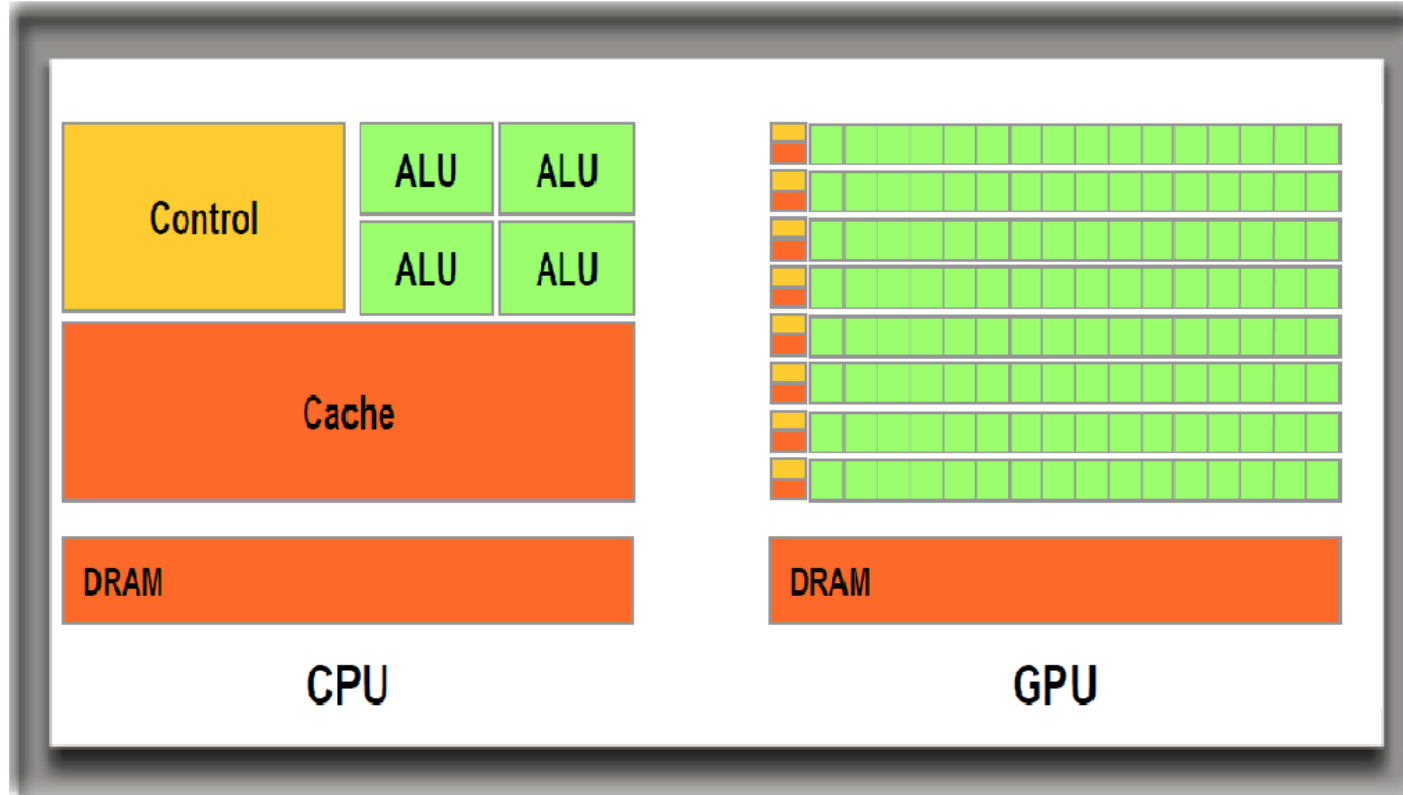


# GPUs vs CPUs, OpenACC vs OpenMP

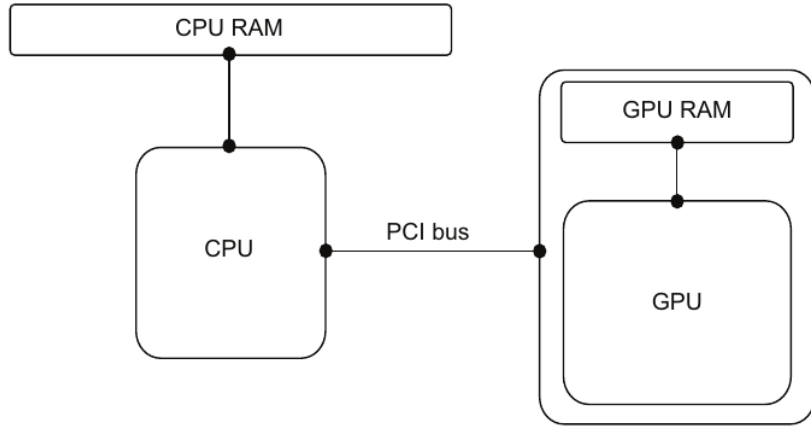
- What is a GPU?
- Its programming model
- OpenACC
- the shallow water model



CPUs are designed to handle a wide-range of tasks, but are limited in the concurrency of tasks that can be running.

GPUs can be an order of magnitude faster due to parallelism, but they are not as versatile as CPUs.

# A closer look, and some labels



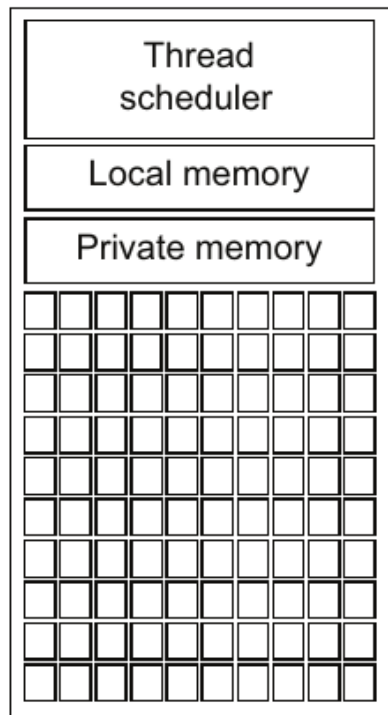
