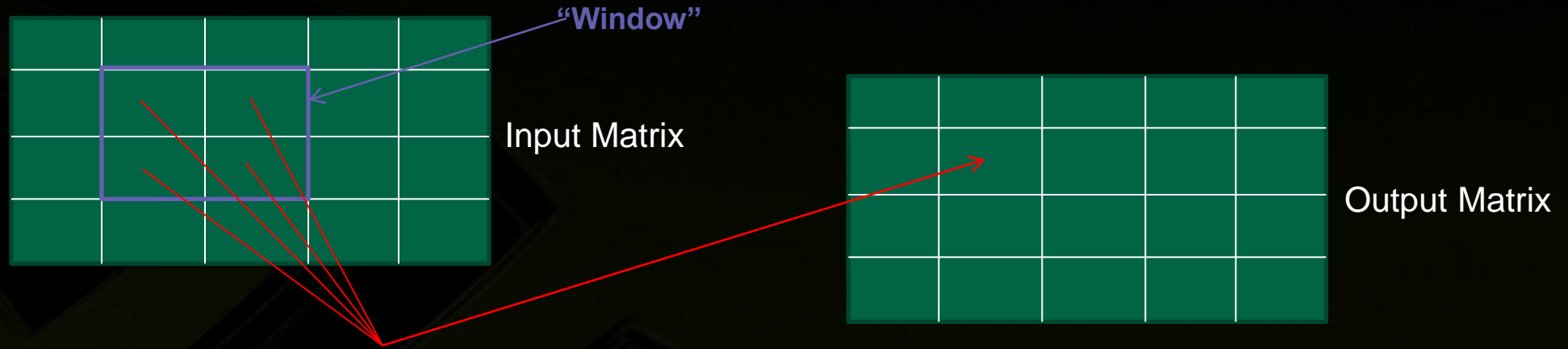# Accelerate your Application with CUDA C

# 2D Minimum Algorithm

- **Consider applying a 2D window to a 2D array of elements**
  - **Each output element is the minimum of input elements within the window**



"Window"

Input Matrix

Output Matrix

- **2D operations like this are found in many fundamental algorithms**
  - **Interpolation,  Convolution, Filtering**
- **Applications in seismic processing, weather simulation, image processing, etc**

http://developer.nvidia.com/cuda

# 2D Minimum: C Version

```c
#define WIN_SIZE 16
#define N 20000

int main() {
  int size=N*N*sizeof(int);
  int i, j, x, y, temp;
  //allocate resources
  int *cell=(int *)malloc(size);  //input
  int *node=(int *)malloc(size);  //output


  initializeArray(cell,N);
  initializeArray(node,N);
```

**Allocate Resources**

**Initialize data**

```c
for(i=0;i<N;i++)
    for(j=0;j<N;j++)
    {
    //find minimum in window
        temp = node[i][j];
        for(x=0;x<WIN_SIZE;x++)
          for(y=0;y<WIN_SIZE;y++)
                if (temp> cell[i+x][j+y])
                    temp = cell[i+x][j+y];
      node[i][j] = temp;
    }

 //free resources
  free(cell); free(node);
}
```

**Loop over dataset**

**Loop over window**

**Find min**

**Cleanup**

http://developer.nvidia.com/cuda

# 2D Minimum:  C Version

```c
#define WIN_SIZE 16
#define N 20000

int main() {
 int size=N*N*sizeof(int);
 int i, j, x, y, temp;
 //alloc
 int *cel
 int *nod
```

```c
for(i=0;i<N;i++)
    for(j=0;j<N;j++)
    {
    //find minimum in window
        temp = node[i][j];
```

| Algorithm | Device | Compute | Speedup |
|-----------|--------|---------|---------|
| C (serial) | Xeon X5650 | 96.1sec | -- |

```c
                temp = cell[i+x][j+y];
        node[i][j] = temp;
    }
```

```c
  //free resources
   free(cell); free(node);
  }
```

```c
initializeArray(cell,N);
initializeArray(node,N);
```

http://developer.nvidia.com/cuda

# CPU + GPU = Big Speedup

**Application Code**

Compute-Intensive Functions

Use CUDA to Parallelize

Rest of Sequential
CPU Code

**GPU**

**CPU**

+

http://developer.nvidia.com/cuda

# Explicit Data Movement

**CPU**
(host)

**GPU**
(device)

**CPU Memory**

| in |
| --- |
| out |

**PCI Bus**

**GPU Memory**

| d_in |
| --- |
| d_out |

http://developer.nvidia.com/cuda

# 2D Window: CUDA – Main Program

## C

```c
int main() {
  int size=N*N*sizeof(int);
  //allocate resources
  int *cell=(int*)malloc(size); //input
  int *node=(int*)malloc(size); //output

  initializeArray(cell,N);
  initializeArray(node,N);
// nested for loops
  for...
   for ...
    for ...
     for ....
//free resources
  free(in); free(out);
}
```

