



vect-acld OK  
stencil-1D OK  
matrix mul OK  
matrix transpose OK

Stencil Shared memory

```
-- shared -- int Smem[1024];  
int index = BlockDim.x * BlockId.x + threadIdx.x  
if index < N :  
    Smem[threadIdx.x + Radius] = input[index + Radius]  
else  
    Smem[threadIdx.x + Radius] = 0  
-- syncthreads()  
if threadIdx.x < Radius  
    Smem[threadIdx.x] = input[index]  
if threadIdx.x > blockDim.x - Radius  
    Smem[threadIdx.x + Radius*2] = input[index + Radius*2]  
-- syncthreads()  
int sum = 0  
for int i = 0 ; i < Radius*2 ; i++  
    sum += Smem[threadIdx.x + i]  
output[index + Radius] = sum  
if index <= N || index >= N - Radius  
    output[index] = 0
```

## Vec Add

```
int global_id = blockId.x * blockDim.x + threadIdx.x
```

```
if global_id < N
```

```
output[global_id] = vec_A[global_id] + vec_B[global_id]
```

## Stencil

```
int global_id = blockId.x * blockDim.x + threadIdx.x
```

```
if global_id >= Radius || global_id < N - Radius  
    int sum = 0
```

```
for int i = -Radius, i < Radius; i++
```

```
    sum += input[global_id - i]
```

```
output[global_id] = sum
```

## Matrix Mul

```
int Row = BlockId.y * blockDim.y + threadIdx.y
```

```
int Col = blockIdx.x * blockDim.x + threadIdx.x
```

```
int Pvalue = 0
```

```
for int K = 0; K < N; K++
```

```
Pvalue += matrix_A[Row * N + K] * matrix_B[Col + K * N]
```

```
output[Row * N + Col] = Pvalue
```

## matrix transpose

```
int Row = blockIdx.y * blockDim.y + threadIdx.y
```

```
int Col = blockIdx.x * blockDim.x + threadIdx.x
```

```
if Row < N && Col < N
```

```
res[Row * N + Col] = input[Col * N + Row]
```

-- Shared -- int sMem[1024] ✓

int index = BlockId.x \* blockDim.x + threadIdx.x ✓

if index < N :

sMem[threadId.x + Radius] = input[index + Radius] ✓

else :

sMem[threadId.x + Radius] = 0 ✓

-- sync threads () ✓

if threadIdx.x < Radius : ✓

sMem[threadId.x] = input[index] ✓

if threadIdx.x >= blockDim.x - Radius ✓

sMem[threadId.x + 2\*Radius] = input[index + 2\*Radius] ✓

-- sync threads () ✓

int sum = 0 ✓

for int i = 0 ; i < Radius\*2 ; i++ ✓

sum += sMem[threadId.x + i] ✓

output[index + Radius] = sum ✗

if index > N - Radius  $\wedge$  index < N ✓

output[index] = 0 ✓

## shm Matrix Mul

-- Shared -- int smem A [TILE][TILE]

-- Shared -- int smem B [TILE][TILE]

int bx = Block Id. x      int by = Block Id. y

int tx = Thread Id. x      int ty = Thread Id. y

int Raw = by \* TLE + ty

int Col = bx \* TLE + tx

int Pvalue = 0

for int m = 0, m < width/tile, width++ {

    smem A [ty][tx] = input A [Raw\*N + m\*tile + tx]

    smem B [ty][tx] = input B [Col + (m\*tile+ty)\*N]

-- syncthreads()

    for int k = 0, k < Tile, k++ {

        Pvalue += smem A [ty][k] \* smem B [k][tx]

    } -- syncthreads()

}

output [Raw\*N + Col] = Pvalue

## shmem Matrix Transpose

```
-- shmem -- int smem [BlockDim][BlockDim]  
int Row = Block Id.x * Block Dim.x + thread Id.x  
int Col = Block Id.y * Block Dim.y + thread Id.y  
if Row < N & Col < N:  
    smem [thread Id.y][thread Id.x] = input [Row * N + Col]  
-- syncthreads()  
Col = block Id.x * Block Dim.x + thread Id.x  
Row = block Id.y * Block Dim.y + thread Id.y  
d-matrix-out [Row * N + Col] = smem [thread Id.x][thread Id.y]
```

-- shared -- int memA [TILE][TILE]

-- shared -- int memB [TILE][TILE]

int ty = threadIdx.y; int tx = threadIdx.x;

int by = BlockId.y; int bx = BlockId.x;

int Col = bx \* TILE .x + tx;

int Row = by \* TILE + ty;

int value = 0

for (int k=0; k < width/TILE; k++) {

smemA[ty][tx] = matA [Row\*N + k\*TILE + tx]

smemB[ty][tx] = matB [Col + (k\*TILE + ty)\*N]

-- syncthreads()

for (int m=0; m < TILE; m++) {

value += smemA[ty][m] \* smemB[m][tx]

}

-- syncthreads()

}

output [Row\*N + Col] = value

-- shared -- int smem [1024]

int global\_id = BlockId.x \* BlockDim.x + threadId.x

smem [threadId.x] = input [global\_id]

