

# MODERN HARDWARE MODELS — HOW TO EXPLOIT SPECIFIC FEATURES WITH CUDA

Siegfried Höfingier

VSC Research Center, TU Wien

March 4, 2020

→ <https://tinyurl.com/cuda4dummies/i/12/notes-12.pdf>

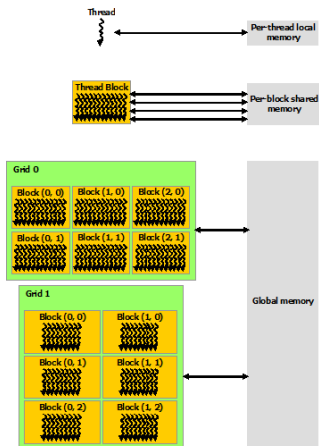
MEMORY HIERARCHY

TAKE HOME MESSAGES

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, PROGRAMMING MODEL CONT.

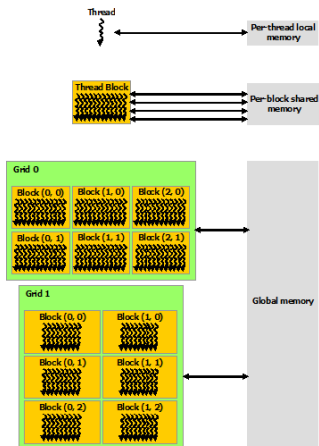
- Threads have access to several memory spaces



→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, PROGRAMMING MODEL CONT.

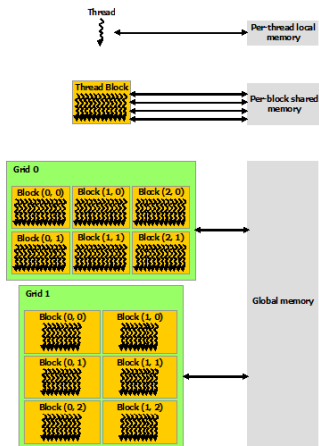


- Threads have access to several memory spaces
- Private local memory for each individual thread (fast)

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, PROGRAMMING MODEL CONT.

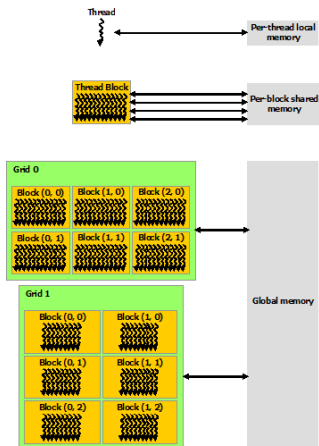


- Threads have access to several memory spaces
- Private local memory for each individual thread (fast)
- Shared memory for all threads within a thread block (pretty fast)

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, PROGRAMMING MODEL CONT.

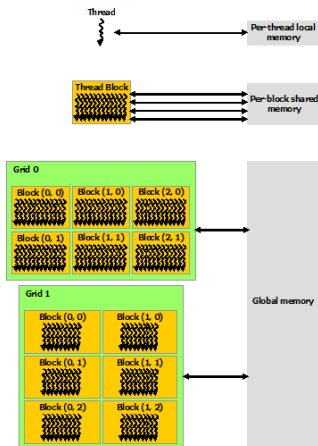


- Threads have access to several memory spaces
- Private local memory for each individual thread (fast)
- Shared memory for all threads within a thread block (pretty fast)
- Global memory accessible to all threads (slow)

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, PROGRAMMING MODEL CONT.

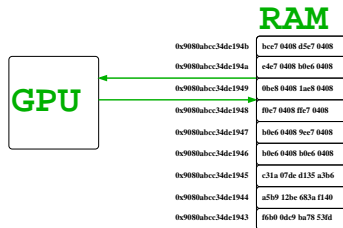


- Threads have access to several memory spaces
- Private local memory for each individual thread (fast)
- Shared memory for all threads within a thread block (pretty fast)
- Global memory accessible to all threads (slow)
- Special ROMs, constant and texture memory

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING



### Unified Memory:

GPUs of compute  
capability  $\geq 6.x$  and  $\geq$  CUDA 8.0

### Separate Device/Host Memory:

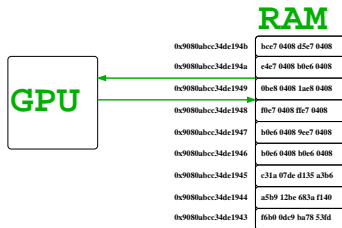
GPUs of compute  
capability  $< 6.x$  and  $<$  CUDA 8.0

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>



# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING



Unified Memory:

GPUs of compute capability  $\geq 6.x$  and  $\geq$  CUDA 8.0

Separate Device/Host Memory:

GPUs of compute capability  $< 6.x$  and  $<$  CUDA 8.0

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- Unified memory (managed) for recent GPUs (Pascal class or newer) greatly simplifies programming4/porting2 the GPU

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- Unified memory (managed) for recent GPUs (Pascal class or newer) greatly simplifies programming4/porting2 the GPU
- Single common address space accessible from any processor in a system

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- Unified memory (managed) for recent GPUs (Pascal class or newer) greatly simplifies programming/porting to the GPU
- Single common address space accessible from any processor in a system
- Strict distinction into device- and host-memory becoming less important to the programmer

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- Unified memory (managed) for recent GPUs (Pascal class or newer) greatly simplifies programming/porting to the GPU
- Single common address space accessible from any processor in a system
- Strict distinction into device- and host-memory becoming less important to the programmer
- Single-pointer-to-data model, no longer any need to copy forth/back data into/from GPU memory

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- Unified memory (managed) for recent GPUs (Pascal class or newer) greatly simplifies programming/porting to the GPU
- Single common address space accessible from any processor in a system
- Strict distinction into device- and host-memory becoming less important to the programmer
- Single-pointer-to-data model, no longer any need to copy forth/back data into/from GPU memory
- On-demand page migration (hardware supported) — Page Migration Engine