## CUDA C

```c
int main() {
  int size=N*N*sizeof(int);
  //allocate resources
  int *cell=(float*)malloc(size);
  int *node=(float*)malloc(size);
  int *d_cell; cudaMalloc(&d_cell,size);
  int *d_node; cudaMalloc(&d_node,size);
  initializeArray(cell,N); initializeArray(node,N);

  cudaMemcpy(d_cell,cell,size,cudaMemcpyHostToDevice);
  cudaMemcpy(d_node,node,size,cudaMemcpyHostToDevice);
  compute_win2D<<<nblocks, nthreads>>>(d_node,d_cell);
  //free resources
  free(cell); free(node);
  cudaFree(d_cell); cudaFree(d_node);
}
```

http://developer.nvidia.com/cuda

# 2D Window: CUDA – Main Program

## C

```
int main() {
  int size=N*N*sizeof(int);
  //allocate resources
  int *cell=(int*)malloc(size); //input
  int *node=(int*)malloc(size); //output

  initializeArray(cell,N);
  initializeArray(node,N);
// nested for loops
  for…
   for …
    for …
     for ….
//free resources
  free(in); free(out);

}
```

## CUDA C

```
int main() {
  int size=N*N*sizeof(int);
  //allocate resources
  int *cell=(float*)malloc(size);
  int *node=(float*)malloc(size);
  int *d_cell; cudaMalloc(&d_cell,size);
  int *d_node; cudaMalloc(&d_node,size);
  initializeArray(cell,N); initializeArray(node,N);

  cudaMemcpy(d_cell,cell,size,cudaMemcpyHostToDevice);
  cudaMemcpy(d_node,node,size,cudaMemcpyHostToDevice);
  compute_win2D<<<nblocks, nthreads>>>(d_node,d_cell);
  //free resources
  free(cell); free(node);
  cudaFree(d_cell); cudaFree(d_node);

}
```

**Allocate Device Memory**

http://developer.nvidia.com/cuda

# 2D Window: CUDA – Main Program

## C

```
int main() {
  int size=N*N*sizeof(int);
  //allocate resources
  int *cell=(int*)malloc(size); //input
  int *node=(int*)malloc(size); //output


  initializeArray(cell,N);
  initializeArray(node,N);
// nested for loops
  for...
   for ...
    for ...
     for ....
//free resources
  free(in); free(out);
}
```

## CUDA C

```
int main() {
  int size=N*N*sizeof(int);
  //allocate resources
  int *cell=(float*)malloc(size);
  int *node=(float*)malloc(size);
  int *d_cell; cudaMalloc(&d_cell,size);
  int *d_node; cudaMalloc(&d_node,size);
  initializeArray(cell,N); initializeArray(node,N


  cudaMemcpy(d_cell,cell,size,cudaMemcpyHostToDevice);
  cudaMemcpy(d_node,node,size,cudaMemcpyHostToDevice);
  compute_win2D<<<nblocks, nthreads>>>(d_node,d_cell);
  //free resources
  free(cell); free(node);
  cudaFree(d_cell); cudaFree(d_node);
}
```

**Copy Data to the Device**

http://developer.nvidia.com/cuda

# 2D Window: CUDA – Main Program

## C

```c
int main() {
 int size=N*N*sizeof(int);
 //allocate resources
 int *cell=(int*)malloc(size); //input
 int *node=(int*)malloc(size); //output

 initializeArray(cell,N);
 initializeArray(node,N);
// nested for loops
 for...
  for ...
   for ...
    for ....
//free resources
 free(in); free(out);
}
```

## CUDA C

```c
int main() {
 int size=N*N*sizeof(int);
 //allocate resources
 int *cell=(float*)malloc(size);
 int *node=(float*)malloc(size);
 int *d_cell; cudaMalloc(&d_cell,size);
 int *d_node; cudaMalloc(&d_node,size);
 initializeArray(cell,N); initializeArray(node,N);


 cudaMemcpy(d_cell,cell,size,cudaMemcpyHostToDevice);
 cudaMemcpy(d_node,node,size,cudaMemcpyHostToDevice);
 compute_win2D<<<nblocks, nthreads>>>(d_node,d_cell);
 //free resources
 free(cell); f
 cudaF              _node);
}
```