-- syncthreads ()

for (int i = 1; i <= BlockDim.x; i++) {

index = threadId.x \* i \* 2

if (index < BlockDim.x) {

smem [threadId.x] += smem [threadId.x + i]

}

-- syncthreads ()

}

if (threadId.x == 0)

output [BlockId.x] = smem [0]

}

## Esercizio 2

Abbiamo da fare il confronto fra due architetture hardware:

- Shared memory: è un'architettura hardware che dispone di una memoria singola accessibile da tutti i processori qualiasi sia la sua locazione fisica. I vantaggi della shared memory sono la rapidità di accesso ai dati vista la vicinanza di questa memoria ai processori e la semplicità di programmazione infatti tutto lo spazio degli indirizzi è condiviso e accessibile da qualsiasi processore. Gli svantaggi invece sono la scalabilità, infatti con l'aumentare del numero di processori il traffico per accedere alla memoria aumenta. Inoltre possono peggiorare le performance per gestire le race condition ed è difficile per il programmatore che ne è responsabile mantenere un eventuale code coherency.
- Distributed memory: è un'architettura hardware in cui ciascun processore dispone della propria memoria locale e può accedere a quelle degli altri processori comunicando con essi seguendo le varie opportunità messe a disposizione. I vantaggi sono la possibilità di scalare in maniera semplice poiché basta aggiungere un processore con la sua memoria. Non è inoltre necessario gestire la code coherency in quanto ogni processore è indipendente. Gli svantaggi sono la difficoltà di trasferire queste strutture in strutture adatte alla shared memory. Inoltre il programmatore ha tutte le responsabilità di gestire la comunicazione tra i vari processori per raggiungere le altre memorie.

## Esercizio 3

$$\text{speedup} = \frac{1}{(1-x) + \frac{x}{100}}$$

$$80 = \frac{1}{(1-x) + \frac{x}{100}}$$

$$80 \left( (1-x) + \frac{x}{100} \right) = 1$$

$$80 - 80x + \frac{80}{100}x = 1$$

$$\frac{396}{5}x = 79$$

$$x = 0,997$$

99.7% parallelizzabile

0.3% per rimanere sequenziale



-- shared \_\_ int snum [0xa]

int ind = BlockId.x \* blockDim.x + threadIdx.x

if (ind < BlockId.x - 1) {

snum [threadId.x + 1] = input [ind]

}

snum [0] = 0

}

-- syncthreads ()

for (offset = 1 ; offset < blockDim.x ; offset += 2) {

if (threadId.x >= offset) {

LeftValue = snum [threadId.x - offset]

RightValue = snum [threadId.x]

}

-- syncthreads ()

if (threadId.x >= offset) {

snum [threadId.x] = LeftValue + RightValue

-- syncthreads

}

output [ind] = snum [threadId.x]

-- shared \_\_ int smem [1024]

int ind = BlockId.x \* blockDim.x + threadIdx.x

smem [threadIdx.x] = input [ind]

-- syncThreads ()

for (int i = 1; i < blockDim.x; i++) {

int offset = threadIdx.x \* (i \* 2)

if (offset < blockDim.x)

smem [offset] += smem [offset + i]

-- syncThreads ()

}

if (threadIdx.x == 0)

output [BlockId.x] = smem [0]

-- shared \_\_ int smem [1024]

int ind = BlockId.x \* blockDim.x + threadIdx.x

if (threadIdx.x < blockDim.x - 1)

smem[threadIdx.x+1] = input[ind]

else

smem[0] = 0

-- syncthreads()

for (int offset=1 ; offset < blockDim.x ; offset+=2) {

int LeftValue; RightValue

if (threadIdx.x > offset)

LeftValue = smem[threadIdx.x - offset]

RightValue = smem[threadIdx.x]

-- syncthreads()

if (threadIdx.x > offset)

smem[threadIdx.x] = LeftValue + RightValue

-- syncthreads

}

output[ind] = smem[threadIdx.x]

-- shared -- int smem [1024]

int ind = BlockId.x \* blockDim.x + threadIdx.x

if threadIdx.x < blockDim.x - 1

smem[threadIdx.x + 1] = input[index]

else

smem[0] = 0

-- syncThreads()

for (int offset = 1; offset < blockDim.x; offset \*= 2) {

int leftvalue; int rightvalue

if (threadIdx.x >= offset)

leftvalue = smem[threadIdx.x - offset]

rightvalue = smem[threadIdx.x]

-- syncThreads()

if (threadIdx.x >= offset)

smem[threadIdx.x] = Leftvalue + Rightvalue

-- syncThreads()

}

output[ind] = smem[threadIdx.x]

3

$$\text{Speedup} = \frac{\frac{1}{1 - \frac{\text{Fracparallel}}{\text{Fracparallel}} + \frac{\text{Fracparallel}}{\text{Speedupparallel}}}}{\text{Speedupparallel}}$$

$$80 = \frac{1}{(1-x) + \frac{x}{100}}$$

$$80 - 80x + \frac{80}{100}x = 1$$

$$79 = 80x - \frac{80}{100}$$

$$79 = \frac{8000 - 80x}{100}$$

$$79 = \frac{7920}{100}x$$

$$x = 0,997$$

Il 99.7% del codice deve essere parallelo  
Lo 0.3% del codice deve essere sequenziale

4

1024 thread  $\rightarrow$  32 warp

1° iterazione  $\rightarrow$  512 thread lavorano

2° iterazione  $\rightarrow$  256 thread ..