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- Unified memory (managed) for recent GPUs (Pascal class or newer) greatly simplifies programming<sup>4</sup>/porting<sup>2</sup> the GPU
- Single common address space accessible from any processor in a system
- Strict distinction into device- and host-memory becoming less important to the programmer
- Single-pointer-to-data model, no longer any need to copy forth/back data into/from GPU memory
- On-demand page migration (hardware supported) — Page Migration Engine
- Simply invoked via `cudaMallocManaged()`

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

### Multiple Thread Blocks Matrix Addition

```
__global__ void MatAdd( float **A, float **B, float **C )
{
    int i, j;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    j = (blockIdx.y * blockDim.y) + threadIdx.y;
    C[i][j] = A[i][j] + B[i][j];
}

int main()
{
    int i;
    dim3 threadsPerBlock, numBlocks;
    float **A, **B, **C;
    ...
    // unified memory allocation, B and C analogous
    cudaMallocManaged(&A, N * sizeof(float *));
    for (i = 0; i < N; i++) {
        cudaMallocManaged(&A[i], N * sizeof(float));
    }
    // kernel invocation
    MatAdd <<< numBlocks, threadsPerBlock >>> (A, B, C);
    cudaDeviceSynchronize();
    ...
    cudaFree(A);
}
```

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>



# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

### Multiple Thread Blocks Matrix Addition

```
__global__ void MatAdd( float **A, float **B, float **C )
{
    int i, j;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    j = (blockIdx.y * blockDim.y) + threadIdx.y;
    C[i][j] = A[i][j] + B[i][j];
}

int main()
{
    int i;
    dim3 threadsPerBlock, numBlocks;
    float **A, **B, **C;
    ...
    // unified memory allocation, B and C analogous
    cudaMallocManaged(&A, N * sizeof(float *));
    for (i = 0; i < N; i++) {
        cudaMallocManaged(&A[i], N * sizeof(float));
    }
    // kernel invocation
    MatAdd <<< numBlocks, threadsPerBlock >>> (A, B, C);
    cudaDeviceSynchronize();
    ...
    cudaFree(A);
}
```

Standard  
malloc-like  
allocation  
on the host

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

## Multiple Thread Blocks Matrix Addition

```
__global__ void MatAdd( float **A, float **B, float **C )
{
    int i, j;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    j = (blockIdx.y * blockDim.y) + threadIdx.y;
    C[i][j] = A[i][j] + B[i][j];
}

int main()
{
    int i;
    dim3 threadsPerBlock, numBlocks;
    float **A, **B, **C;
    ...
    // unified memory allocation, B and C analogous
    cudaMallocManaged(&A, N * sizeof(float *));
    for (i = 0; i < N; i++) {
        cudaMallocManaged(&A[i], N * sizeof(float));
    }
    // kernel invocation
    MatAdd <<< numBlocks, threadsPerBlock >>> (A, B, C);
    cudaDeviceSynchronize();
    ...
    cudaFree(A);
}
```

Standard  
malloc-like  
allocation  
on the host

host pointers directly  
usable in  
kernel code

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

## Multiple Thread Blocks Matrix Addition

```
__global__ void MatAdd( float **A, float **B, float **C )
{
    int i, j;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    j = (blockIdx.y * blockDim.y) + threadIdx.y;
    C[i][j] = A[i][j] + B[i][j];
}

int main()
{
    int i;
    dim3 threadsPerBlock, numBlocks;
    float **A, **B, **C;
    ...
    // unified memory allocation, B and C analogous
    cudaMallocManaged(&A, N * sizeof(float *));
    for (i = 0; i < N; i++) {
        cudaMallocManaged(&A[i], N * sizeof(float));
    }
    // kernel invocation
    MatAdd <<< numBlocks, threadsPerBlock >>> (A, B, C);
    cudaDeviceSynchronize();
    ...
    cudaFree(A);
}
```

Standard  
malloc-like  
allocation  
on the host

host pointers  
directly  
usable in  
kernel code

ensure  
proper ker-  
nel comple-  
tion

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

## Multiple Thread Blocks Matrix Addition

```
__global__ void MatAdd( float **A, float **B, float **C )
{
    int i, j;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    j = (blockIdx.y * blockDim.y) + threadIdx.y;
    C[i][j] = A[i][j] + B[i][j];
}

int main()
{
    int i;
    dim3 threadsPerBlock, numBlocks;
    float **A, **B, **C;
    ...
    // unified memory allocation, B and C analogous
    cudaMallocManaged(&A, N * sizeof(float *));
    for (i = 0; i < N; i++) {
        cudaMallocManaged(&A[i], N * sizeof(float));
    }
    // kernel invocation
    MatAdd <<< numBlocks, threadsPerBlock >>> (A, B, C);
    cudaDeviceSynchronize();
    ...
    cudaFree(A);
}
```

Standard  
malloc-like  
allocation  
on the host

free memory  
when done

host pointers  
directly  
usable in  
kernel code

ensure  
proper ker-  
nel comple-  
tion

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

### 2nd Form — Managed Global Memory

```
#define ARRAYDIM 268435456

// global managed declaration on GPU
__device__ __managed__ float x[ARRAYDIM], y[ARRAYDIM], z[ARRAYDIM];

__global__ void KrnlDmmy()
{
    int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    x[i] = y[i] + z[i];
    return;
}

int main()
{
    ...
    // kernel invocation
    KrnlDmmy <<<< numBlocks, threadsPerBlock >>>> ();
    cudaDeviceSynchronize();
    ...
}
```