- CPU memory levels

- Registers
- Cache memory
- Main memory
- Network memory

- GPU memory levels

- Registers (GPU)
- Local memory (GPU)
- Global memory (GPU)
- Main memory (HOST)
- Network memory (REMOTE)

- The key to performance: Data handling and latency hiding

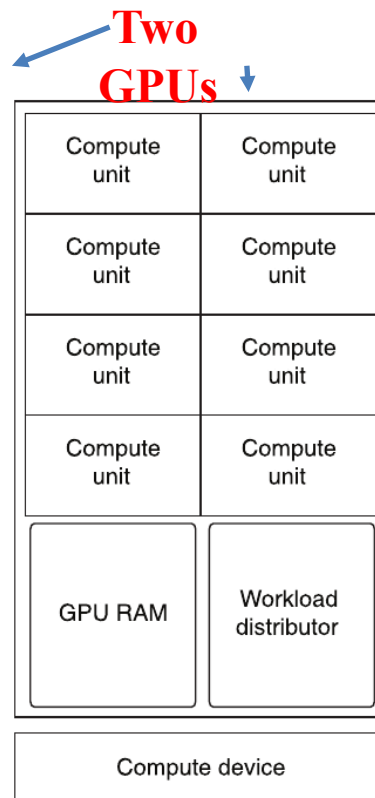
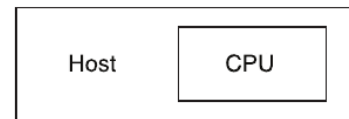
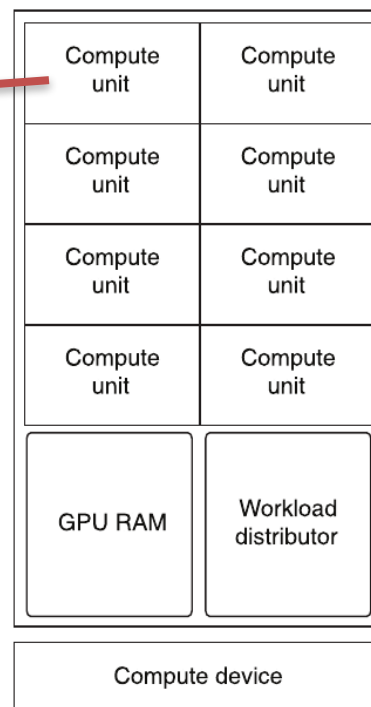


Simplified block diagram of a compute unit (CU) with a large number of processing elements (PEs).