3° ..  $\rightarrow$  128 .. ..

4° ..  $\rightarrow$  64 .. ..

5° ..  $\rightarrow$  32 .. ..

Alla 5° iterazione lavora solo con warp di thread

Le altre 31 sono ferme

10

Senza shared

$$\frac{\# \text{ computation}}{\# \text{ communication}} = \frac{1}{2} = \frac{1}{2}$$

Con shared

$$\frac{\# \text{ computation}}{\# \text{ communication}} = \frac{n\text{thread}^* 1}{n\text{thread}^* 2} = \frac{1}{2}$$

Senza shared

$$\frac{\# \text{ computation}}{\# \text{ communication}} = \frac{2}{2} = 1$$

Con shared

$$\frac{\# \text{ computation}}{\# \text{ communication}} = \frac{256 * (2 * 16)}{256 * 2} = 16$$

↳ riduzione di bandwidth 16

①  $\frac{4,90}{1,43} = 3,42$

②  $\frac{5,63}{1,35} = 4,11$

③  $\frac{1,43}{0,95} = 2,06$        $\frac{1,39}{0,96} = 1,40$

$$\frac{1}{1 - 0,16 + \frac{0,16}{10}} = 1,96 \text{ speedup}$$

$$\frac{1}{0,5 + \frac{0,5}{5}} = 1,66 \text{ speedup} = \frac{5}{3} \text{ speedup}$$

$$\text{cost} = \frac{5}{3} + \frac{2}{3} = \frac{7}{3} \text{ costo}$$

Ho più costo che speedup introdotto, non ne vale la pena

$$\frac{1}{1-x + \frac{x}{100}} = 80$$

$$80 - 80x + \frac{80}{100}x = 1$$

$$\frac{8000 - 80}{100} x = 79$$

$$x = \frac{79}{79,2} = 0,992$$

$$x = 0,03$$

-- shared -- int smem [1024]

int index = BlockId.x \* blockDim.x + threadIdx.x

if (threadIdx.x < blockDim.x - 1)

smem[threadIdx.x + 1] = input[index]

else

smem[0] = 0

-- syncthreads()

for (int offset = 1; offset < blockDim.x; offset += 2) {

int Leftvalue, Rightvalue;

if (threadIdx.x > offset)

Leftvalue = smem[threadIdx.x - offset]

Rightvalue = smem[threadIdx.x]

-- syncthreads()

if (threadIdx.x > offset)

smem[threadIdx.x] = Leftvalue + Rightvalue

-- syncthreads()

}

output[index] = smem[threadIdx.x]

-- shared -- int smem [1024]

int ind = BlockId.x \* blockDim.x + threadIdx

smem[threadIdx] = input[ind]

-- syncthreads()

int step=1

for (int limit = blockDim.x/2; limit>0; limit/=2) {

if (threadIdx < limit) {

RightValue = (threadIdx + 1) \* (step\*2) - 1

LeftValue = RightValue - step

smem[RightValue] += smem[LeftValue]

}

step \*= 2

-- syncthreads()

}

if (threadIdx == 0)

smem[BlockDim.x - 1] = 0

-- syncthreads()

step = blockDim.x / 2

for (int limit = 1; limit < blockDim.x; limit\*=2) {

RightValue = (threadIdx+1) \* (step\*2) - 1

LeftValue = RightValue - step

int tmp = smem[LeftValue]

smem[LeftValue] = smem[RightValue]

smem[RightValue] += tmp

}

step /= 2

-- Syncthreads()

}

output[ind] = smem[threadIdx]

-- Shared \_\_ int smem [cozu]

int index = BlockId.x \* blockDim.x + threadIdx.x

smem [threadIdx.x] = input [index]

-- syncThreads()

int step = 1

for (int limit = BlockId.x / 2; limit > 0; limit /= 2) {

if (threadIdx.x < limit)

Rightvalue = (threadIdx.x + 1) \* (step \* 2) - 1

Leftvalue = Rightvalue - step

smem [Rightvalue] += smem [Leftvalue]

step \*= 2

-- syncThreads()

}

if (threadIdx.x == 0)

smem [BlockDim.x - 1] = 0

-- syncThreads ()

int step = blockDim / 2

for (int limit = 1; limit < blockDim / 2; limit \*= 2) {

if (threadIdx.x < limit)

Rightvalue = (threadIdx.x + 1) \* (i \* 2) - 1

Leftvalue = Rightvalue - step

int temp = smem [Leftvalue]

smem [Leftvalue] = smem [Rightvalue]

smem [Rightvalue] += temp

step \*= 2

-- syncThreads()

}

output [index] = input [threadIdx.x]

-- shared -- int smem[1024]

int index = BlockId.x \* blockDim.x + threadIdx.x

smem[threadIdx.x] = input[index]

