

# Introduction aux GPUs, à la programmation CUDA

12 juin 2024

Jussieu, Paris

Samuel Thibault

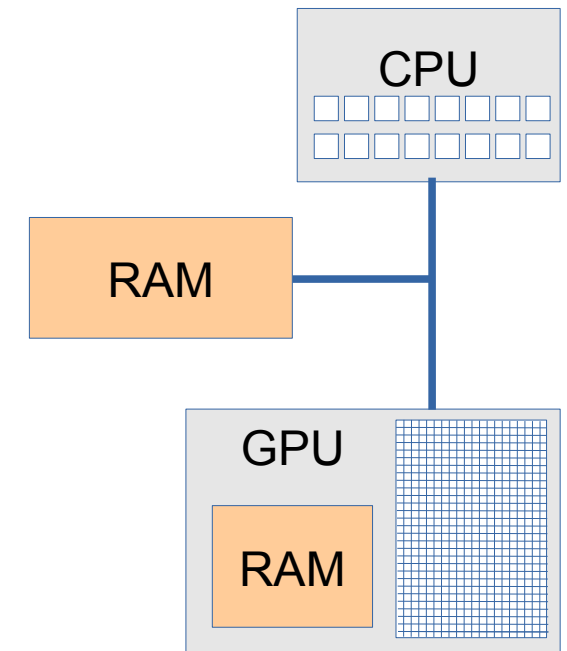


*Inria*

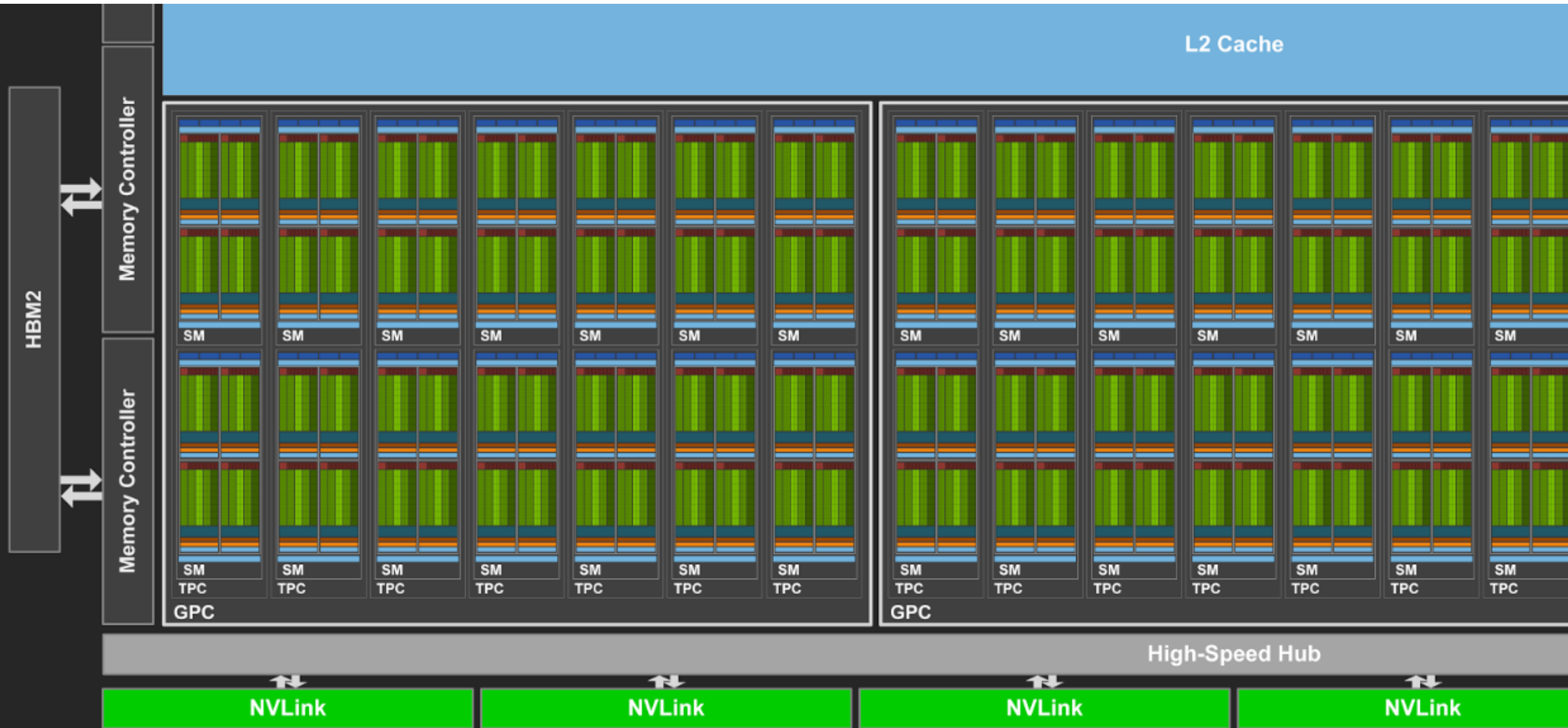
anr<sup>®</sup>

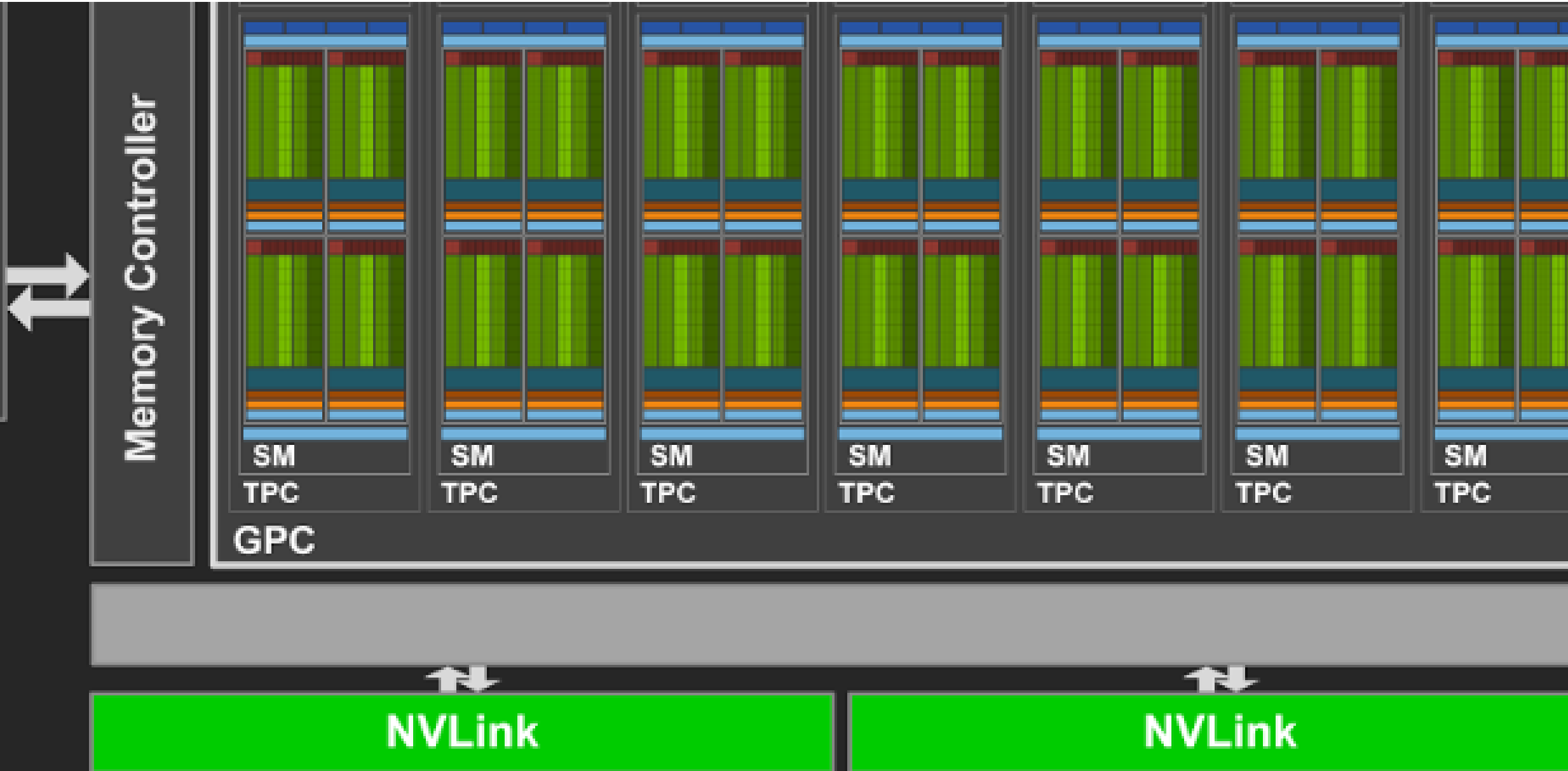


- Programmation particulière
  - Parallélisme massif !
  - Mémoire séparée (la plupart du temps)