**Call Cuda Function**

**Launch Parameters**

**Device Pointers**

http://developer.nvidia.com/cuda

# 2D Window: CUDA – Main Program

## C

```c
int main() {
  int size=N*N*sizeof(int);
  //allocate resources
  int *cell=(int*)malloc(size); //input
  int *node=(int*)malloc(size); //output

  initializeArray(cell,N);
  initializeArray(node,N);
// nested for loops
  for...
    for ...
     for ...
      for ....
//free resources
  free(in); free(out);
}
```

## CUDA C

```c
int main() {
  int size=N*N*sizeof(int);
  //allocate resources
  int *cell=(float*)malloc(size);
  int *node=(float*)malloc(size);
  int *d_cell; cudaMalloc(&d_cell,size);
  int *d_node; cudaMalloc(&d_node,size);
  initializeArray(cell,N); initializeArray(node,N);

  cudaMemcpy(d_cell,cell,size,cudaMemcpyHostToDevice);
  cudaMemcpy(d_node,node,size,cudaMemcpyHostToDevice);
  compute_win2D<<<nblocks, nthreads>>>(d_node.d_cell);
  //free resources
  free(cell); free(node);
  cudaFree(d_cell); cudaFree(d_node);
}
```

Cleanup

http://developer.nvidia.com/cuda

# Parallel Execution Model

- A CUDA C function is executed by **many parallel threads**
- Threads are organized as a **grid** of independent thread **blocks**

**Grid**

| Block 0 | Block 1 | ... | Block N_B-2 | Block N_B-1 |

# 2D Window: CUDA – Kernel Function

## C

```
for(i=0;i<N;i++)
    for(j=0;j<N;j++)
    {
    //find minimum in window
        temp = node[i][j];
        for(x=0;x<WIN_SIZE;x++)
         for(y=0;y<WIN_SIZE;y++)
                if (temp> cell[i+x][j+y])
                    temp = cell[i+x][j+y];
      node[i][j] = temp;
    }
```

## CUDA C

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
  int idx=blockIdx.x*blockDim.x+threadIdx.x;
  int idy=blockIdx.y*blockDim.y+threadIdx.y;
  int temp, x, y;
  if((idx<N)&&(idy<N)) {
//find minimum in window
        temp = knode[idx][idy];
        for(x=0;x<WIN_SIZE;x++)
          for(y=0;y<WIN_SIZE;y++)
                if (temp> kcell[idx+x][idy+y])
                    temp = kcell[idx+x][idy+y];
      knode[i][j] = temp;
}
```

http://developer.nvidia.com/cuda

# 2D Window: CUDA – Kernel Function

## C

```
for(i=0;i<N;i++)
    for(j=0;j<N;j++)
    {
    //find minimum in window
        temp = node[i][j];
        for(x=0;x<WIN_SIZE;x++)
         for(y=0;y<WIN_SIZE;y++)
                if (temp> cell[i+x][j+y])
                    temp = cell[i+x][j+y];
        node[i][j] = temp;
    }
```

## CUDA C

Add __global__ Keyword

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
    int  idx=blockIdx.x*blockDim.x+threadIdx.x;
    int  idy=blockIdx.y*blockDim.y+threadIdx.y;
    int  temp,  x, y;
    if((idx<N)&&(idy<N)) {
    //find minimum in window
        temp = knode[idx][idy];
        for(x=0;x<WIN_SIZE;x++)
         for(y=0;y<WIN_SIZE;y++)
                if (temp> kcell[idx+x][idy+y])
                    temp = kcell[idx+x][idy+y];
        knode[i][j] = temp;
}
```
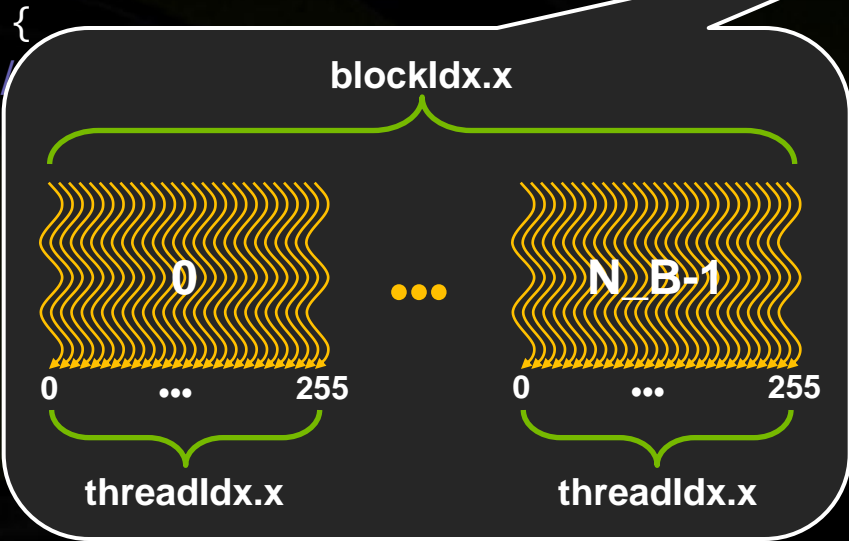
http://developer.nvidia.com/cuda