-- syncthreads()

int step = 8

for (int limit = blockDim.x / 2 ; limit > 0 ; limit /= 2) {

if (threadIdx.x < limit)

Rightvalue = (threadIdx.x + 1) \* (i \* 2) - 1

Leftvalue = Rightvalue - step

smem[Rightvalue] += smem[Leftvalue]

step \*= 2

--syncthreads()

}

if (threadIdx.x == 0)

smem[BlockId.x - 1] = 0

-- syncthreads

int step = blockDim.x / 2

for (int limit = 1 ; limit <= blockDim.x ; limit \*= 2) {

if (threadIdx.x < limit)

Rightvalue = (threadIdx.x + 1) \* (i \* 2) - 1

Leftvalue = Rightvalue - step

int tmp = smem[Leftvalue]

smem[Leftvalue] = smem[Rightvalue]

smem[Rightvalue] += tmp

step /= 2

-- syncthreads()

}

output[index] = smem[threadIdx.x]

- - shared -- int smem [1024]

int index = BlockId.x \* blockDim.x + threadIdx.x

smem [threadId.x] = input [index]

-- syncthreads()

int step = 1

for (int limit = BlockId.x / 2; limit > 0; limit /= 2) {

if (threadId.x < limit)

int Rightvalue = (threadId.x + 1) \* (i \* 2) - 1

int Leftvalue = Rightvalue - step

smem [Rightvalue] += smem [Leftvalue]

step \*= 2

-- syncthreads()

}

if (threadId.x == 0)

smem [BlockDim.x - 1] = 0

-- syncthreads

int step = blockDim.x / 2

for (int limit = 1; limit <= blockDim.x / 2; limit \*= 2)

if (limit < step)

int Rightvalue = (threadId.x + 1) \* (i \* 2) - 1

int Leftvalue = Rightvalue - step

int tmp = smem [Leftvalue]

smem [Leftvalue] = smem [Rightvalue]

smem [Rightvalue] += tmp

step /= 2

-- syncthreads()

}

output [index] = smem [threadId.x]

$$y = \frac{1}{(1-x) + \frac{x}{20}}$$

$$z = \frac{1}{(1-x) + \frac{x}{20}}$$

$$2 - 2x + \frac{2x}{20} = 1$$

$$1 = 2x - \frac{2x}{20}$$

$$1 = \frac{40 - 2x}{20}$$

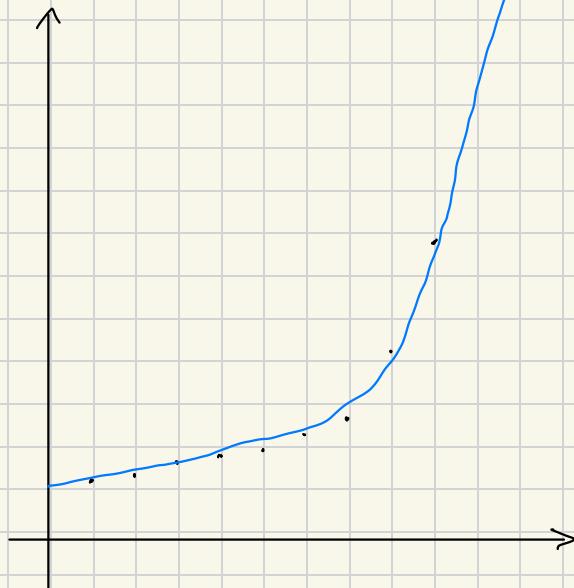
$$1 = \frac{38}{20}x$$

$$x = \frac{20}{38}$$

$$x = 0,52$$

$$y = \frac{1}{(1-x) + \frac{x}{20}}$$

$$y = \frac{1}{(1-x) + \frac{x}{20}} = \frac{1}{1 - \frac{19}{20}x}$$



$$10 = \frac{1}{(1-x) + \frac{x}{20}}$$

$$1 = 10 - 10x + \frac{x}{2}$$

$$9 = \frac{19}{2}x$$

$$x = \frac{9 \cdot 2}{19} = 0,95$$

$$x = \frac{1}{(1-0,7) + \frac{0,7}{40}}$$

$$x = \frac{1}{0,3 + \frac{0,7}{40}} = 3,14$$

$$3,14 = \frac{1}{(1-x) + \frac{x}{20}}$$

$$3,14 - 3,14x + \frac{3,14}{20}x = 1$$

$$2,14 = 2,98x$$

$$x = 0,718$$

$$\frac{1}{1-0,9 + \frac{0,9}{5}} = 3,57$$

$$\frac{1}{1-\underline{1} + \frac{1}{\underline{1}}} = 1 \quad \text{Speedup } 3,52$$

$$CPI_{\text{org}} = 0,25 * 4 + 0,75 * 1,33 = 1 + 1 = 2$$

$$CPI_{\text{mod}_1} = 2 - 0,02 (20 - 2) = 1,64$$

$$CPI_{\text{mod}_2} = 2 - 0,25 (4 - 2,5) = 1,625$$

$\nearrow$   
meglio

$$0,43 * 1 + 0,21 * 2 + 0,12 * 2 + 0,24 * 3 = 1,57$$

$$\frac{(0,43 - (0,43 - 0,25)) * 1 + (0,21 - (0,43 + 0,25)) * 2 + (0,12 * 0,25) * 2 + 0,12 * 2 + 0,24 * 3}{1 - (0,43 * 0,25)}$$

$$= 1,908$$

$$CP_{\text{orig}} = 1 * 0,3 + 1 * 0,7 = 1$$

$$CP_{\text{new}} = 1,05 * 0,1 + 1 * 0,7 = 0,805$$

↳ ottimizzato

-- shared - int smem [1024]

int index = BlockId.x \* blockDim.x + threadIdx.x

smem[threadIdx.x] = input[index]