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

## 2nd Form — Managed Global Memory

```
#define ARRAYDIM 268435456

// global managed declaration on GPU
__device__ __managed__ float x[ARRAYDIM], y[ARRAYDIM], z[ARRAYDIM];

__global__ void KrnlDmmy()
{
    int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    x[i] = y[i] + z[i];
    return;
}

int main()
{
    ...
    // kernel invocation
    KrnlDmmy <<<< numBlocks, threadsPerBlock >>>> ();
    cudaDeviceSynchronize();
    ...
}
```

Global variables directly usable on device & host

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

## 2nd Form — Managed Global Memory

```
#define ARRAYDIM 268435456

// global managed declaration on GPU
__device__ __managed__ float x[ARRAYDIM], y[ARRAYDIM], z[ARRAYDIM];

__global__ void KrnlDmmy()
{
    int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    x[i] = y[i] + z[i];
    return;
}

int main()
{
    ...
    // kernel invocation
    KrnlDmmy <<< numBlocks, threadsPerBlock >>> ();
    cudaDeviceSynchronize();
    ...
}
```

Global variables directly usable on device & host

void kernel call

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

## 2nd Form — Managed Global Memory

```
#define ARRAYDIM 268435456

// global managed declaration on GPU
__device__ __managed__ float x[ARRAYDIM], y[ARRAYDIM], z[ARRAYDIM];

__global__ void KrnlDmmy()
{
    int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    x[i] = y[i] + z[i];
    return;
}

int main()
{
    ...
    // kernel invocation
    KrnlDmmy <<< numBlocks, threadsPerBlock >>> ();
    cudaDeviceSynchronize();
    ...
}
```

Global variables directly usable on device & host

void kernel call

ensure proper kernel completion

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>



### Unified Memory Example Version 1

```
#define ARRAYDIM 268435456

__global__ void KrnlDmmy(float *x, float *y, float *z)
{
    int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    x[i] = (float) i;
    y[i] = (float) (i + 1);
    z[i] = x[i] + y[i];
    return;
}

int main()
{
    ...
    // unified memory allocation, b and c analogous
    cudaMallocManaged(&a, ARRAYDIM * sizeof(float));
    // kernel invocation
    KrnlDmmy <<<< numBlocks, threadsPerBlock >>>> (a, b, c);
    cudaDeviceSynchronize();
    ...
}
```

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

## Unified Memory Example Version 1

```
#define ARRAYDIM 268435456

__global__ void KrnlDmmy(float *x, float *y, float *z)
{
    int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    x[i] = (float) i;
    y[i] = (float) (i + 1);
    z[i] = x[i] + y[i];
    return;
}

int main()
{
    ...
    // unified memory allocation, b and c analogous
    cudaMallocManaged(&a, ARRAYDIM * sizeof(float));
    // kernel invocation
    KrnlDmmy <<<< numBlocks, threadsPerBlock >>>> (a, b, c);
    cudaDeviceSynchronize();
    ...
}
```

Kernel  
initializa-  
tion and  
calcula-  
tion

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

## Unified Memory Example Version 1

```
#define ARRAYDIM 268435456

__global__ void KrnlDmmy(float *x, float *y, float *z)
{
    int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    x[i] = (float) i;
    y[i] = (float) (i + 1);
    z[i] = x[i] + y[i];
    return;
}

int main()
{
    ...
    // unified memory allocation, b and c analogous
    cudaMallocManaged(&a, ARRAYDIM * sizeof(float));
    // kernel invocation
    KrnlDmmy <<<< numBlocks, threadsPerBlock >>>> (a, b, c);
    cudaDeviceSynchronize();
    ...
}
```

Kernel  
initializa-  
tion and  
calcula-  
tion

1 GB  
arrays

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

```
[sh@n71-006]$ nvcc ./unified_memory_example_1.cu
[sh@n71-006]$ ./a.out
[sh@n71-006]$ nvprof ./a.out

==126896== NVPROF is profiling process 126896, command: ./a.out
==126896== Profiling application: ./a.out
==126896== Profiling result:
   Type  Time(%)   Time     Calls   Avg       Min       Max  Name
GPU activities: 100.00% 558.88ms   1 558.88ms 558.88ms 558.88ms KrnlDmmy(float*, float*, flo
  API calls: 66.84% 558.90ms   1 558.90ms 558.90ms 558.90ms cudaDeviceSynchronize
            20.99% 175.51ms   3 58.504ms 16.390us 175.45ms cudaMallocManaged
            12.04% 100.67ms   3 33.557ms 33.476ms 33.610ms cudaFree
            0.07% 616.14us  94 6.5540us 600ns 248.06us cuDeviceGetAttribute
            0.04% 364.67us   1 364.67us 364.67us 364.67us cuDeviceTotalMem
            0.01% 71.540us   1 71.540us 71.540us 71.540us cudaLaunch
            0.01% 69.450us   1 69.450us 69.450us 69.450us cuDeviceGetName
            0.00% 5.7300us   3 1.9100us 250ns 5.0900us cudaSetupArgument
            0.00% 4.5400us   3 1.5130us 630ns 2.8700us cuDeviceGetCount
            0.00% 2.3300us   2 1.1650us 650ns 1.6800us cuDeviceGet
            0.00% 1.2900us   1 1.2900us 1.2900us 1.2900us cudaConfigureCall

==126896== Unified Memory profiling result:
Device "GeForce GTX 1080 (0)"
  Count  Avg Size  Min Size  Max Size  Total Size  Total Time  Name
    5275      -      -      -      -      548.8451ms  Gpu page fault groups
```

→ [https://tinyurl.com/cuda4dummies/i/12/unified\\_memory\\_example\\_1.cu](https://tinyurl.com/cuda4dummies/i/12/unified_memory_example_1.cu)

