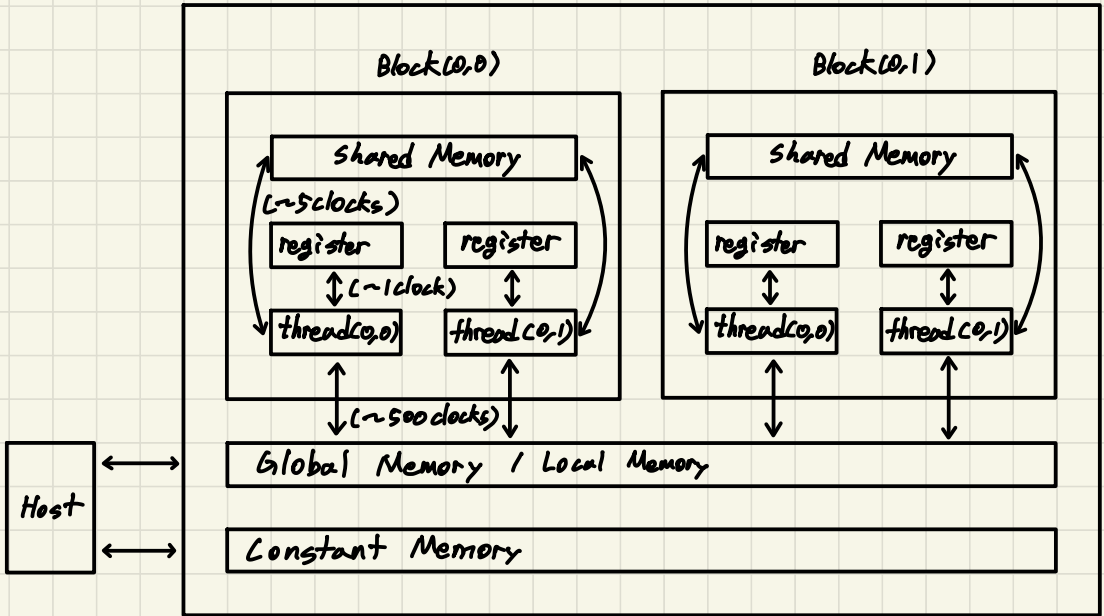


CUDA memory hierarchy

Grid



1. max 1024 threads per 1 block
2. 1 block per sm
3. (32 threads = 1 warp) per sm

```
--global-- void kernelFunc(float* dst, const float* src) {
    float p = src[threadIdx.x]; // register per thread
    float heap[10]; // local per thread
    --shared-- float partial-sum[1024]; // shared Memory per Block
}
```

TILED MATRIX MULTIPLICATION Basic

rhs matrix

thread (0,0)	thread (0,1)						
thread (1,0)	thread (1,1)						
thread (0,0)	thread (0,1)						
thread (1,0)	thread (1,1)						
thread (0,0)	thread (0,1)						
thread (1,0)	thread (1,1)						

lhs matrix

thread (0,0)	thread (0,1)	thread (0,0)	thread (0,1)	thread (0,0)	thread (0,1)
thread (1,0)	thread (1,1)	thread (1,0)	thread (1,1)	thread (1,0)	thread (1,1)

blockIdx(0,0) blockIdx(0,1)

blockDim.x

result matrix

blockDim.y

thread (0,0)	thread (0,1)						
thread (1,0)	thread (1,1)						

cols

rows

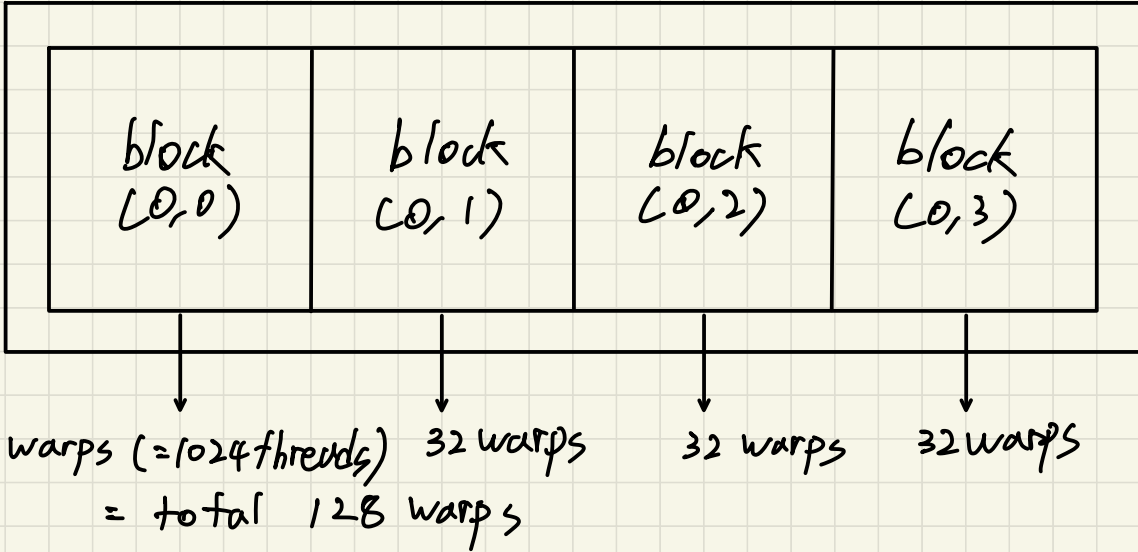
pseudo code

```

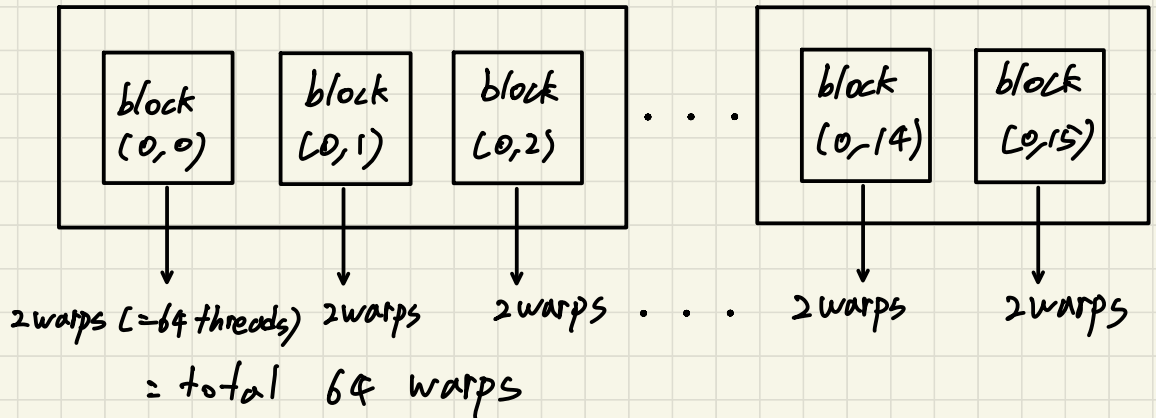
int x = blockIdx.x * blockDim.x + threadIdx.x;
int y = blockIdx.y * blockDim.y + threadIdx.y;
double sum = 0.0;
for (int i = 0; i < lhs_cols; i++)
{
    sum += lhs[y * lhs_cols + i] * rhs[i * rhs_cols + x];
}
result[y * rows + x] = sum;
    
```

Result matrix dimension 256×1
TILED MATRIX MULTIPLICATION

TILE WIDTH = 32



TILE WIDTH = 8



TILED MATRIX MULTIPLICATION

using shared memory (TILE WIDTH = 2)

