

Se quiere implementar en CUDA una función de reducción para sumar N números.

Supremo $N = 8192$ elementos

Max 512 t/B • Max 8 B • 16 SM → 24 Warps

Diagrama de árbol de ejecución:

```
graph TD
    0 --> 1
    1 --> 2
    2 --> 3
    3 --> 4
    4 --> 5
    5 --> 6
    6 --> 7
    7 --> 13
    13 --> 22
    22 --> 28
```

- 28
- @ Indique una planificación de bloques y threads para realizar esta reducción, ¿cuántos warps se van a utilizar? ¿cuántos warps por SM se están planificando?
- 24 warps
- 30 warps = 128

ways for SM rectifier planification:

$$512 \text{ F/B} \cdot \frac{1 \text{ W}}{32 \text{ A}} = 16 \text{ W/B max} \quad 16 \frac{\text{W}}{\text{B}} \cdot 80 \text{ max} = 1280$$

2103 F 512 F

$$\frac{24 \text{ W}}{\text{SM}} \cdot \frac{32 \text{ T}}{\text{W}} = 768 \text{ T/SM}$$

$$\frac{512T}{80} = 64 \frac{T}{B} \rightarrow 2W/B$$

Se emplean 16 waps, 8 bloques, 64 T/B

- 6) se pretende realizar cambios optimizando la planificación de warps, ¿qué cambios realizarian para conseguirlo? ¿cuál es el menor número de SM para esa planificación?
- Podriamos dinamicamente unir los threads ejecutando instrucciones en los warps para reducir la divergencia en ejecución

$$\frac{8192 \text{ V}}{768 \text{ V/SM}} = 10.6 \rightarrow 11 \text{ SM minimum}$$

- ⑦ Explique el esquema de la reducción propuesta y comportamiento en eficiencia

-- shared -- float psum C3
thread id

unsigned int t = threadId * x; x

for (unsigned int stride = 1; stride < blockDim.x; stride += 2)

```
-- sync thread 1;
```

$$\sqrt{1 + (2 * stroke)} = 0$$

✓ $(t \times (2 * stride) == 0)$
 $pSum[t] += pSum[t + stride]$

Stride es la distancia entre elementos a sumar, pivote 1, luego 2, 4, 8... Cuando te encuentras en un elemento divisible entre el doble de distancia \rightarrow estar en el 1º de los elementos a sumar aunque p sum de la pos del elemento te sumas el de la pos + stride y así hasta que el stride supera blockDim (has sumado todos los elementos del bloque)

1 procesador escalar \rightarrow n sumas $T(n) = \lg(n)$ $E = \frac{T(1)}{T(n)} = \frac{n}{\lg(n)}$

$$T(n) = \lg(n) \quad E = \frac{T(1)}{T(n)} = \frac{n}{\lg(n)}$$

Proponer un esquema alternativo de reducción

```
-- shared -- float psum[3];
```

```
t = threadIdx.x
```

```
for (stride = blockDim.x/2; stride > 1; stride >> 1) {
```

```
-- syncthreads()
```

```
if (t < stride)
```

```
    psum[t] += psum[t+stride]
```

```
}
```

Adiferencia del anterior, en este los threads activos siempre están juntos lo que disminuye la divergencia en ejecución mejorando la eficiencia del programa



De esta manera la mitad continúan a la siguiente iteración como paraban antes pero esta vez todos juntos reduciendo activos de manera