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

```
[sh@n71-006]$ nvcc ./unified_memory_example_1.cu
[sh@n71-006]$ ./a.out
[sh@n71-006]$ nvprof ./a.out

==126896== NVPROF is profiling process 126896, command: ./a.out
==126896== Profiling application: ./a.out
==126896== Profiling result:
   Type  Time(%)   Time     Calls   Avg       Min       Max   Name
GPU activities: 100.00% 558.88ms | 1 558.88ms 558.88ms 558.88ms | KrnlDmmy(float*, float*, flo
  API calls: 66.84% 558.90ms | 1 558.90ms 558.90ms 558.90ms | cudaDeviceSynchronize
              20.99% 175.51ms | 3 58.504ms 16.390us 175.45ms | cudaMallocManaged
              12.04% 100.67ms | 3 33.557ms 33.476ms 33.610ms | cudaFree
              0.07% 616.14us | 94 6.5540us 600ns 248.06us | cuDeviceGetAttribute
              0.04% 364.67us | 1 364.67us 364.67us 364.67us | cuDeviceTotalMem
              0.01% 71.540us | 1 71.540us 71.540us 71.540us | cudaLaunch
              0.01% 69.450us | 1 69.450us 69.450us 69.450us | cuDeviceGetName
              0.00% 5.7300us | 3 1.9100us 250ns 5.0900us | cudaSetupArgument
              0.00% 4.5400us | 3 1.5130us 630ns 2.8700us | cuDeviceGetCount
              0.00% 2.3300us | 2 1.1650us 650ns 1.6800us | cuDeviceGet
              0.00% 1.2900us | 1 1.2900us 1.2900us 1.2900us | cudaConfigureCall

==126896== Unified Memory profiling result:
Device "GeForce GTX 1080 (0)"
  Count  Avg Size  Min Size  Max Size  Total Size  Total Time  Name
  5275   -        -         -         -         548.8451ms  Gpu page fault groups
-----
```

→ [https://tinyurl.com/cuda4dummies/i/12/unified\\_memory\\_example\\_1.cu](https://tinyurl.com/cuda4dummies/i/12/unified_memory_example_1.cu)

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- While straightforward profiling with `nvprof` is nice (especially in terms of unified memory analysis) the results showed a significant number of page faults

→ <https://devblogs.nvidia.com/unified-memory-cuda-beginners>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- While straightforward profiling with `nvprof` is nice (especially in terms of unified memory analysis) the results showed a significant number of page faults
- In a real application, the GPU is likely to use the data more often and for more complex computations, so the overhead from page faulting may become negligible; to quantify, we...

→ <https://devblogs.nvidia.com/unified-memory-cuda-beginners>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- While straightforward profiling with `nvprof` is nice (especially in terms of unified memory analysis) the results showed a significant number of page faults
  - In a real application, the GPU is likely to use the data more often and for more complex computations, so the overhead from page faulting may become negligible; to quantify, we...
1. Could separate the actual calculation from the initialization with the help of a second kernel

→ <https://devblogs.nvidia.com/unified-memory-cuda-beginners>



# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- While straightforward profiling with `nvprof` is nice (especially in terms of unified memory analysis) the results showed a significant number of page faults
  - In a real application, the GPU is likely to use the data more often and for more complex computations, so the overhead from page faulting may become negligible; to quantify, we...
1. Could separate the actual calculation from the initialization with the help of a second kernel
  2. Could repeat the kernel doing the calculation many times

→ <https://devblogs.nvidia.com/unified-memory-cuda-beginners>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- While straightforward profiling with `nvprof` is nice (especially in terms of unified memory analysis) the results showed a significant number of page faults
  - In a real application, the GPU is likely to use the data more often and for more complex computations, so the overhead from page faulting may become negligible; to quantify, we...
1. Could separate the actual calculation from the initialization with the help of a second kernel
  2. Could repeat the kernel doing the calculation many times
  3. Could use unified memory prefetching to explicitly move the data to the GPU

→ <https://devblogs.nvidia.com/unified-memory-cuda-beginners>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

```
[sh@n71-006]$ nvcc ./unified_memory_example_2.cu
[sh@n71-006]$ nvprof ./a.out

==3316== NVPROF is profiling process 3316, command: ./a.out
==3316== Profiling application: ./a.out
==3316== Profiling result:
   Type  Time(%)   Time         Calls      Avg      Min       Max   Name
GPU activities:  63.49%  341.33ms |      1   341.33ms  341.33ms  341.33ms | KrnlDmmyInit(float*, float*,
                36.51%  196.25ms |      1   196.25ms  196.25ms  196.25ms | KrnlDmmyCalc(float*, float*,
API calls:      62.34%  537.60ms |      2   268.80ms  196.26ms  341.34ms | cudaDeviceSynchronize
                25.62%  220.97ms |      3    73.657ms  16.760us  220.91ms | cudaMallocManaged
                11.75%  101.31ms |      3    33.770ms  33.658ms  33.845ms | cudaFree
                0.11%   950.09us |      2    475.05us  75.991us  874.10us | cudaLaunch
                0.07%   607.36us |     94    6.4610us   540ns   244.02us | cuDeviceGetAttribute
                0.07%   578.76us |      1    578.76us  578.76us  578.76us | cuDeviceGetName
                0.04%   386.76us |      1    386.76us  386.76us  386.76us | cuDeviceTotalMem
                0.00%   8.7200us |      6    1.4530us   230ns   5.4600us | cudaSetupArgument
                0.00%   4.0800us |      3    1.3600us   640ns   2.6100us | cuDeviceGetCount
                0.00%   2.3400us |      2    1.1700us  1.1400us  1.2000us | cudaConfigureCall
                0.00%   2.0800us |      2    1.0400us   620ns   1.4600us | cuDeviceGet

==3316== Unified Memory profiling result:
Device "GeForce GTX 1080 (0)"
  Count  Avg Size  Min Size  Max Size  Total Size  Total Time  Name
-----
   5503      -      -      -      -      533.2729ms  Gpu page fault groups
-----
```