-- syncthreads()

int step = 1

for (int limit = blockDim.x / 2 ; limit > 0 ; limit /= 2) {

if (threadIdx.x < limit)

int Rightvalue = (threadIdx.x + 1) \* (step \* 2) - 1

int Leftvalue = Rightvalue - step

smem[Rightvalue] += smem[Leftvalue]

Step \*= 2

-- syncthreads()

}

if threadIdx.x == 0

smem[BlockDim.x - 1] = 0

-- syncthreads()

int step = blockDim.x / 2

for (int limit = 1 ; limit <= blockDim.x ; limit \*= 2)

if (threadIdx.x < limit)

int Rightvalue = (threadIdx.x + 1) \* (step \* 2) - 1

int Leftvalue = Rightvalue - step

int tmp = smem[Leftvalue]

smem[Leftvalue] = smem[Rightvalue]

smem[Rightvalue] += tmp

Step /= 2

-- syncthreads()

3

output[index] = smem[threadIdx.x]

$$\frac{\# \text{ computation}}{\# \text{ communication}} = \frac{36}{7}$$

$$\text{GFLOPS} = \frac{100}{4} \cdot \frac{36}{7} = 128.5 \text{ Gflops} < 200 \text{ GFlops di picco} \quad \leftarrow \text{memory bound, la computazione è limitata dai trasferimenti di memoria}$$

$$\text{G-FLOPs} = \frac{2500}{4} \cdot \frac{36}{7} = 3214 \text{ Gflops} > 300 \text{ GFlops di picco} \quad \leftarrow \text{compute bound, la computazione è limitata dalla potenza del Kernel}$$

$$\frac{109970178}{0,094 \cdot 10^9} = 1,17 \text{ GFLOPS}$$

$$3 \quad 80 = \frac{1}{1-x + \frac{x}{100}}$$

$$80 - 80x + \frac{80}{100}x = 1$$

$$79 = \frac{7920}{100}x$$

$$x = \frac{79}{792} = 0,997$$

$x = 99,7\%$ . codice necessariamente parallelo

codice sequenziale 0,3 %

4

$$\text{CPI}_{avg} = 0,25 * 4 + 0,75 * 1,33 \approx 2$$

$$\text{CPI}_1 = 2 - 0,02(20-2) = 1,64$$

$$\text{CPI}_2 = 2 - 0,25(4-2,5) = 1,625$$

meglio secondo ho CPI più basso

Usando l'algoritmo di reduce assumendo che abbia la miglior implementazione per evitare la divergenza in un warp ho:

1024 thread all'inizio di cui solo 512 sopravvivono il controllo

step 2      512  $\rightarrow$  256

256  $\rightarrow$  128

128  $\rightarrow$  64

64  $\rightarrow$  32

perciò alla 5° iterazione avrà un warp attivo e 31 warps inattivi

-- shared -- int smem[1024]

int index = BlockIdx.x \* blockDim.x + threadIdx.x

smem[threadIdx.x] = input[index]

-- syncthreads()

int step = 1

for (int limit = blockDim.x / 2; limit > 0; limit /= 2) {

if (threadIdx.x < limit) {

int Rightvalue = (threadIdx.x + 1) \* (step \* 2) - 1

int Leftvalue = Rightvalue - step

smem[Rightvalue] += smem[Leftvalue]

}

step \*= 2

-- syncthreads()

}

if (threadIdx.x == 0)

smem[BlockDim.x - 1] = 0

-- syncthreads()

int step = blockDim.x / 2

for (int limit = 1; limit < blockDim.x / 2; limit \*= 2) {

if (threadIdx.x < limit) {

int Rightvalue = (threadIdx.x + 1) \* (step \* 2) - 1

int Leftvalue = Rightvalue - step

int tmp = smem[Leftvalue]

smem[Leftvalue] = smem[Rightvalue]

smem[Rightvalue] += tmp

}

step /= 2

-- syncthreads()

}

output[index] = smem[threadIdx.x]

3.

$$80 = \frac{x}{(1-x) + \frac{x}{100}}$$

$$80 - 80x + \frac{80}{100}x = 1$$

$$\frac{7920}{100}x = 79$$

$$x = \frac{79}{79,2}$$

$x = 0,997$  codice parallelo

0,3% codice che può rimanere seq

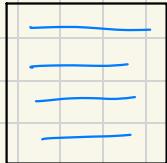
4

$$CPI_{\text{orig}} = 0,25 * 4 + 0,75 * 1,33 = 1 + 1 = 2$$

$$CPI_{\text{mod}} = 2 - 0,02(20 - 2) = 1,64$$

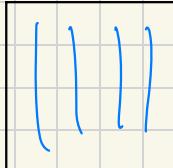
$$CPI_{\text{mod}} = 2 - 0,25(4 - 2,5) = 1,625 \leftarrow \text{è la soluzione migliore}$$