**(Shared memory)**

**(Registers)**

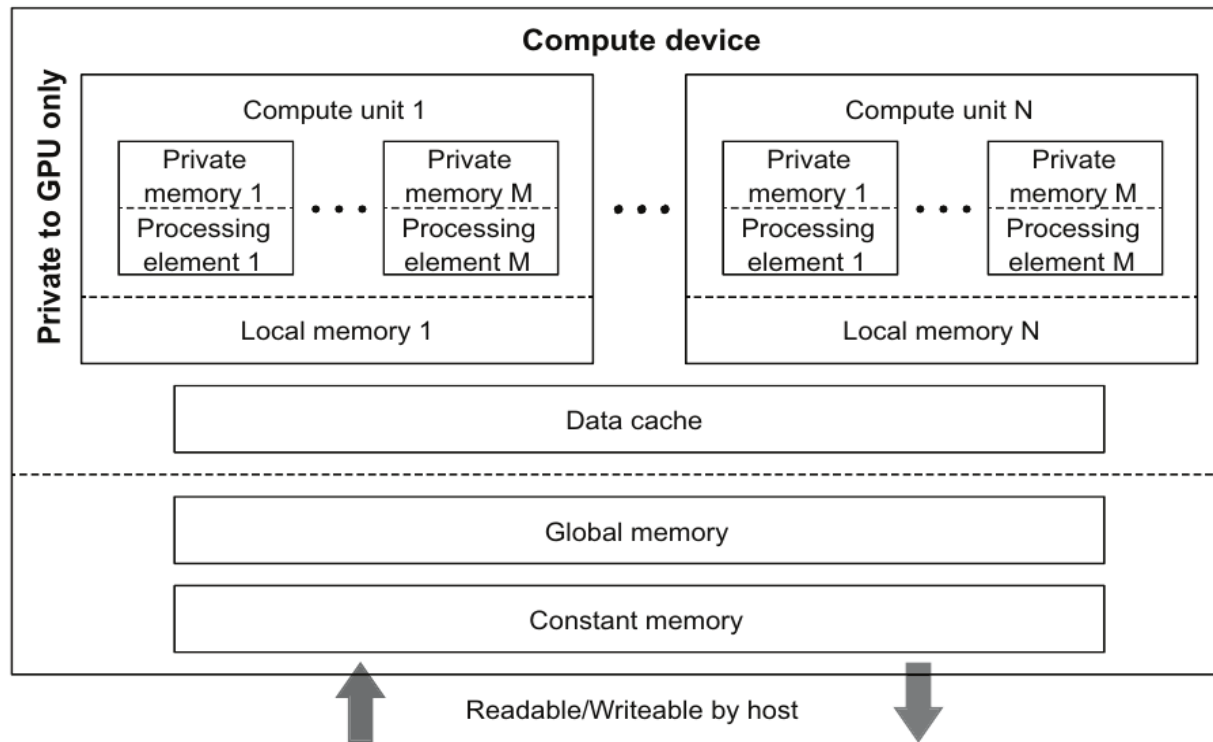
**Processing elements  
(Compute cores)**



A GPU compute device has multiple CUs. (CU, *compute unit*, is the term agreed to by the community for the OpenCL standard.) NVIDIA calls them *streaming multiprocessors* (SMs), and Intel refers to them as *subslices*.

$$\text{Peak Theoretical Flops (GFlops/s)} = \text{Clock rate MHz} \times \text{Compute Units} \times \text{Processing units} \times \text{Flops/cycle}$$