# 2D Window: CUDA – Kernel Function

## C

## CUDA C

**Replace Outer Loops With an Index Calculation**

```
for(i=0;i<N;i++)
    for(j=0;j<N;j++)
    {
```



blockIdx.x

0        ...        N_B-1

0    ...    255        0    ...    255

threadIdx.x        threadIdx.x

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
    int  idx=blockIdx.x*blockDim.x+threadIdx.x;
    int  idy=blockIdx.y*blockDim.y+threadIdx.y;
    int  temp,  x, y;
    if((idx<N)&&(idy<N)) {
//find minimum in window
        temp = knode[idx][idy];
        for(x=0;x<WIN_SIZE;x++)
          for(y=0;y<WIN_SIZE;y++)
               if (temp> kcell[idx+x][idy+y])
                   temp = kcell[idx+x][idy+y];
        knode[i][j] = temp;
}
```
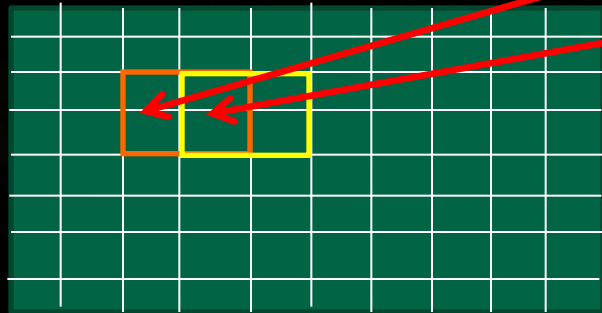
# 2D Window: CUDA – Kernel Function

## C

## CUDA C

```
for(i=0;i<N;i++)
    for(j=0;j<N;j++)
    {
    //find minimum
        temp = nod
        for(x=0;x<W
            for(y=0;y<
                if (tem
                    temp = cell[i+x][j+y];
        node[i][j] = temp;
    }
```

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
    int  idx=blockIdx.x*blockDim.x+threadIdx.x;
    int  idy=blockIdx.y*blockDim.y+threadIdx.y;
    int  temp x, y;

                if (temp> kcell[idx+x][idy+y])
                    temp = kcell[idx+x][idy+y];
        knode[i][j] = temp;
}
```

| Algorithm | Device | Compute | Speedup |
|-----------|--------|---------|---------|
| C (serial) | X5650 | 96.1sec | -- |
| CUDA | M2090 | 8.33sec | 11.5x |

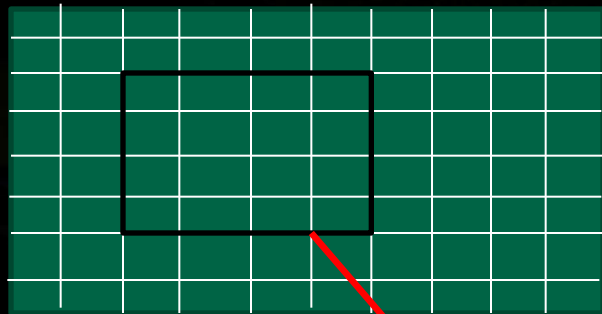http://developer.nvidia.com/cuda

# Observation: Data Reuse

Thread(i)

Thread(i+1)

Neighboring threads read the same elements.