M

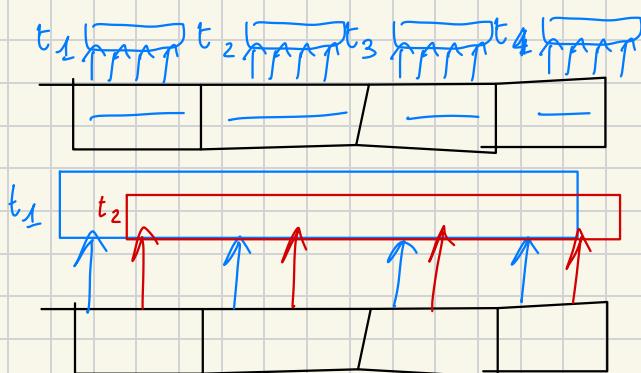


X

N



Raw Major Layout



8 Utilizziamo un algoritmo in cui la warp divergence è ridotta al minimo

1024 di cui 512 buonano alla prima iterazione

2° 256    "    "    "

3° 128    "    "    "

4° 64    "    "    "

5° 32    "    "    "

1 warp efficace alla 5° iterazione

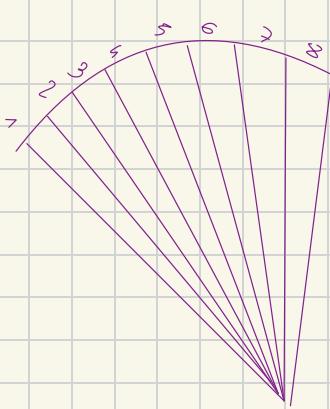
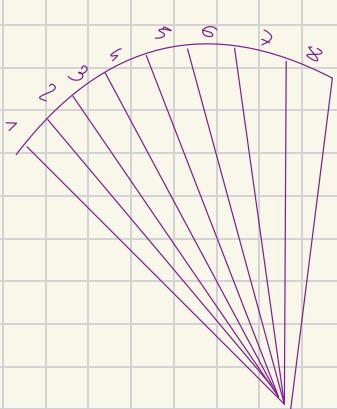
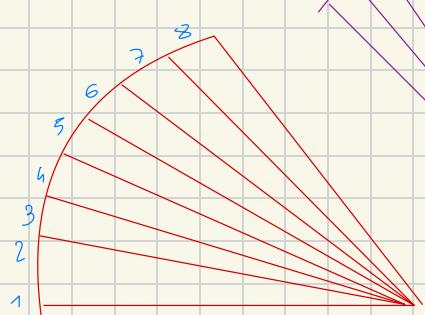
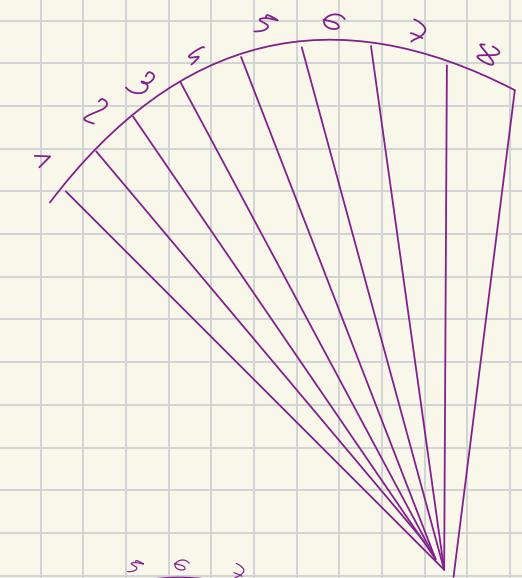
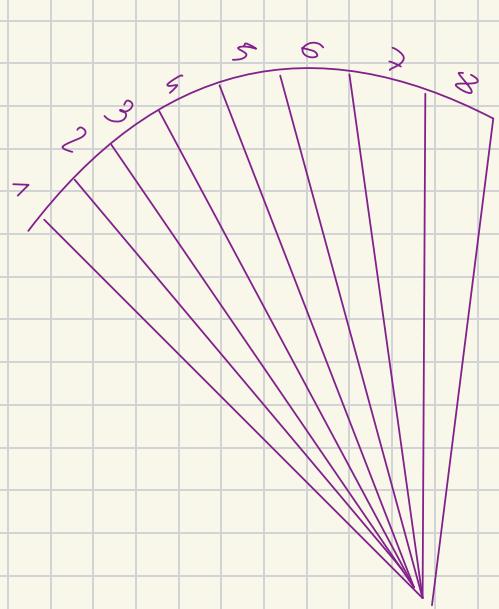
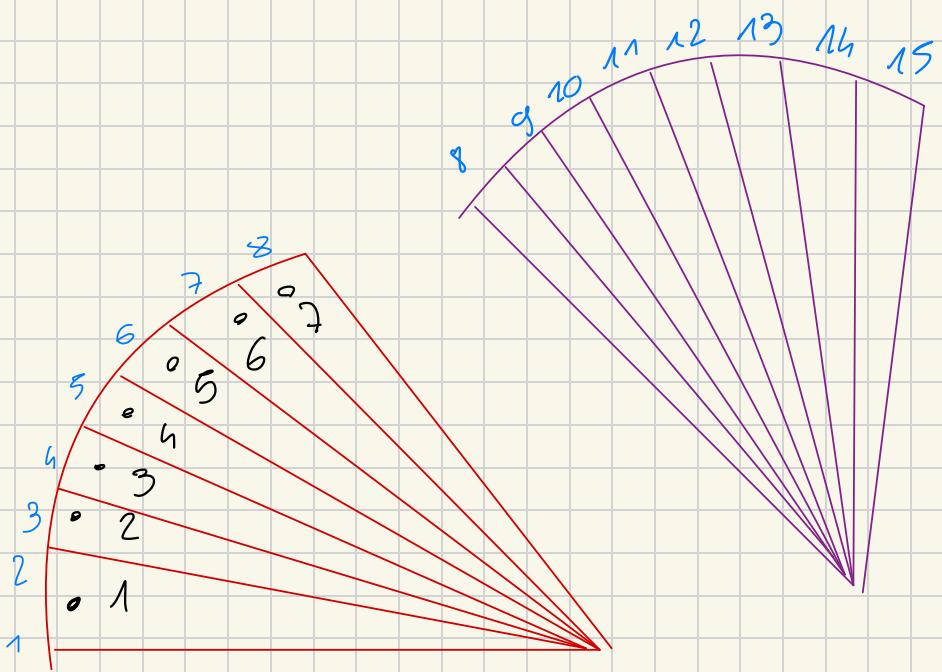
3-1 inefficiente