- Développement particulier
  - Compilation différente
  - Lancement par offload logiciel
    - On ne peut pas juste recompiler
  - Outils de débogage différents











Massive parallelism !

Example: V100

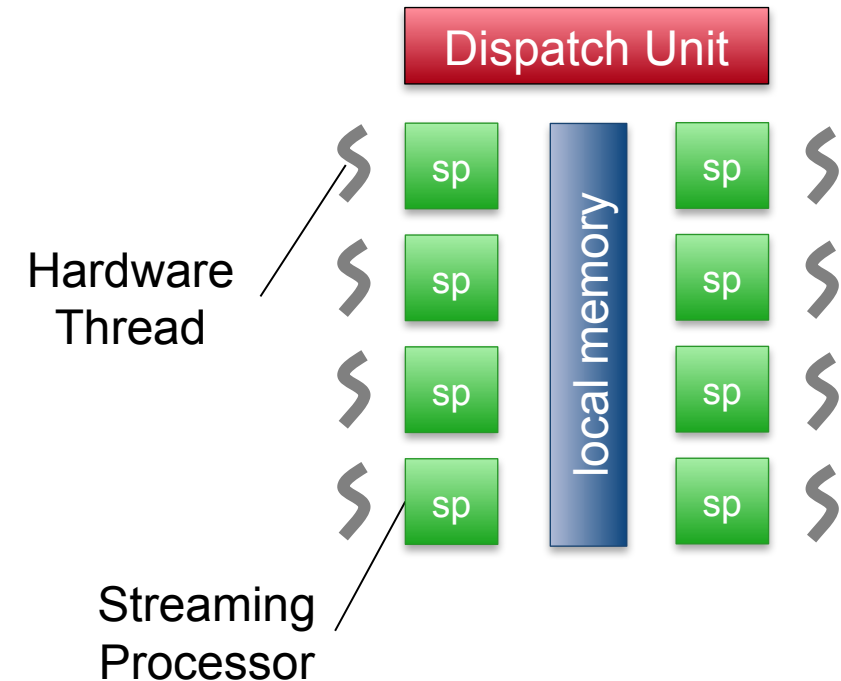
- 5,120 CUDA cores, structured in SM, TPC, GPC, ...
  - Parallelism to be structured
  - And massive to overlap memory access
- SP: ~16TFlop/s
- DP: ~8TFlop/s
- HBM Memory: ~900GB/s
- (and not counting the tensor cores)

# GPU execution model

Basic block: Streaming Multiprocessor (SM)

= Cluster of streaming processors (~core)

- Local memory sharing
- Synchronization
- 64KB of registers
- hardware threads
  - Creation/destruction is free !

