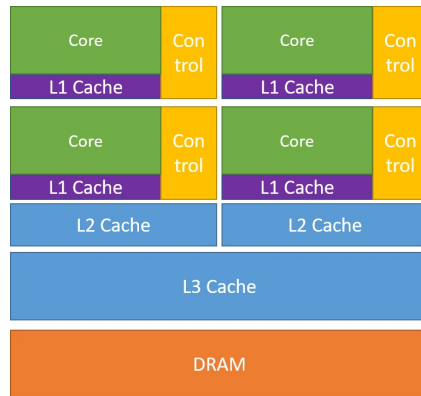


How does a GPU work?



CPU



OPTIMIZED FOR
REDUCING LATENCY

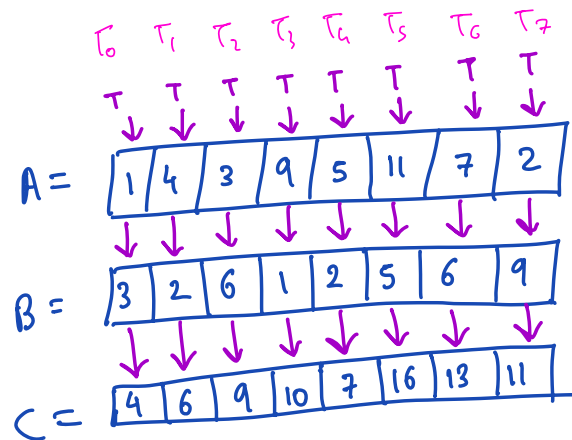


GPU



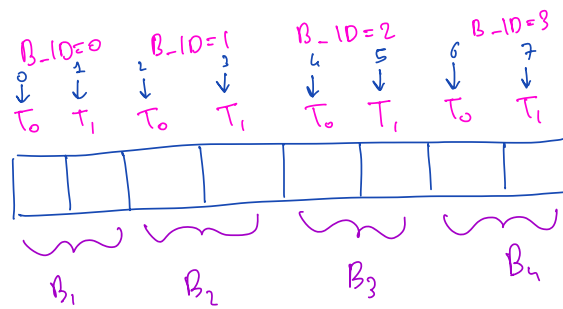
OPTIMIZED FOR
MAXIMIZING
THROUGHPUT

Vector addition



```
__global__ void cuda_vector_add_single_block(int *out, int *a, int *b, int n)
{
    int i = threadIdx.x;
    if (i < n)
    {
        out[i] = a[i] + b[i];
    }
}
```

```
// run the kernel without blocks
cuda_vector_add_single_block<<<1, N>>>>(d_out, d_a, d_b, N);
```