9

$$\frac{\# \text{ computation}}{\# \text{ communication}} = \frac{36}{7}$$

GFLOPS =  $\frac{100}{4} \cdot \frac{36}{7} = 128,5 < 200$  GFLOPS di picco perciò la computazione è memory bound perché limitata dai trasferimenti da memoria

GFLOPS =  $\frac{2500}{4} \cdot \frac{36}{7} = 3214 > 300$  Gflops di picco perciò la computazione è kernel bound perché limitata dalla capacità di computazione della GPU



3 Mat mul ✓

8 Reduce Norm ✓

10 Prefix Noise ✓

9 Reduce Nb branch div ✓

5 Stencil 1D Shmem ✓

7 Transp Shmem ✓

2 Stencil ✓

11 Prefix Up/Down ✓

1 Summa ✓

4 Mat Transp ✓

6 Mat Shmem

```
-- shared -- int smemA[TILE][TILE]  
-- shared -- int smemB[TILE][TILE]
```

```
int ty = threadIdx.y; int tx = threadIdx.x  
int bx = BlockId.x; int by = BlockId.y
```

```
int Raw = by * TILE + ty
```

```
int Col = bx * TILE + tx
```

```
int Pwhe = 0
```

```
for (int m = 0; m < width/TILE; m++)
```

```
smemA[ty][tx] = input[Raw*N + m*TILE+tx]
```

```
smemB[ty][tx] = input[Col + (m*TILE+ty)*N]
```

```
-- syncThreads()
```

```
for (int i = 0; i < TILE; i++)
```

```
Pwhe += smemA[ty][i] * smemB[i][tx]
```

```
-- syncThreads()
```

```
}
```

```
output[Raw*N+Col] = Pwhe
```

2 8, 9, 10, 11, 12, 13, 14, 15  
16, 17, 18, 19, 20, 21, 22, 23

---

8 72

16 80

24 88

32

40

48

56

64

---

0

64

128

:

:

1

2

3

4

5

6

7

8

$$\frac{\# \text{ computation}}{\# \text{ communication}} = \frac{2}{2} = 1$$

$$\frac{\# \text{ computation}}{\# \text{ communication}} = \frac{\# \text{ num thread } (16 * 2)}{\# \text{ num threads } * 2} \rightarrow \begin{array}{l} \text{ogni thread fa} \\ 16 * 2 \text{ operazioni} \end{array}$$

= 16

ogni thread fa 2 operazioni  
di comunicazione

---

Esempio 0 slide

(a)  $\frac{4.90}{1.63} = 3.02$

(b)  $\frac{5,69}{1,35} = 4.21$

(c)  $\frac{1,43}{0,695} = 2,05 \quad \frac{1,35}{0,96} = 1,40$

---

1

$$\text{Speedup} = \frac{1}{1 - 0,4 + \frac{0,4}{10}} = 1,56$$


---

2

$$\text{speedup} = \frac{1}{0,5 + \frac{0,5}{5}} = 1,66 = \frac{5}{3}$$

$$\text{Costo} = 5 \cdot \frac{1}{3} + \frac{2}{3} = \frac{7}{3} = 2.33$$

non ne vale la pena, il costo è maggiore dello speedup

---

3

$$80 = \frac{1}{(1-x) + \frac{x}{100}}$$

$$80 - 80x + \frac{80x}{100} = 1$$

$$\frac{7920}{100}x = 79$$

$$792x = 79$$

$$x = 0,997$$

$$99.7\%$$

dove essere parallelo, 0,3% sequenziale

(4)

$$\text{Speedup} = \frac{1}{(1-x) + \frac{x}{20}} = \frac{1}{\frac{20 - 20x + x}{20}} = \frac{1}{1 - \frac{19}{20}x}$$