Global Memory

# Optimization: Use Shared Memory

Global Memory

SMEM COPY

| 1 | 2 | BLK_SIZE |
| 3 | 4 | WSIZE |

Shared Memory

http://developer.nvidia.com/cuda

# CUDA C: Optimized Kernel

```c
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
  __shared__ int smem[BLK_SIZE+WSIZE][BLK_SIZE+WSIZE];
  int  idx=blockIdx.x*blockDim.x+threadIdx.x;
  int  idy=blockIdx.y*blockDim.y+threadIdx.y;
  int temp, x, y;
  if((idx<N)&&(idy<N)) {

  smem[threadIdx.x][threadIdx.y]=kcell[idx][idy];
  if (threadIdx.y > (N-WSIZE))
    smem[threadIdx.x][threadIdx.y + WSIZE]=kcell[idx][idy+WSIZE];
  if (threadIdx.x >(N-WSIZE))
    smem[threadIdx.x + WSIZE][threadIdx.y]=kcell[idx+WSIZE][idy];
  if ((threadIdx.x >(N-WSIZE)) && (threadIdx.y > (N-WSIZE)))
   smem[threadIdx.x+WSIZE][threadIdx.y+WSIZE]=
             kcell[idx+WSIZE][idy+WSIZE];
//wait for all threads to finish read
  __syncthreads();

//find minimum in window
  temp = knode[idx][idy];
  for(x=0;x<WSIZE;x++)
   for(y=0;y<WSIZE;y++)
     if (temp> smem[threadIdx.x+x][threadIdx.y+y])
       temp = smem[threadIdx.x+x][threadIdx.y+y];
     knode[i][j] = temp;
  }
}
```
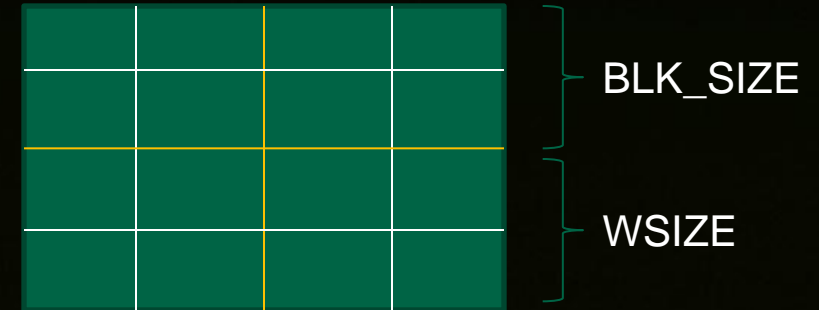
# CUDA C: Optimized Kernel

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
  __shared__ int smem[BLK_SIZE+WSIZE][BLK_SIZE+WSIZE];
  int  idx=blockIdx.x*blockDim.x+threadIdx.x;
  int  idy=blockIdx.y*blockDim.y+threadIdx.y;
  int  temp,  x, y;
  if((idx<N)&&(idy<N)) {

  smem[threadIdx.x][threadIdx.y]=kcell[idx][idy];
  if (threadIdx.y > (N-WSIZE))
    smem[threadIdx.x][threadIdx.y + WSIZE]=kcell[idx][idy+WSIZE];
  if (threadIdx.x >(N-WSIZE))
    smem[threadIdx.x + WSIZE][threadIdx.y]=kcell[idx+WSIZE][idy];
  if ((threadIdx.x >(N-WSIZE)) && (threadIdx.y > (N-WSIZE)))
   smem[threadIdx.x+WSIZE][threadIdx.y+WSIZE]=
              kcell[idx+WSIZE][idy+WSIZE];
//wait for all threads to finish read
  __syncthreads();
```

**Allocate Shared Memory for Each Block**

BLK_SIZE

WSIZE

Shared Memory

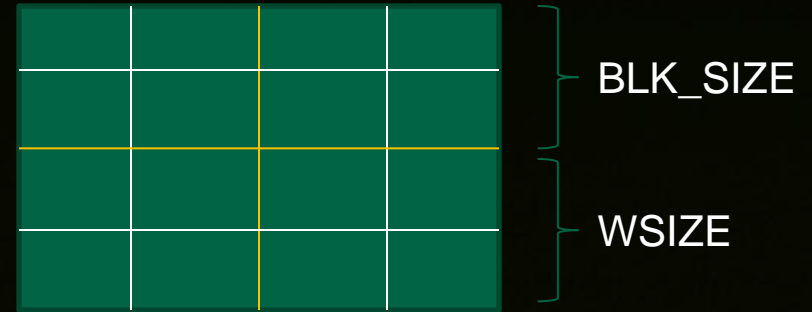http://developer.nvidia.com/cuda

# CUDA C: Optimized Kernel

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
  __shared__ int smem[BLK_SIZE+WSIZE][BLK_SIZE+WSIZE];
  int  idx=blockIdx.x*blockDim.x+threadIdx.x;
  int  idy=blockIdx.y*blockDim.y+threadIdx.y;
  int  temp, x, y;
  if((idx<N)&&(idy<N)) {

  smem[threadIdx.x][threadIdx.y]=kcell[idx][idy];
  if (threadIdx.y > (N-WSIZE))
    smem[threadIdx.x][threadIdx.y + WSIZE]=kcell[idx][idy+WSIZE];
  if (threadIdx.x >(N-WSIZE))
    smem[threadIdx.x + WSIZE][threadIdx.y]=kcell[idx+WSIZE][idy];
  if ((threadIdx.x >(N-WSIZE)) && (threadIdx.y > (N-WSIZE)))
   smem[threadIdx.x+WSIZE][threadIdx.y+WSIZE]=
            kcell[idx+WSIZE][idy+WSIZE];
//wait for all threads to finish read
  __syncthreads();
```