→ [https://tinyurl.com/cuda4dummies/i/12/unified\\_memory\\_example\\_2.cu](https://tinyurl.com/cuda4dummies/i/12/unified_memory_example_2.cu)

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

Compute  
kernel much  
faster now !

```
[sh@n71-006]$ nvcc ./unified_memory_example_2.cu
[sh@n71-006]$ nvprof ./a.out
```

```
==3316== NVPROF is profiling process 3316, command: ./a.out
==3316== Profiling application: ./a.out
==3316== Profiling result:
```

Type	Time(%)	Time	Calls	Avg	Min	Max	Name
GPU activities:	63.49%	341.33ms	1	341.33ms	341.33ms	341.33ms	KrnlDmmyInit(float*, float*,
	36.51%	196.25ms	1	196.25ms	196.25ms	196.25ms	KrnlDmmyCalc(float*, float*,
API calls:	62.34%	537.60ms	2	268.80ms	196.26ms	341.34ms	cudaDeviceSynchronize
	25.62%	220.97ms	3	73.657ms	16.760us	220.91ms	cudaMallocManaged
	11.75%	101.31ms	3	33.770ms	33.658ms	33.845ms	cudaFree
	0.11%	950.09us	2	475.05us	75.991us	874.10us	cudaLaunch
	0.07%	607.36us	94	6.4610us	540ns	244.02us	cuDeviceGetAttribute
	0.07%	578.76us	1	578.76us	578.76us	578.76us	cuDeviceGetName
	0.04%	386.76us	1	386.76us	386.76us	386.76us	cuDeviceTotalMem
	0.00%	8.7200us	6	1.4530us	230ns	5.4600us	cudaSetupArgument
	0.00%	4.0800us	3	1.3600us	640ns	2.6100us	cuDeviceGetCount
	0.00%	2.3400us	2	1.1700us	1.1400us	1.2000us	cudaConfigureCall
	0.00%	2.0800us	2	1.0400us	620ns	1.4600us	cuDeviceGet

```
==3316== Unified Memory profiling result:
```

```
Device "GeForce GTX 1080 (0)"
```

Count	Avg Size	Min Size	Max Size	Total Size	Total Time	Name
5503	-	-	-	-	533.2729ms	Gpu page fault groups

→ [https://tinyurl.com/cuda4dummies/i/12/unified\\_memory\\_example\\_2.cu](https://tinyurl.com/cuda4dummies/i/12/unified_memory_example_2.cu)

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

Compute  
kernel much  
faster now !

```
[sh@n71-006]$ nvcc ./unified_memory_example_2.cu
[sh@n71-006]$ nvprof ./a.out
```

```
==3316== NVPROF is profiling process 3316, command: ./a.out
==3316== Profiling application: ./a.out
==3316== Profiling result:
```

	Type	Time(%)	Time	Calls	Avg	Min	Max	Name
GPU activities:	63.49%	341.33ms		1	341.33ms	341.33ms	341.33ms	KrnlDmyInit(float*, float*, float*, float*)
				1	196.25ms	196.25ms	196.25ms	KrnlDmyCalc(float*, float*, float*, float*)
				2	268.80ms	196.26ms	341.34ms	cudaDeviceSynchronize
				3	73.657ms	16.760us	220.91ms	cudaMallocManaged
				3	33.770ms	33.658ms	33.845ms	cudaFree
				2	475.05us	75.991us	874.10us	cudaLaunch
				94	6.4610us	540ns	244.02us	cuDeviceGetAttribute
				1	578.76us	578.76us	578.76us	cuDeviceGetName
				1	386.76us	386.76us	386.76us	cuDeviceTotalMem
				6	1.4530us	230ns	5.4600us	cudaSetupArgument
				3	1.3600us	640ns	2.6100us	cuDeviceGetCount
				2	1.1700us	1.1400us	1.2000us	cudaConfigureCall
				2	1.0400us	620ns	1.4600us	cuDeviceGet
								0.00% 4.0800us
								0.00% 2.3400us
								0.00% 2.0800us

Still very low effective memory bandwidth,  
$$\frac{3 \times 268435456 \times 4}{196.25 \times 10^{-3}} = 16.4 \text{ GB/s}$$
  
from theoretical 320GB/s

```
==3316== Unified Memory profiling result:
Device "GeForce GTX 1080 (0)"
```

Count	Avg Size	Min Size	Max Size	Total Size	Total Time	Name
5503	-	-	-	-	533.2729ms	Gpu page fault groups

→ [https://tinyurl.com/cuda4dummies/i/12/unified\\_memory\\_example\\_2.cu](https://tinyurl.com/cuda4dummies/i/12/unified_memory_example_2.cu)