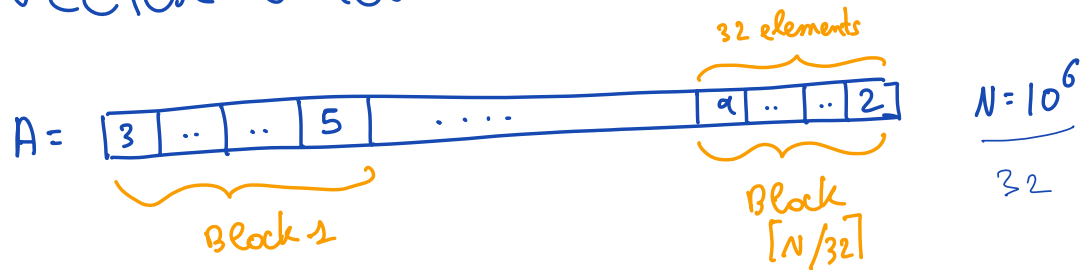
$$N = 8$$

4 cores

$$\text{Num-Blocks} = \frac{N}{2 \rightarrow \text{Block-Size}} = 4$$

$$\text{element-id} = i = \text{B-ID} \times \text{Block-Size} + \text{Tid}$$

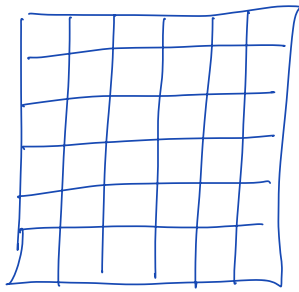
Vector addition with blocks



In this case we will have $\lfloor \frac{N}{32} \rfloor$ blocks
Block size = 32

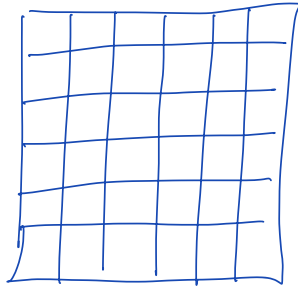
Matrix addition with blocks



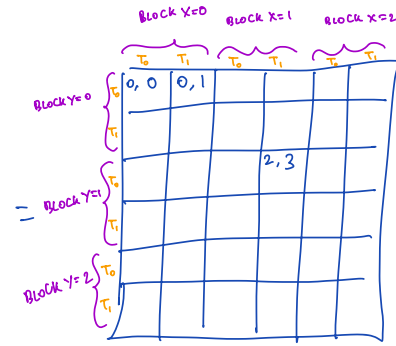


A

+

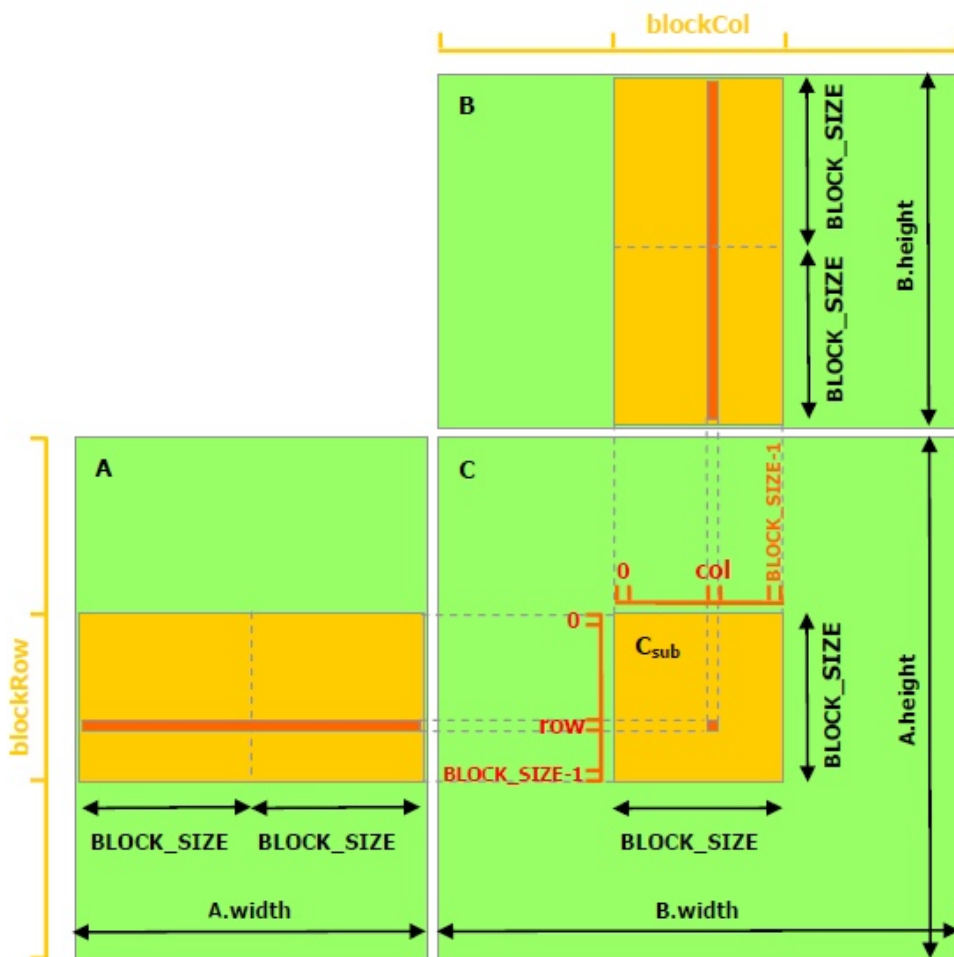


B



C

Shared memory



Many more topics

- Control divergence
- Occupancy
- Thread synchronization