Calculate Global Memory Indices

BLK_SIZE

WSIZE

Shared Memory

http://developer.nvidia.com/cuda
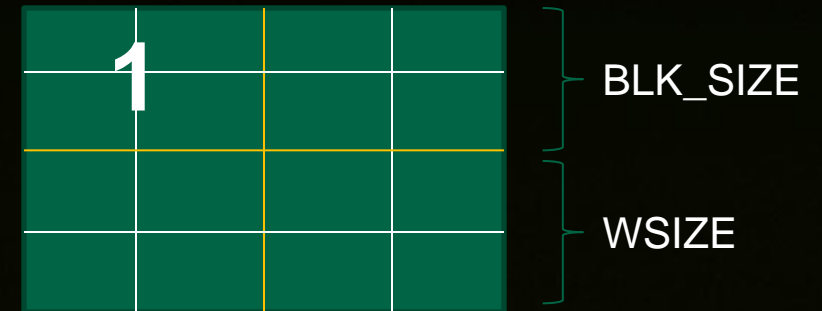
# CUDA C: Optimized Kernel

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
    __shared__ int smem[BLK_SIZE+WSIZE][BLK_SIZE+WSIZE];
    int  idx=blockIdx.x*blockDim.x+threadIdx.x;
    int  idy=blockIdx.y*blockDim.y+threadIdx.y;
    int  temp,  x, y;
    if((idx<N)&&(idy<N)) {

    smem[threadIdx.x][threadIdx.y]=kcell[idx][idy];
    if (threadIdx.y > (N-WSIZE))
        smem[threadIdx.x][threadIdx.y + WSIZE]=kcell[idx][idy+WSIZE];
    if (threadIdx.x >(N-WSIZE))
        smem[threadIdx.x + WSIZE][threadIdx.y]=kcell[idx+WSIZE][idy];
    if ((threadIdx.x >(N-WSIZE)) && (threadIdx.y > (N-WSIZE)))
        smem[threadIdx.x+WSIZE][threadIdx.y+WSIZE]=
                    kcell[idx+WSIZE][idy+WSIZE];
    //wait for all threads to finish read
    __syncthreads();
```

**Load Interior Region 1 to Shared Memory**



Shared Memory

BLK_SIZE

WSIZE

http://developer.nvidia.com/cuda
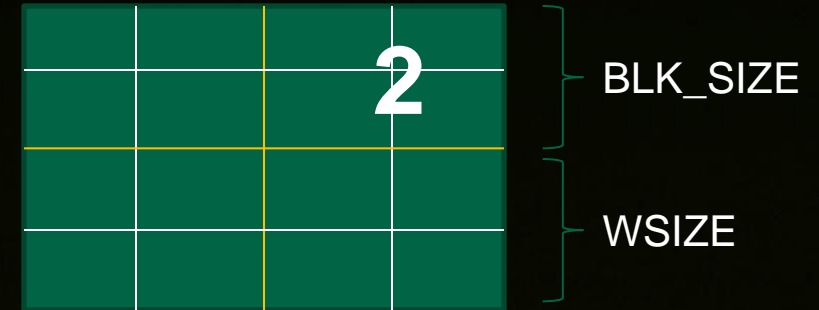
# CUDA C: Optimized Kernel

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
  __shared__ int smem[BLK_SIZE+WSIZE][BLK_SIZE+WSIZE];
  int  idx=blockIdx.x*blockDim.x+threadIdx.x;
  int  idy=blockIdx.y*blockDim.y+threadIdx.y;
  int  temp,  x, y;
  if((idx<N)&&(idy<N)) {

  smem[threadIdx.x][threadIdx.y]=kcell[idx][idy];
  if (threadIdx.y > (N-WSIZE))
     smem[threadIdx.x][threadIdx.y + WSIZE]=kcell[idx][idy+WSIZE];
  if (threadIdx.x >(N-WSIZE))
     smem[threadIdx.x + WSIZE][threadIdx.y]=kcell[idx+WSIZE][idy];
  if ((threadIdx.x >(N-WSIZE)) && (threadIdx.y > (N-WSIZE)))
   smem[threadIdx.x+WSIZE][threadIdx.y+WSIZE]=
              kcell[idx+WSIZE][idy+WSIZE];
  //wait for all threads to finish read
   __syncthreads();
```
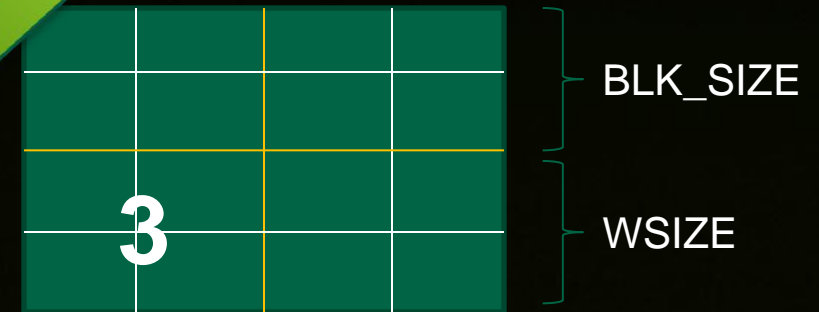
Load Halo Region 2 to Shared Memory