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

```
[sh@n71-006]$ nvcc ./unified_memory_example_3.cu
[sh@n71-006]$ nvprof ./a.out

==3497== NVPROF is profiling process 3497, command: ./a.out
==3497== Profiling application: ./a.out
==3497== Profiling result:
   Type  Time(%)   Time     Calls   Avg        Min        Max   Name
GPU activities:  82.04%  1.48704s | 100  14.870ms | 13.222ms  177.58ms | KrnlDmmyCalc(float*, float*,
API calls:      17.96%  325.61ms | 1    325.61ms | 325.61ms  325.61ms | KrnlDmmyInit(float*, float*,
API calls:      82.58%  1.81309s | 101  17.951ms | 13.225ms  325.63ms | cudaDeviceSynchronize
API calls:      0.00%  0.00000s | 3    90.630ms | 16.100us  271.84ms | cudaMallocManaged
API calls:      0.00%  0.00000s | 3    33.626ms | 33.537ms  33.672ms | cudaFree
API calls:      0.00%  0.00000s | 101  85.217us | 42.150us  846.06us | cudaLaunch
API calls:      0.00%  0.00000s | 94   6.2330us | 550ns    246.25us | cuDeviceGetAttribute
API calls:      0.00%  0.00000s | 1    346.65us | 346.65us  346.65us | cuDeviceTotalMem
API calls:      0.00%  0.00000s | 303   298ns    200ns    4.9500us | cudaSetupArgument
API calls:      0.00%  0.00000s | 1    68.361us | 68.361us  68.361us | cuDeviceGetName
API calls:      0.00%  0.00000s | 101   465ns    280ns    3.6800us | cudaConfigureCall
API calls:      0.00%  0.00000s | 3    1.4700us | 640ns    2.9300us | cuDeviceGetCount
API calls:      0.00%  0.00000s | 2    1.0400us | 590ns    1.4900us | cuDeviceGet

==3497== Unified Memory profiling result:
Device "GeForce GTX 1080 (0)"
  Count  Avg Size  Min Size  Max Size  Total Size  Total Time  Name
-----  - - - - -
    5509         -         -         -         -    501.2754ms  Gpu page fault groups
```

100× greatly improves effective memory bandwidth, 217GB/s and compute performance

→ [https://tinyurl.com/cuda4dummies/i/12/unified\\_memory\\_example\\_3.cu](https://tinyurl.com/cuda4dummies/i/12/unified_memory_example_3.cu)

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

```
[sh@n71-006]$ nvcc ./unified_memory_example_4.cu
```

```
[sh@n71-006]$ nvprof ./a.out
```

```
==9756== NVPROF is profiling process 9756, command: ./a.out
```

```
==9756== Profiling application: ./a.out
```

```
==9756== Profiling result:
```

	Type	Time(%)	Time	Calls	Avg	Min	Max	Name
GPU activities:	79.51%	1.32179s	100	13.218ms	13.213ms	13.224ms	KrnlDmmyCalc(float*, float*,	
	20.49%	340.61ms	1	340.61ms	340.61ms	340.61ms	KrnlDmmyInit(float*, float*,	
API calls:	82.06%	1.66280s	101	16.463ms	13.215ms	340.62ms	cudaDeviceSynchronize	
	11.24%	227.77ms	3	75.924ms	27.941us	227.68ms	cudaMallocManaged	
	4.01%	90.52ms	3	33.179ms	32.044ms	33.897ms	cudaFree	
			3	9.0126ms	305.04us	25.975ms	cudaMemPrefetchAsync	
			101	79.679us	49.111us	553.30us	cudaLaunch	
			94	6.1240us	590ns	239.69us	cuDeviceGetAttribute	
			1	361.41us	361.41us	361.41us	cuDeviceTotalMem	
			303	314ns	190ns	7.7900us	cudaSetupArgument	
			1	69.211us	69.211us	69.211us	cuDeviceGetName	
	0.00%	58.470us	101	578ns	350ns	4.6800us	cudaConfigureCall	
	0.00%	4.4700us	3	1.4900us	700ns	2.6800us	cuDeviceGetCount	

Prefetching: fastest,  
244GB/s, fewest  
page faults

```
==9756== Unified Memory profiling result:
```

```
Device "GeForce GTX 1080 (0)"
```

Count	Avg Size	Min Size	Max Size	Total Size	Total Time	Name
3184	-	-	-	-	351.4517ms	Gpu page fault groups

→ [https://tinyurl.com/cuda4dummies/i/12/unified\\_memory\\_example\\_4.cu](https://tinyurl.com/cuda4dummies/i/12/unified_memory_example_4.cu)

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

```
[sh@n71-006]$ nvcc ./unified_memory_example_5.cu
[sh@n71-006]$ nvprof ./a.out

==11798== NVPROF is profiling process 11798, command: ./a.out
==11798== Profiling application: ./a.out
==11798== Profiling result:
   Type  Time(%)   Time     Calls   Avg        Min        Max  Name
GPU activities: 99.38%  1.32146s | 100  13.215ms | 13.208ms  13.220ms | KrnlDmmyCalc(void)
                0.62%  8.2463ms | 1    8.2463ms | 8.2463ms  8.2463ms | KrnlDmmyInit(void)
API calls:     67.65%  1.33004s | 101  13.169ms | 8.2460ms  13.241ms | cudaDeviceSynchronize
                32.30%  635.07ms | 101  6.2878ms | 68.721us  627.53ms | cudaLaunch
                0.03%  582.78us | 94   6.1990us | 610ns    233.16us | cuDeviceGetAttribute
                0.02%  366.10us | 1    366.10us | 366.10us  366.10us | cuDeviceTotalMem
                0.00%  66.710us | 1    66.710us | 66.710us  66.710us | cuDeviceGetName
                0.00%  59.720us | 101  591ns    310ns    13.860us | cudaConfigureCall
                0.00%  4.2500us | 3    1.4160us | 650ns    2.5900us | cuDeviceGetCount
                0.00%  2.2100us | 2    1.1050us | 630ns    1.5800us | cuDeviceGet
```

Globally managed ex  
aquo, 244GB/s, NO  
page faults

→ [https://tinyurl.com/cuda4dummies/i/12/unified\\_memory\\_example\\_5.cu](https://tinyurl.com/cuda4dummies/i/12/unified_memory_example_5.cu)



### GPU Memory Oversubscription

```
#define DBLEARRAY16GB 2147483648

__global__ void KrnlDmmy(float *a)
{
    unsigned int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    a[i] = (double) i + 5.0e00;
    return;
}

int main()
{
    ...
    // unified memory allocation
    cudaMallocManaged(&c, DBLEARRAY16GB * sizeof(double));
    // kernel invocation
    KrnlDmmy <<< numBlocks, threadsPerBlock >>> (c);
    cudaDeviceSynchronize();
    sgnl = check_array(c);
    cudaFree(c);
    return(0);
}
```