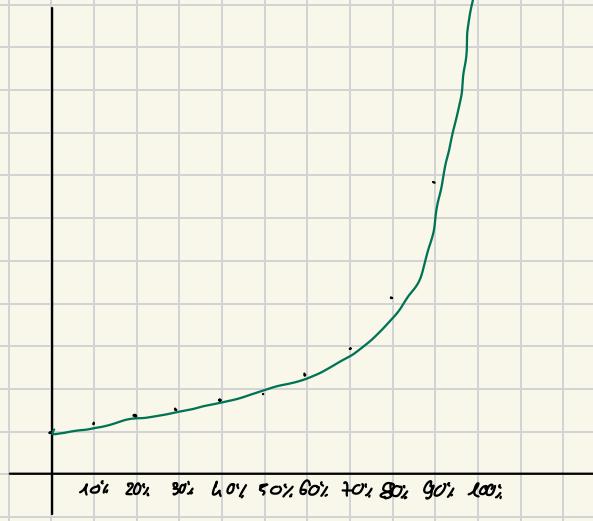
20

$$z = \frac{1}{1 - \frac{19}{20}x}$$

$$z - \frac{19}{10}x = 1$$

$$\frac{19}{10}x = 1$$

$$x = \frac{10}{19} = 0,526$$



$$10 = \frac{1}{1 - \frac{19}{20}x} = 10 - \frac{19}{2}x = 1$$

$$g = \frac{19}{2}x \quad x = 0,95$$

$\hookrightarrow 95\%$  deve essere parallel

4bis

$$\frac{1}{0,3 + \frac{0,17}{40}} = 3,149$$

$$3,149 = \frac{1}{1 - \frac{19}{20}x}$$

$$3,149 - 3,149 \cdot \frac{19}{20}x = 1$$

$$2,99x = 2,149$$

$$x = \frac{2,149}{2,99} = 0,718$$

71,8% di percentuale di parallelismo

1,8% in più non ha senso praticamente nulla

$$\frac{1}{0,1 + \frac{0,8}{5}} = 3,57$$

Non cache = 1

fu uno speed up di 3,57

$$CPI_{orig} = 0,25 * 4 + 0,75 * 1,33 = 2$$

$$CPI_{mod1} = 2 - 0,02(20 - 2) = 1,64$$

$$CPI_{mod2} = 2 - 0,25(4 - 2,5) = 1,625$$

La seconda è meglio perché il CPI è più basso e quindi ci vogliono meno cicli per fare un'istruzione

$$0,43 + 0,42 + 0,26 + 0,48 = 1,57$$

$$\frac{(0,43 - (0,43 \cdot 0,25))^* 1 + (0,21 - (0,43 \cdot 0,25))^* 2 + (0,43 \cdot 0,25)^* 2 + (0,12)^* 2}{0,26^* 3}$$
$$(1 - 0,43 \cdot 0,25)$$

$$\frac{0,3225 + 0,205 + 0,215 + 0,26 + 0,72}{0,8925} = 1,907$$

No il CPI è aumentato quindi modificare non comporta uno speedup né peggiora le performances

$$0,3 * 1 + 0,7 * 1 = 1$$

$$1,05 * 0,2 + 0,7 * 1 = 0,91$$

(a seconda richiesta sono CPI per far un  
tossi quindi è negl/ia)

$$\frac{10^9 \cdot 930178}{0,096 \cdot 10^9} = 1,17 \text{ GFLOPs}$$

-- shared \_\_ int smemA[TILE][TILE]  
-- shared \_\_ int smemB[TILE][TILE]

int ty = threadId.y; int tx = threadId.x  
int by = BlockId.y; int bx = BlockId.x

int Row = by \* TILE + ty  
int Col = bx \* TILE + tx

int Pvalue = 0

for (int m=0; m < Width/TILE; m++) {

smemA[ty][tx] = input[Row \* N + m \* TILE + tx]  
smemB[ty][tx] = input[Col + (m \* TILE + ty) \* N]

-- syncthreads()

for (int i=0; i < TILE; i++)

Pvalue += smemA[ty][i] \* smemB[i][tx]

-- syncthreads()

}

output[Row \* N + Col] = Pvalue

- 3 Mat mul ✓
- 8 Reduce Norm ✓
- 10 Prefix Noise ✓
- 9 Reduce No branch div ✓
- 5 Stencil 1D Shared ✓
- 7 Triang Shmem ✓
- 2 Stencil ✓
- 11 Prefix Up Down ✓
- 1 Scram ✓
- 4 Mat Transp ✓
- 6 Mat Shmem