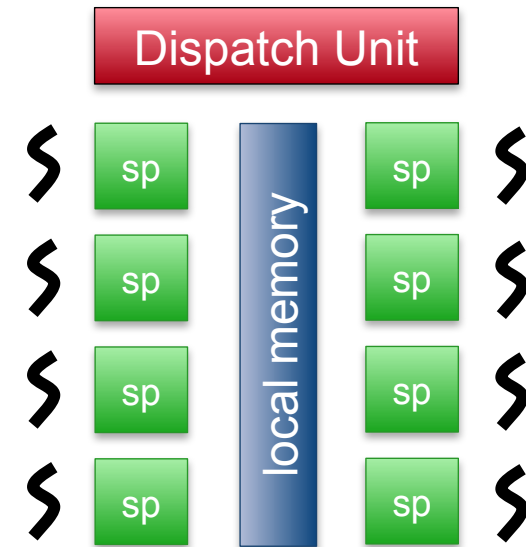




# GPU execution model

Only one instruction dispatch unit per SM

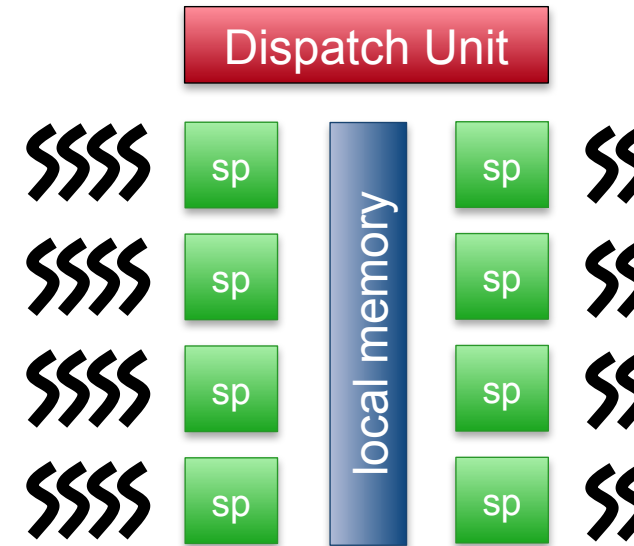
- All cores execute the same instruction at the same clock cycle
- On different data  $\Rightarrow$  SIMD
- Takes several cycles to fetch&decode



# GPU execution model

Only one instruction dispatch unit per SM

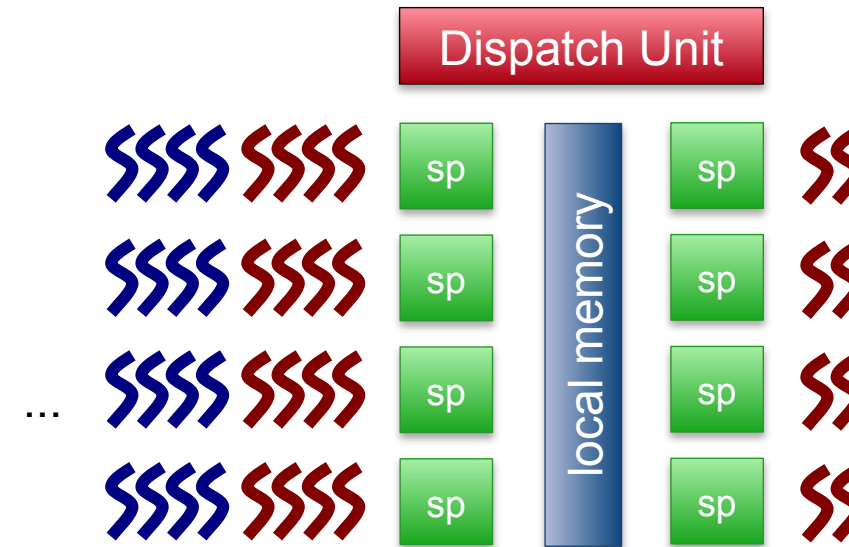
- All cores execute the same instruction at the same clock cycle
  - On different data  $\Rightarrow$  SIMD
  - Takes several cycles to fetch&decode
- $\Rightarrow$  overlap by running yet more threads running the same instruction
- Warp = 32 threads



# GPU execution model

Only one instruction dispatch unit per SM

- Loading data from global memory takes cycles
- ⇒ overlap by running yet more threads
  - 128 are enough to hide memory latency



# GPU execution model

Threads, threads, threads...

- But they are trivial to create
- Example: vector scaling

```
GPU function  global__ void
              _cuda(unsigned n, float *val, float factor)
{
    val[threadIdx.x] *= factor;
}
```

GPU pointer !

```
CPU function  void
              func(unsigned n, float *val, float factor)
{
    vector_mult_cuda<<<1,n,0>>>(n, val, factor);
}
```



# GPU execution model

But hierarchy of SM

- Thread scheduling needs structure : blocks

```

static __global__ void vector_mult_cuda(unsigned n, float *val, float factor)
{
    unsigned i = blockIdx.x*blockDim.x + threadIdx.x;

    if (i < n) val[i] *= factor;
}

extern "C" void Compensate round-up signed n, float *val, float factor)
{
    unsigned threads_per_block = 64;
    unsigned nblocks = (n + threads_per_block - 1) / threads_per_block;

    vector_mult_cuda<<<nblocks, threads_per_block, 0>>>(n, val, *factor);
}

```

Compensate round-up

Two-dimension indexing (can be three)

# GPU execution model

Threads, threads, threads...