2

BLK_SIZE

WSIZE

Shared Memory

http://developer.nvidia.com/cuda

# CUDA C: Optimized Kernel

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
  __shared__ int smem[BLK_SIZE+WSIZE][BLK_SIZE+WSIZE];
  int  idx=blockIdx.x*blockDim.x+threadIdx.x;
  int  idy=blockIdx.y*blockDim.y+threadIdx.y;
  int  temp, x, y;
  if((idx<N)&&(idy<N)) {

  smem[threadIdx.x][threadIdx.y]=kcell[idx][idy];
  if (threadIdx.y > (N-WSIZE))
     smem[threadIdx.x][threadIdx.y + WSIZE]=kcell[idx][idy+WSIZE];
  if (threadIdx.x >(N-WSIZE))
     smem[threadIdx.x + WSIZE][threadIdx.y]=kcell[idx+WSIZE][idy];
  if ((threadIdx.x >(N-WSIZE)) && (threadIdx.y > (N-WSIZE)))
     smem[threadIdx.x+WSIZE][threadIdx.y+WSIZE]=
                  kcell[idx+WSIZE][idy+WSIZE];
//wait for all threads to finish read
   __syncthreads();
```

Load Halo Region 3 to Shared Memory

3

BLK_SIZE

WSIZE

Shared Memory
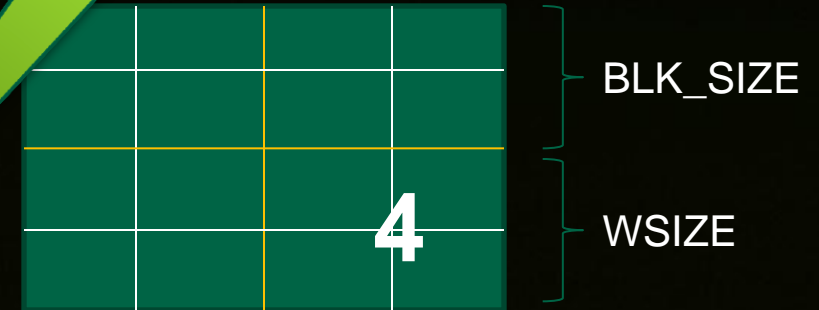
http://developer.nvidia.com/cuda

# CUDA C: Optimized Kernel

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
  __shared__ int smem[BLK_SIZE+WSIZE][BLK_SIZE+WSIZE];
  int  idx=blockIdx.x*blockDim.x+threadIdx.x;
  int  idy=blockIdx.y*blockDim.y+threadIdx.y;
  int  temp,  x, y;
  if((idx<N)&&(idy<N)) {

  smem[threadIdx.x][threadIdx.y]=kcell[idx][idy];
  if (threadIdx.y > (N-WSIZE))
    smem[threadIdx.x][threadIdx.y + WSIZE]=kcell[idx][idy+WSIZE];
  if (threadIdx.x >(N-WSIZE))
    smem[threadIdx.x + WSIZE][threadIdx.y]=kcell[idx+WSIZE][idy];
  if ((threadIdx.x >(N-WSIZE)) && (threadIdx.y > (N-WSIZE)))
   smem[threadIdx.x+WSIZE][threadIdx.y+WSIZE]=
            kcell[idx+WSIZE][idy+WSIZE];
//wait for all threads to finish read
  __syncthreads();
```

Load Halo Region 4 to Shared Memory

BLK_SIZE

**4**

WSIZE

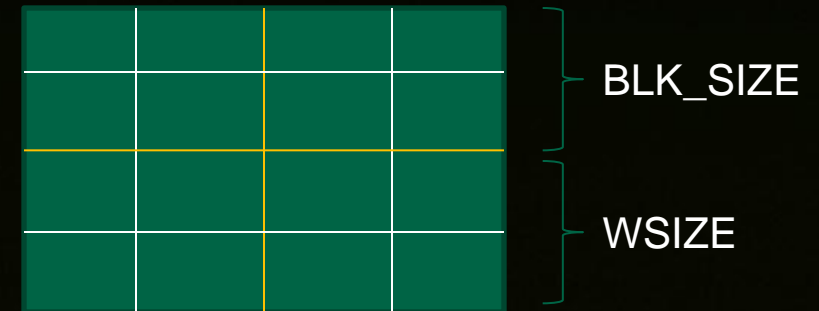Shared Memory

# CUDA C: Optimized Kernel

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
  __shared__ int smem[BLK_SIZE+WSIZE][BLK_SIZE+WSIZE];
  int  idx=blockIdx.x*blockDim.x+threadIdx.x;
  int  idy=blockIdx.y*blockDim.y+threadIdx.y;
  int  temp, x, y;
  if((idx<N)&&(idy<N)) {