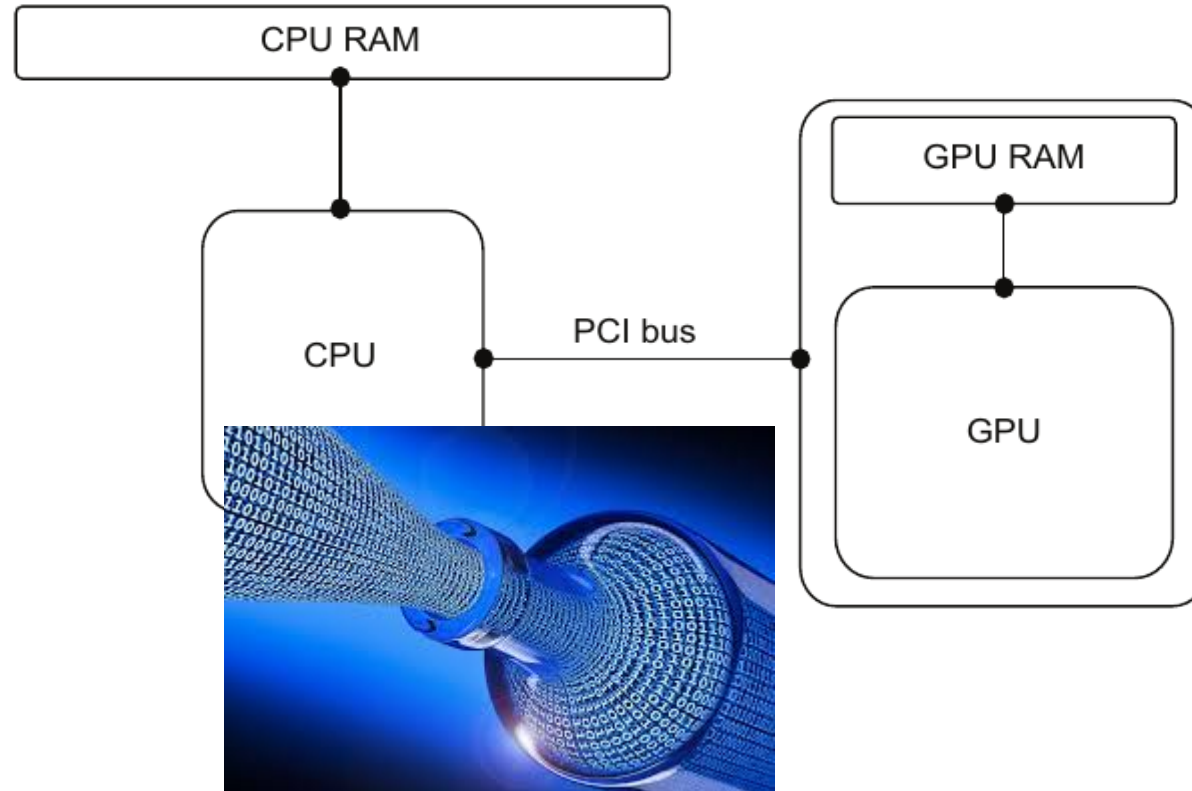
# The Memory



Specialized global memory with higher bandwidth than CPUs

$$\text{Theoretical Bandwidth} = \text{Memory Clock Rate (GHz)} \times \text{Memory bus (bits)} \times (1 \text{ byte}/8 \text{ bits})$$

## Lastly: the Peripheral Component Interconnect (PCI Bus)



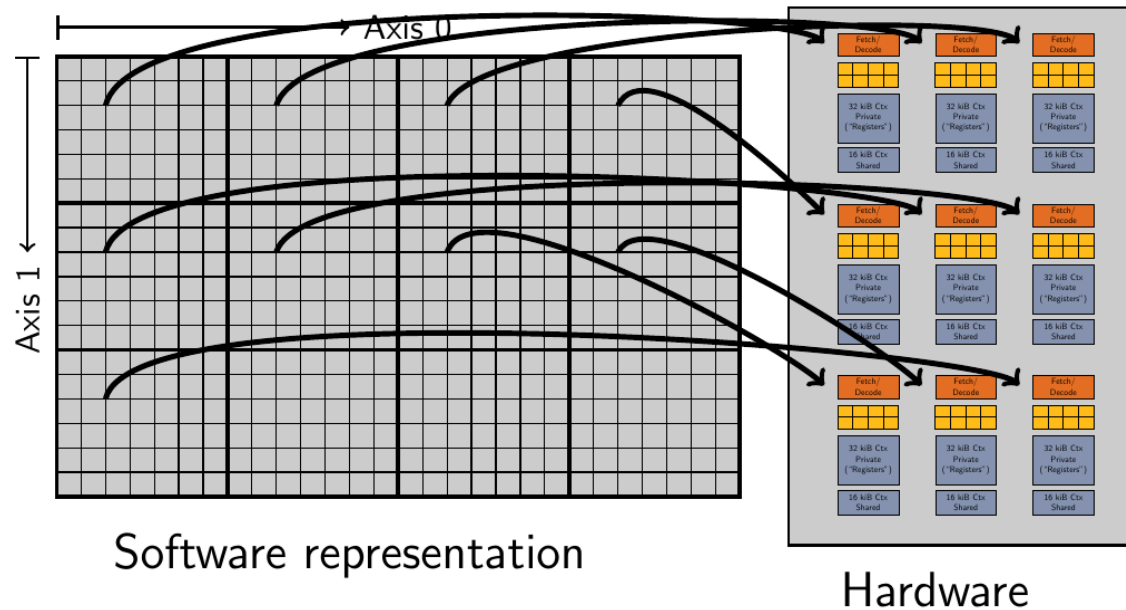
$$\text{Theoretical Bandwidth (GB/s)} = \text{Lanes} \times \text{TransferRate (GT/s)} \times \text{OverheadFactor (Gb/GT)} \times \text{byte/8 bits}$$