- Trivial to create
  - Massive parallelism
  - Forget about for loops
- Thread scheduling needs structure : blocks
  - How large? Depends on
    - the GPU generation,
    - available parallelism, ...

## Thread divergence

```
static __global__ void
vector_mult_cuda(unsigned n, float *val, float factor)
{
    unsigned i = blockIdx.x*blockDim.x + threadIdx.x;

    if (i < n && i % 2 == 0) val[i] *= factor;
}
```

Only half the threads actually work

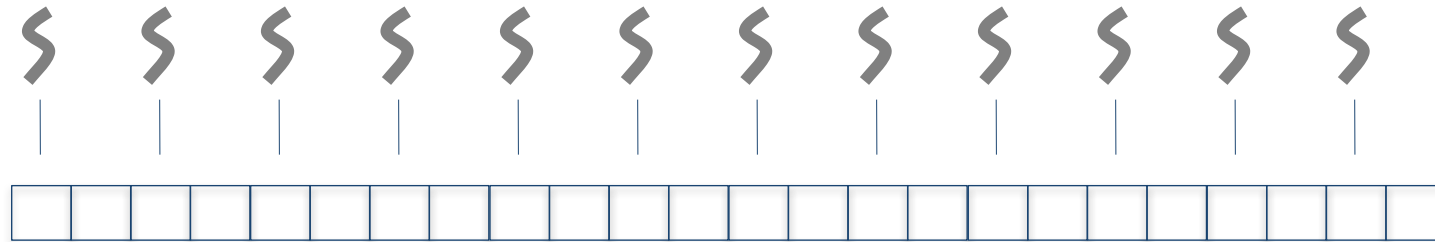
⇒ 2x parallelism loss

Related to warps ; depends on GPU generation

## Memory coalescing

```
static __global__ void vector_mult_cuda(unsigned n, float *val, float factor)
{
    unsigned i = blockIdx.x*blockDim.x + threadIdx.x;

    val[i*2  ] *= factor;
    val[i*2+1] *= factor;
}
```



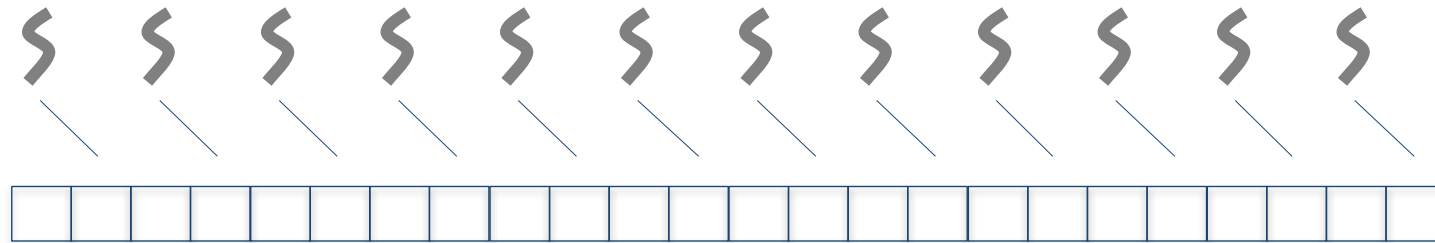
Performance is weak : accesses are non-coalesced



