 128 bytes

 ▶ N*32 bytes move across the bus on a miss

 ▶ Bus utilization: 128 / (N*32) (12.5% worst case N = 32)

 ▶ int c = a[rand()]; –Xptxas –dlcm=cg

addresses from a  warp

| | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|

0    32    64    96    128    160    192    224    256    288    320    352    384    416    448

**Memory addresses**

# GPU MEM OPTIMIZATION GUIDELINES
GPU  MEM

▶ Strive for perfect coalescing

 ▶ (Align starting address - may require padding)
   –

 ▶ A warp should access within a contiguous region
   wrap

▶ Have enough concurrent accesses to saturate the bus

 ▶ Process several elements per thread

  ▶ Multiple loads get pipelined
    l oad
  ▶ Indexing calculations can often be reused

 ▶ Launch enough warps to maximize throughput
   wraps
  ▶ Latency is hidden by switching warps
    wrap

▶ Use all the caches!
   !

**Page 20-25 SHARED MEMORY**

# SHARED MEMORY

 ▶ Uses:

  ▶ Inter-thread communication within a block

  ▶ Cache data to reduce redundant global memory accesses

  ▶ Use it to improve global memory access patterns

 ▶ Organization:

  ▶ 32 banks, 4-byte wide banks
    32        4
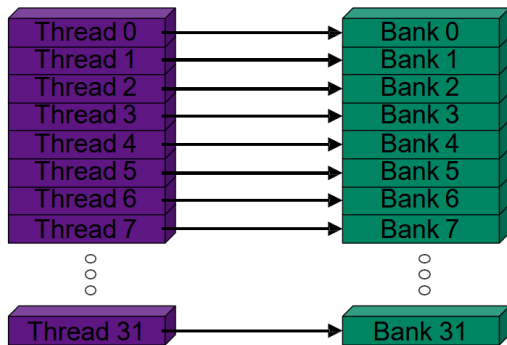  ▶ Successive 4-byte words belong to different banks
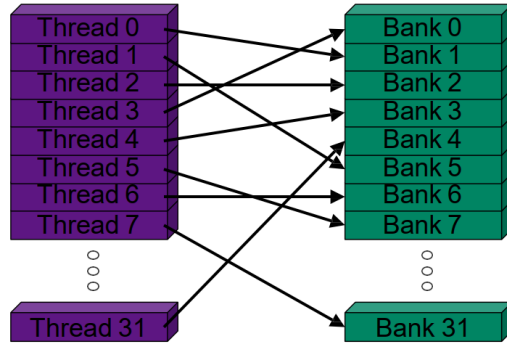    4

# SHARED MEMORY

- Performance:
  - Typically: 4 bytes per bank per 1 or 2 clocks per multiprocessor
  - shared accesses are issued per 32 threads (warp)
    32        wrap
  - serialization: if $N$ threads of 32 access different 4-byte words in the same bank, $N$ accesses are executed serially
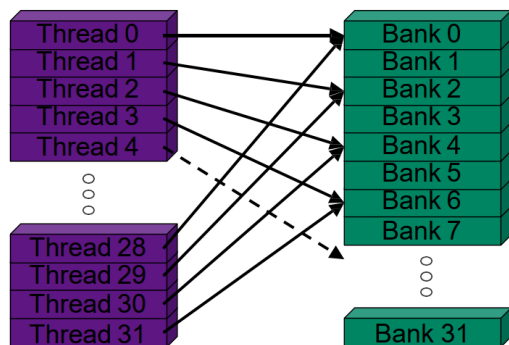    : N (32 ) bank 4 N

# BANK ADDRESSING EXAMPLES

No Bank Conflicts

No Bank Conflicts
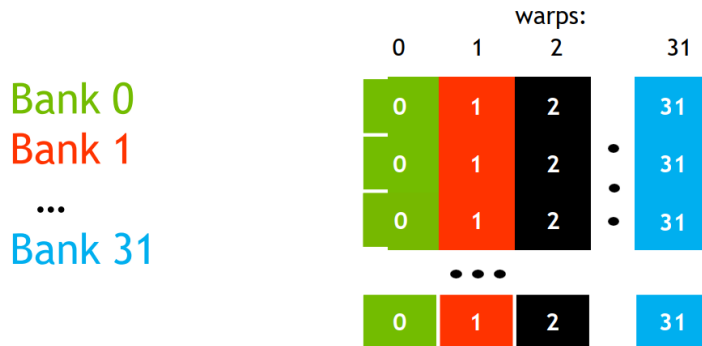
# BANK ADDRESSING EXAMPLES

2-way Bank Conflicts

16-way Bank Conflicts

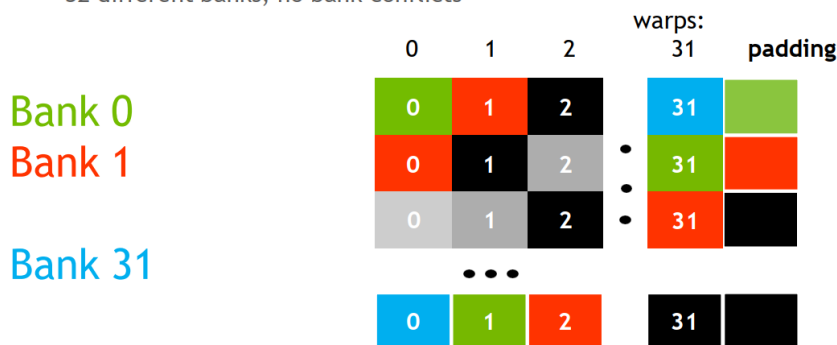# SHARED MEMORY: AVOIDING BANK CONFLICTS

bank

- ▸ **32x32** Shared MEM array

- ▸ Warp accesses a column:

  - ▸ 32-way bank conflicts (threads in a warp access the same bank)



# SHARED MEMORY: AVOIDING BANK CONFLICTS

- ▸ Add a column for padding:

  - ▸ 32x33 SMEM array

- ▸ Warp accesses a column:

  - ▸ 32 different banks, no bank conflicts
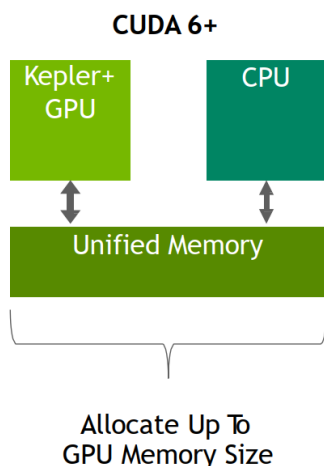


---

《课件-05_Atomics_Reductions_Warp_Shuffle》

**Page12-33 parallel reduction optimization**

---

《课件-06_Managed_Memory》

Page 5-9 UNIFIED MEMORY

# UNIFIED MEMORY

## Reduce Developer Effort

**CUDA 6+**

| Kepler+ GPU | CPU |
|---|---|

Unified Memory

Allocate Up To
GPU Memory Size

**Simpler Programming & Memory Model**
- Single allocation, single pointer, accessible anywhere
- Eliminate need for *explicit* copy
- Simplifies code porting

**Maintain Performance through Data Locality**
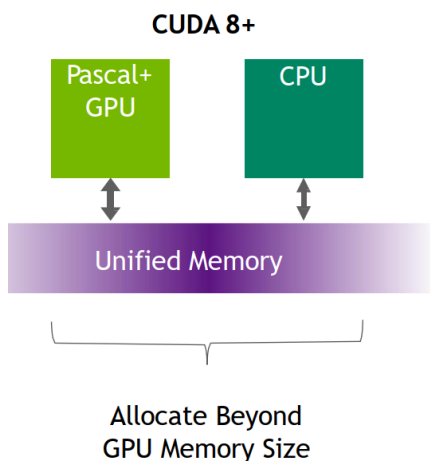- Migrate data to accessing processor
- Guarantee global coherence
- Still allows explicit hand tuning

# CUDA 8+: UNIFIED MEMORY

## Demand Paging For Pascal and Beyond

Pascal

**CUDA 8+**

| Pascal+ GPU | CPU |
|---|---|

Unified Memory

Allocate Beyond
GPU Memory Size

**Enable Large Data Models**
enable
- GPU
- Oversubscribe GPU memory
- Allocate up to system memory size

**Simpler Data Accesss**
- CPU/GPU Data coherence   CPU / GPU
- Unified memory atomic operations

**Tune Unified Memory Performance**
- cudaMemAdvise API
- Usage hints via cudaMemAdvise API
- Explicit prefetching API
- API

# SIMPLIFIED MEMORY MANAGEMENT CODE

## CPU Code

```
void sortfile(FILE *fp, int N) {
  char *data;
  data = (char *)malloc(N);

  fread(data, 1, N, fp);

  qsort(data, N, 1, compare);


  use_data(data);

  free(data);
}
```

## Ordinary CUDA Code

```
void sortfile(FILE *fp, int N) {
  char *data, *d_data;
  data = (char *)malloc(N);
           (&d_data, N);
  fread(data, 1, N, fp);
  cudaMemcpy(d_data, data, N, …); // 1
  qsort<<<...>>>(data,N,1,compare); // 2
  cudaMemcpy(data, d_data, N, …); // 3

  use_data(data);
  cudaFree(d_data);
  free(data);
}
```

# SIMPLIFIED MEMORY MANAGEMENT CODE

**CPU Code**

```
void sortfile(FILE *fp, int N) {
  char *data;
  data = (char *)malloc(N);

  fread(data, 1, N, fp);

  qsort(data, N, 1, compare);


  use_data(data);

  free(data);
}
```

**CUDA Code with Unified Memory**

```
void sortfile(FILE *fp, int N) {
  char *data;
  cudaMallocManaged(&data, N);

  fread(data, 1, N, fp);

  qsort<<<...>>>(data,N,1,compare);
  cudaDeviceSynchronize();

  use_data(data);

  cudaFree(data);
}
```

# UNIFIED MEMORY EXAMPLE

## With On-Demand Paging

```
__global__
void setValue(int *ptr, int index,  int val)
{
  ptr[index] = val;
}


void foo(int size) {
  char *data;
  cudaMallocManaged(&data,  size);        ←——————————  Unified Memory allocation

  memset(data, 0, size);                  ←——————————  Access all values on CPU

  setValue<<<...>>>(data, size/2, 5);     ←——————————  Access one value on GPU
  cudaDeviceSynchronize();

  useData(data);

  cudaFree(data);
}
```