## What we'll use today:



2 Nvidia A30, sliced into 12 vGPUs each, so there are 24 vGPUs with 6GB and 14 CU each. 80k dkk in total.

... just for you, but please log-off if you don't use them ;-)

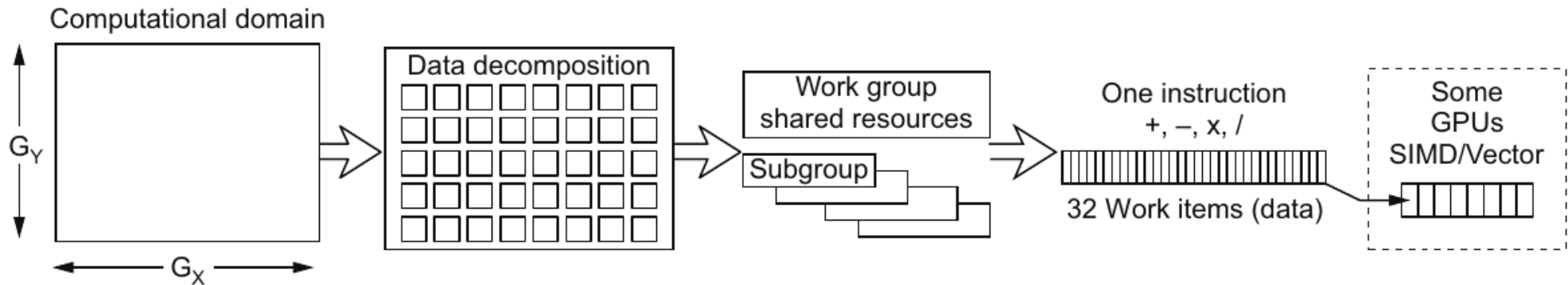


- Program as if an infinite amount of cores
- Data centric programming model
- Extremely parallel programming model
  - One thread per data element