  smem[threadIdx.x][threadIdx.y]=kcell[idx][idy];
  if (threadIdx.y > (N-WSIZE))
    smem[threadIdx.x][threadIdx.y + WSIZE]=kcell[idx][idy+WSIZE];
  if (threadIdx.x >(N-WSIZE))
    smem[threadIdx.x + WSIZE][threadIdx.y]=kcell[idx+WSIZE][idy];
  if ((threadIdx.x >(N-WSIZE)) && (threadIdx.y > (N-WSIZE)))
    smem[threadIdx.x+WSIZE][threadIdx.y+WSIZE]=
            kcell[idx+WSIZE][idy+WSIZE];
  //wait for all threads to finish read
  __syncthreads();
```

Shared Memory

BLK_SIZE

WSIZE

**Wait for All Threads to Finish Writing to Shared Memory**

http://developer.nvidia.com/cuda

# CUDA C: Optimized Kernel

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
  __shared__ int smem[BLK_SIZE+WSIZE][BLK_SIZE+WSIZE];
  int  idx=blockIdx.x*blockDim.x+threadIdx.x;
  int  idy=blockIdx.y*blockDim.y+threadIdx.y;
  int  temp, x, y;
  if((idx<N)&&(idy<N)) {

  smem[threadIdx.x][threadIdx.y]=kcell[idx][idy];
  if (threadIdx.y > (N-WSIZE))
    smem[threadIdx.x][threadIdx.y + WSIZE]=kcell[idx][idy+WSIZE];
  if (threadIdx.x >(N-WSIZE))
    smem[threadIdx.x + WSIZE][threadIdx.y]=kcell[idx+WSIZE][idy];
  if ((threadIdx.x >(N-WSIZE)) && (threadIdx.y > (N-WSIZE)))
    smem[threadIdx.x+WSIZE][threadIdx.y+WSIZE]=
                kcell[idx+WSIZE][idy+WSIZE];
//wait for all threads to finish read
    __syncthreads();
```

```
//find minimum in window
  temp = knode[idx][idy];
  for(x=0;x<WSIZE;x++)
    for(y=0;y<WSIZE;y++)
      if (temp> smem[threadIdx.x+x][threadIdx.y+y])
        temp = smem[threadIdx.x+x][threadIdx.y+y];
      knode[i][j] = temp;
  }
}
```

**Find Minimum: Read Input from Shared Memory, Accumulate into a Register**

http://developer.nvidia.com/cuda

# CUDA C: Optimized Kernel

```
__global__ void compute_win2D(int knode[][N], int kcell[][N])
{
  __shared__ int smem[BLK_SIZE+WSIZE][BLK_SIZE+WSIZE];
  int  idx=blockIdx.x*blockDim.x+threadIdx.x;
  int  idy=blockIdx.y*blockDim.y+threadIdx.y;
  int temp, x, y;
  if((idx<N)&&(idy<N)) {

  smem[threadIdx.x][threadIdx.y]=kcell[idx][idy];
  if (threadIdx.y > (N-WSIZE))
    smem[threadIdx.x][threadIdx.y + WSIZE]=kcell[idx][idy+WSIZE];
  if (threadIdx.x >(N-WSIZE))
    smem[threadIdx.x + WSIZE][threadIdx.y]=kcell[idx+WSIZE][idy];
  if ((threadIdx.x >(N-WSIZE)) && (threadIdx.y > (N-WSIZE)))
   smem[threadIdx.x+WSIZE][threadIdx.y+WSIZE]=
              kcell[idx+WSIZE][idy+WSIZE];
//wait for all threads to finish read
  __syncthreads();
```

```
//find minimum in window
  temp = knode[idx][idy];
  for(x=0;x<WSIZE;x++)
   for(y=0;y<WSIZE;y++)
     if (temp> smem[threadIdx.x+x][threadIdx.y+y])
       temp = smem[threadIdx.x+x][threadIdx.y+y];
     knode[i][j] = temp;
  }
}
```

**Write Minimum to Global Memory**

**Questions?**