# CUDA memory hierarchy

Grid

Block(0,0)

Shared Memory

(~5clocks)

register | register

↕(~1clock) ↕

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

Block(0,1)

Shared Memory

register | register

↕ ↕

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

↓(~500 clocks)↓

Global Memory / Local Memory

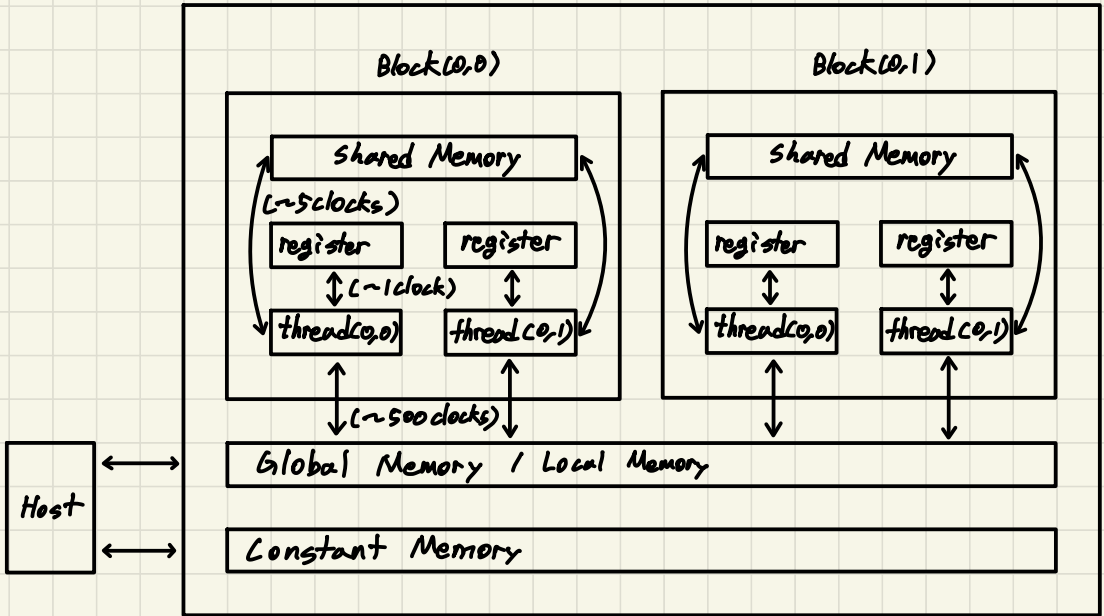Constant Memory

Host

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 = 0; // register per thread
    float arr[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) | | | | | |

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

blockDim.x

**lhs matrix**

**result matrix**

blockDim.y

cols

rows

## pseudo code

int x = blockIdx.x × blockDim.x + threadIdx.x

int y = blockIdx.y × blockDim.y + threadIdx.y

result[y, x] = sum ( lhs[y, :] × rhs[:, x])

# Calculate the number of warp
## (result matrix dimension 256 x 1)

## TILE WIDTH = 32

| block (0,0) | block (0,1) | block (0,2) | block (0,3) |
|---|---|---|---|

32 warps (=1024 threads)   32 warps      32 warps      32 warps

= total 128 warps

## TILE WIDTH = 8

| block (0,0) | block (0,1) | block (0,2) | . . . | block (0,14) | block (0,15) |
|---|---|---|---|---|---|

2 warps (=64 threads)  2 warps      2 warps . . .   2 warps       2 warps

= total 64 warps

# TILED MATRIX MULTIPLICATION
## using shared memory ( TILE WIDTH = 2 )

rhs_mat (32×32)

$N_{00}$ $N_{01}$
$N_{10}$ $N_{11}$
$N_{20}$ $N_{21}$
$N_{30}$ $N_{31}$
$N_{40}$ $N_{41}$
$N_{50}$ $N_{51}$
$N_{60}$ $N_{61}$
$N_{70}$ $N_{71}$

lhs_mat ( 16×32 )

$M_{00}$ $M_{01}$ $M_{02}$ $M_{03}$ $M_{04}$ $M_{05}$ $M_{06}$ $M_{07}$
$M_{10}$ $M_{11}$ $M_{12}$ $M_{13}$ $M_{14}$ $M_{15}$ $M_{16}$ $M_{17}$

$P_{00}$ $P_{01}$
$P_{10}$ $P_{11}$

) TILE WIDTH

result mat (16×32)

**step 1**
copy to shared Memory
( using _sync threads)

$M_{00}$ $M_{01}$
$M_{10}$ $M_{11}$

$N_{00}$ $N_{01}$
$N_{10}$ $N_{11}$

. . .

**step 4**
copy to shared Memory
( using _sync threads)

$M_{06}$ $M_{07}$
$M_{16}$ $M_{17}$

$N_{60}$ $N_{61}$
$N_{70}$ $N_{71}$

. . .

**calculation**

$P_{00} \pm M_{00} \times N_{00} + M_{01} \times N_{10}$
$P_{01} \pm M_{00} \times N_{01} + M_{01} \times N_{11}$
$P_{10} \pm M_{10} \times N_{00} + M_{11} \times N_{10}$
$P_{11} \pm M_{10} \times N_{01} + M_{11} \times N_{11}$

**calculation**

$P_{00} \pm M_{06} \times N_{60} + M_{07} \times N_{70}$
$P_{01} \pm M_{06} \times N_{61} + M_{07} \times N_{71}$
$P_{10} \pm M_{16} \times N_{60} + M_{17} \times N_{70}$
$P_{11} \pm M_{16} \times N_{61} + M_{17} \times N_{71}$