# GPU programming model



## Common programming abstractions across various GPU models.



**Terminology can be confusing:**

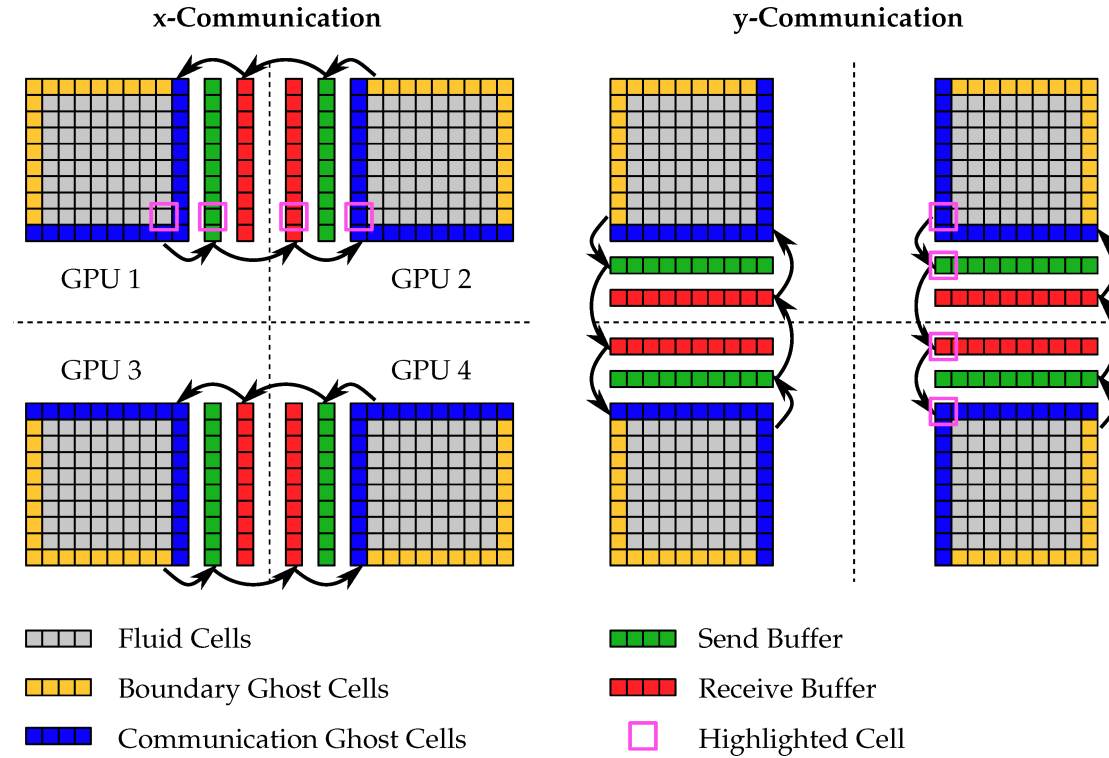
**Ndrange** = grid = loop bounds

**Work group** = thread