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

### GPU Memory Oversubscription

```
#define DBLEARRAY16GB 2147483648

__global__ void KrnlDmmy(float *a)
{
    unsigned int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    a[i] = (double) i + 5.0e00;
    return;
}

int main()
{
    ...
    // unified memory allocation
    cudaMallocManaged(&c, DBLEARRAY16GB * sizeof(double));
    // kernel invocation
    KrnlDmmy <<< numBlocks, threadsPerBlock >>> (c);
    cudaDeviceSynchronize();
    sgnl = check_array(c);
    cudaFree(c);
    return(0);
}
```

16 GB array

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

## GPU Memory Oversubscription

```
#define DBLEARRAY16GB 2147483648

__global__ void KrnlDmmy(float *a)
{
    unsigned int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    a[i] = (double) i + 5.0e00;
    return;
}

int main()
{
    ...
    // unified memory allocation
    cudaMallocManaged(&c, DBLEARRAY16GB * sizeof(double));
    // kernel invocation
    KrnlDmmy <<< numBlocks, threadsPerBlock >>> (c);
    cudaDeviceSynchronize();
    sgnl = check_array(c);
    cudaFree(c);
    return(0);
}
```

Kernel calculation  
(8 GB on-board)

16 GB array

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

## GPU Memory Oversubscription

```
#define DBLEARRAY16GB 2147483648

__global__ void KrnlDmmy(float *a)
{
    unsigned int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    a[i] = (double) i + 5.0e00;
    return;
}

int main()
{
    ...
    // unified memory allocation
    cudaMallocManaged(&c, DBLEARRAY16GB * sizeof(double));
    // kernel invocation
    KrnlDmmy <<< numBlocks, threadsPerBlock >>> (c);
    cudaDeviceSynchronize();
    sgnl = check_array(c);
    cudaFree(c);
    return(0);
}
```

Simple correctness check

Kernel calculation  
(8 GB on-board)

16 GB array

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

```
[sh@n71-006]$ nvcc ./unified_memory_oversubscription.cu
[sh@n71-006]$ nvprof ./a.out
```

```
==19259== NVPROF is profiling process 19259, command: ./a.out
expecting 2305843018877370368 while receiving 2305843018877370368.000000
==19259== Profiling application: ./a.out
==19259== Profiling result:
```

Type	Time(%)	Time	Calls	Avg	Min	Max	Name
GPU activities:	100.00%	4.01542s	1	4.01542s	4.01542s	4.01542s	KrnlDmmy(double*)
API calls:	76.93%	4.01561s	1	4.01561s	4.01561s	4.01561s	cudaDeviceSynchronize
	17.36%	906.20ms	1	906.20ms	906.20ms	906.20ms	cudaFree
	5.69%	296.86ms	1	296.86ms	296.86ms	296.86ms	cudaMallocManaged
	0.01%	604.35us	94	6.4290us	550ns	250.47us	cuDeviceGetAttribute
	0.01%	364.77us	1	364.77us	364.77us	364.77us	cuDeviceTotalMem
	0.00%	142.05us	1	142.05us	142.05us	142.05us	cudaLaunch
	0.00%	70.930us	1	70.930us	70.930us	70.930us	cuDeviceGetName
	0.00%	7.6000us	1	7.6000us	7.6000us	7.6000us	cudaSetupArgument
	0.00%	4.9100us	3	1.6360us	810ns	3.2200us	cuDeviceGetCount
	0.00%	2.8600us	1	2.8600us	2.8600us	2.8600us	cudaConfigureCall

Correctness  
test ok

Many page  
faults

```
==19259== Unified Memory profiling result:
```

```
Device "GeForce GTX 1080 (0)"
```

Count	Avg Size	Min Size	Max Size	Total Size	Total Time	Name
52018	322.53KB	4.0000KB	2.0000MB	16.00000GB	2.348844s	Device To Host
40520	-	-	-	-	5.497636s	Gpu page fault groups

```
Total CPU Page faults: 28113
```

→ [https://tinyurl.com/cuda4dummies/i/12/unified\\_memory\\_oversubscription.cu](https://tinyurl.com/cuda4dummies/i/12/unified_memory_oversubscription.cu)

Much more  
time on  
CPU

# MEMORY HIERARCHY

CUDA C-PROGRAMMING GUIDE, PROGRAMMING MODEL CONT.

How did this  
work in the  
days before  
CUDA man-  
aged unified  
memory...



→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, UNIFIED MEMORY PROGRAMMING CONT.

- On pre-Pascal-class GPUs `cudaMallocManaged()` behaves like `cudaMalloc()`
- On pre-Pascal-class GPUs managed allocations are automatically visible to all GPUs in a system via peer-to-peer capabilities
- Managed memory is allocated in GPU memory as long as all GPUs have peer-to-peer support enabled, otherwise the driver will migrate all managed allocations to the CPU/host memory (“zero-copy” memory)
- P2P GPUs will experience PCIe bandwidth restrictions

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, DEVICE MEMORY

On pre-Pascal-class GPUs a distinction is made between host- and device-memory. Kernels operate out of device memory, so the CUDA runtime provides functions to allocate, deallocate, and copy data into device memory as well as transfer data between host memory and device memory. A typical sequence of operations is:

1. Declare and allocate host and device memory
2. Initialize host data
3. Transfer data from the host to the device
4. Execute one or more kernels
5. Transfer results from the device to the host