block = loop block

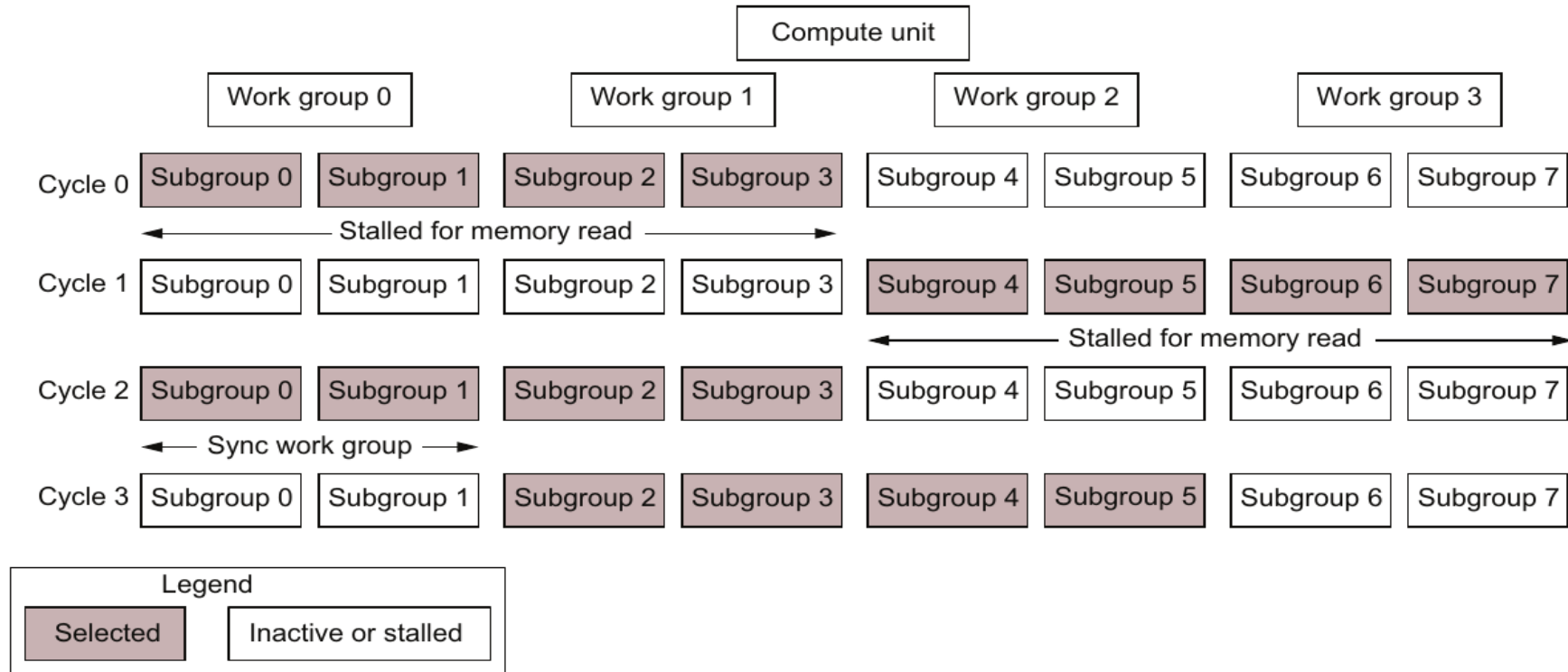
**Subgroup** = wavefront = SIMD length

**Work item** = thread

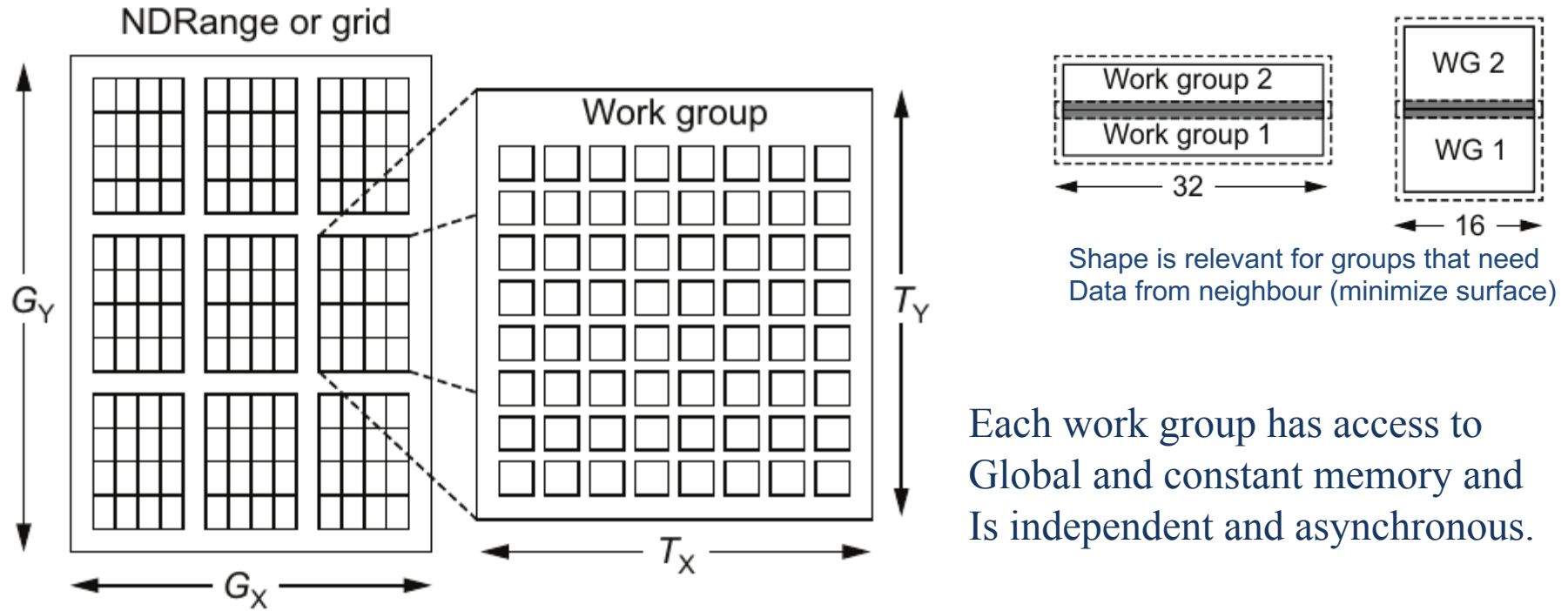
**Key: little coordination needed within operations**



Communication for stencil operations through ghost cells



**Latency hiding/ reduces the need for deep memory hierarchy**



Data decomposition of a  $1024 \times 1024$  2D computational domain

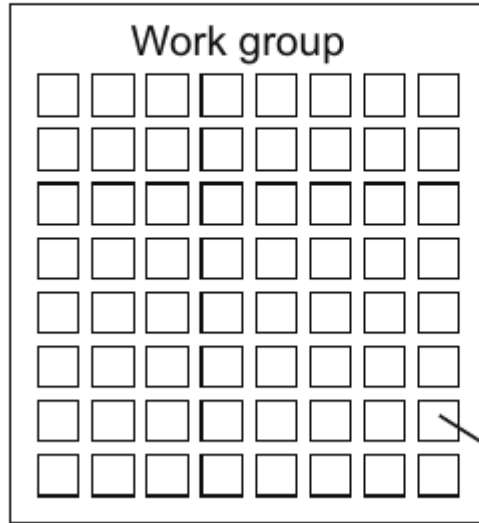
If we set a tile size of  $16 \times 8$ , the data decomposition will be as follows:

$$NT_x = G_x / T_x = 1024 / 16 = 64$$

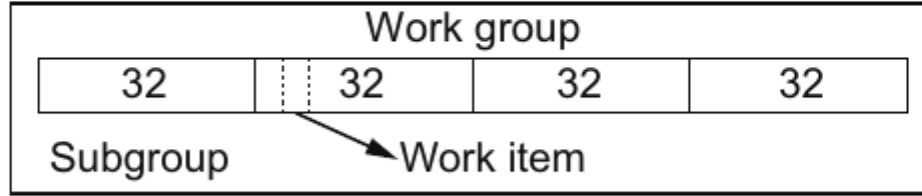
$$NT_y = G_y / T_y = 1024 / 8 = 128$$



# Workgroups



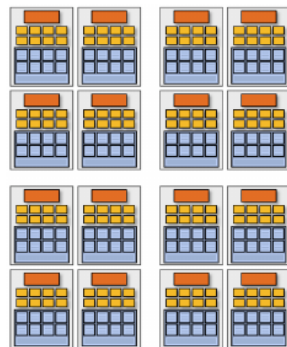
**Each subgroup (warp) is 32 work items (threads).  
This work group size is 128, but can be up to 1,024.**



The work groups spreads out the work across the threads on a compute unit and have access to local, shared memory.

... and finally, each work item may be able to do a vector or SIMD

- Do we care about the number of cores?
- Do we care about the number ALU's per core?



- NO !!!  
Program as if an infinitely amount of cores  
Program as if an infinitely amount of ALU's per core

# Practicalities

```
#pragma acc <directive> [clause]
#pragma omp <directive> [clause]
```

Many similarities between OpenACC and OpenMP, there attempts to merge both.

OpenMP (1997) long focused on CPUs.  
OpenACC (2011) specifically designed to Handle GPUs.



It appears to me that in OpenACC begin have less opportunities to make mistake OpenMP experts have more control.

Steps	CPU	GPU
Original	a[1000] for or do loop for or do loop free a	
1. Offload work to GPU	a[1000] <b>#work pragma</b>  <b>#work pragma</b> free a	for or do loop (128 threads)  for or do loop (128 threads)
2. Manage data movement to GPU	<b>#data pragma{</b> <b>#work pragma</b> <b>#work pragma</b> <b>}</b>	a[1000] for or do loop (128 threads) for or do loop (128 threads) free a
3. Optimize GPU kernels	<b>#data pragma{</b> <b>#work pragma</b> <b>#work pragma</b> <b>}</b>	a[1000] for or do loop ( <b>256 threads</b> ) for or do loop ( <b>64 threads</b> ) free a

# Benefits of using OpenACC

- achieve parallelization on GPUs without having to learn an accelerator language such as CUDA
- similar philosophy to OpenMP and very easy to learn

```
#pragma acc kernels
{
    for (int i = 0; i < N; i++) {
        y[i] = a * x[i] + y[i];
    }
}
```

```
#pragma acc parallel loop
{
    for (int i = 0; i < N; i++) {
        y[i] = a * x[i] + y[i];
    }
}
```



## Examples of some useful directives:

```
#pragma acc parallel [clause [[,] clause]...] new-line
```

```
#pragma acc parallel loop present(data)
```

```
#pragma acc parallel loop collapse(2) present(a, b)
```

### **To move data around (sparingly ;-)**

```
#pragma acc data copy(myworld)
```

```
#pragma acc update self(myworld)
```

A full list of directives can be had from <https://www.openacc.org/resources>,  
The reference card under Guides.