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

→ <https://devblogs.nvidia.com/easy-introduction-cuda-c-and-c>



# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, DEVICE MEMORY CONT.

```
#define ARRAYDIM 268435456
__global__ void KrnlDmmy(float *x, float *y, float *z)
{
    int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    z[i] = x[i] + y[i];
    return;
}

int main()
{
    ...
    // host memory allocation, hb, hc analogous
    ha = (float *) malloc(ARRAYDIM * sizeof(float));
    // device memory allocation, db, dc analogous
    cudaMalloc(&da, ARRAYDIM * sizeof(float));
    // host to device memory transfer
    cudaMemcpy(da, ha, ARRAYDIM * sizeof(float), cudaMemcpyHostToDevice);
    // kernel invocation
    KrnlDmmy <<< numBlocks, threadsPerBlock >>> (da, db, dc);
    // device to host memory transfer
    cudaMemcpy(hc, dc, ARRAYDIM * sizeof(float), cudaMemcpyDeviceToHost);
    // free device/host memory
    cudaFree(da);
    free(ha);
    ...
}
```

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, DEVICE MEMORY CONT.

```
#define ARRAYDIM 268435456
__global__ void KrnlDmmy(float *x, float *y, float *z)
{
    int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    z[i] = x[i] + y[i];
    return;
}

int main()
{
    ...
    // host memory allocation, hb, hc analogous
    ha = (float *) malloc(ARRAYDIM * sizeof(float));
    // device memory allocation, db, dc analogous
    cudaMalloc(&da, ARRAYDIM * sizeof(float));
    // host to device memory transfer
    cudaMemcpy(da, ha, ARRAYDIM * sizeof(float), cudaMemcpyHostToDevice);
    // kernel invocation
    KrnlDmmy <<< numBlocks, threadsPerBlock >>> (da, db, dc);
    // device to host memory transfer
    cudaMemcpy(hc, dc, ARRAYDIM * sizeof(float), cudaMemcpyDeviceToHost);
    // free device/host memory
    cudaFree(da);
    free(ha);
    ...
}
```

Kernel  
memory  
alloca-  
tion &  
set up

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, DEVICE MEMORY CONT.

```
#define ARRAYDIM 268435456
__global__ void KrnlDmmy(float *x, float *y, float *z)
{
    int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    z[i] = x[i] + y[i];
    return;
}

int main()
{
    ...
    // host memory allocation, hb, hc analogous
    ha = (float *) malloc(ARRAYDIM * sizeof(float));
    // device memory allocation, db, dc analogous
    cudaMalloc(&da, ARRAYDIM * sizeof(float));
    // host to device memory transfer
    cudaMemcpy(da, ha, ARRAYDIM * sizeof(float), cudaMemcpyHostToDevice);
    // kernel invocation
    KrnlDmmy <<< numBlocks, threadsPerBlock >>> (da, db, dc);
    // device to host memory transfer
    cudaMemcpy(hc, dc, ARRAYDIM * sizeof(float), cudaMemcpyDeviceToHost);
    // free device/host memory
    cudaFree(da);
    free(ha);
    ...
}
```

Kernel  
memory  
alloca-  
tion &  
set up

Std  
kernel  
call

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, DEVICE MEMORY CONT.

```
#define ARRAYDIM 268435456
__global__ void KrnlDmmy(float *x, float *y, float *z)
{
    int i;
    i = (blockIdx.x * blockDim.x) + threadIdx.x;
    z[i] = x[i] + y[i];
    return;
}

int main()
{
    ...
    // host memory allocation, hb, hc analogous
    ha = (float *) malloc(ARRAYDIM * sizeof(float));
    // device memory allocation, db, dc analogous
    cudaMalloc(&da, ARRAYDIM * sizeof(float));
    // host to device memory transfer
    cudaMemcpy(da, ha, ARRAYDIM * sizeof(float), cudaMemcpyHostToDevice);
    // kernel invocation
    KrnlDmmy <<< numBlocks, threadsPerBlock >>> (da, db, dc);
    // device to host memory transfer
    cudaMemcpy(hc, dc, ARRAYDIM * sizeof(float), cudaMemcpyDeviceToHost);
    // free device/host memory
    cudaFree(da);
    free(ha);
    ...
}
```

Kernel  
memory  
alloca-  
tion &  
set up

Std  
kernel  
call

Device  
memory  
back-  
transfer

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, DEVICE MEMORY CONT.

- Special forms `cudaMallocPitch()` and `cudaMalloc3D()` for 2D and 3D arrays
- Optimized for best performance when accessing via pointers
- Also good for device copies `cudaMemcpy2D()` and `cudaMemcpy3D()`
- Returned pitch (or stride) must be used to access array elements
- Optimally padded to meet alignment requirements
- Additional types of global memory, e.g.  
`__device__ float *devPointer`

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

# MEMORY HIERARCHY

## CUDA C-PROGRAMMING GUIDE, LOCAL MEMORY, REGISTERS

- Local memory for certain automatic variables
- Access to local memory is similar to global memory, i.e. high latency and low bandwidth
- Organized such that consecutive 32-bit words are accessed by consecutive thread IDs
- Analyzable with **Nsight**
- Fastest memory is registers, L1, for small, statically indexed arrays, e.g. `A[16]`

→ <https://docs.nvidia.com/cuda/cuda-c-programming-guide>

→ <https://stackoverflow.com/questions/10297067/in-a-cuda-kernel-how-do-i-store-an-array-in-local-thread-local-memory>

# TAKE HOME MESSAGES

- Unified memory — a big improvement

# TAKE HOME MESSAGES

- Unified memory — a big improvement
- Single address space also facilitates easier handling of 2D/3D arrays on the GPU



# TAKE HOME MESSAGES

- Unified memory — a big improvement
- Single address space also facilitates easier handling of 2D/3D arrays on the GPU
- Straightforward profiling with nvprof

# TAKE HOME MESSAGES

- Unified memory — a big improvement
- Single address space also facilitates easier handling of 2D/3D arrays on the GPU
- Straightforward profiling with nvprof
- Page migration engine facilitates seamless data transfer for array sizes exceeding largely the amount of RAM available on board of the GPU