## Memory coalescing

```
static __global__ void vector_mult_cuda(unsigned n, float *val, float factor)
{
    unsigned i = blockIdx.x*blockDim.x + threadIdx.x;

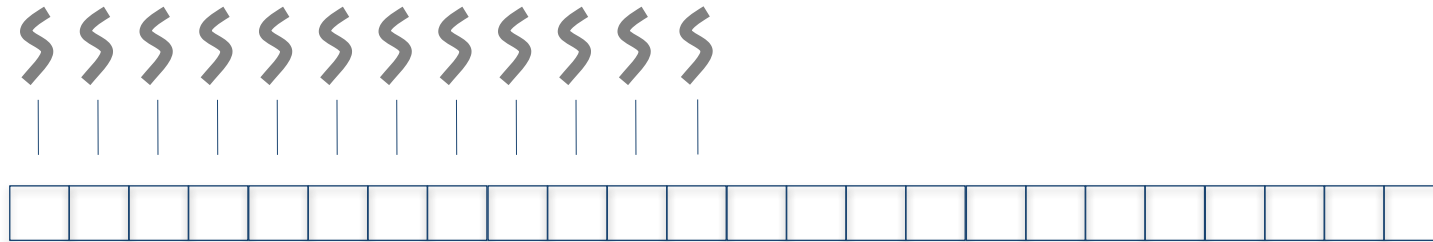
    val[i*2 ] *= factor;
    val[i*2+1] *= factor;
}
```



Performance is weak : accesses are non-coalesced

## Memory coalescing

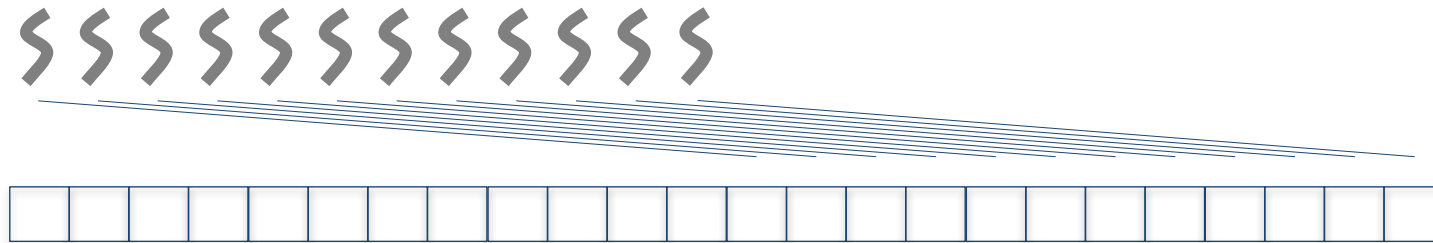
```
static __global__ void vector_mult_cuda(unsigned n, float *val, float factor)
{
    unsigned i = blockIdx.x*blockDim.x + threadIdx.x;
    unsigned nthreads = blockDim.x * blockDim.y;
    val[          i] *= factor;
    val[nthreads + i] *= factor;
}
```



Accesses are coalesced

## Memory coalescing

```
static __global__ void vector_mult_cuda(unsigned n, float *val, float factor)
{
    unsigned i = blockIdx.x*blockDim.x + threadIdx.x;
    unsigned nthreads = blockDim.x * blockDim.y;
    val[          i] *= factor;
    val[nthreads + i] *= factor;
}
```



Accesses are coalesced

## Memory coalescing

- Depends on the GPU generation
- Better use structures of arrays
  - E.g. all threads get the temperature, then the pressure, etc.
- For the global memory
- Not for the local shared memory
  - Has to be allocated explicitly

Summary : requires

- Massive parallelism
- Taking care of divergence
- Taking care of memory coalescing
  - Or taking care of using the local shared memory
- And other details I don't have time to explain :)

If you can avoid writing kernels it's best...

But we need to get data onto the GPU...

# Data management

- Allocate memory on the GPU
- Copy between GPU and CPU
  - For copy to be efficient, use `cudaMallocHost` for CPU memory
    - Pinned memory for DMA

```
void scal_func(unsigned n, float *val_cpu, float factor)
{
    float *val_gpu;
    cudaMalloc(&val_gpu, n * sizeof(float));
    cudaMemcpy(val_gpu, val_cpu, n * sizeof(float), cudaMemcpyHostToDevice);
    scal_cuda_func(n, val_gpu, factor);
    cudaMemcpy(val_cpu, val_gpu, n * sizeof(float), cudaMemcpyDeviceToHost);
}
```

# Data management

But transfer takes time

- Overlap cudaMemcpy with another kernel execution
- Async variants of cuda operations, and poll or wait
- Use Streams for coordination

```
void scal_func(unsigned n, float *val_cpu, float factor)
{
    float *val_gpu;
    cudaMalloc(&val_gpu, n * sizeof(float));
    cudaMemcpy(val_gpu, val_cpu, n * sizeof(float), cudaMemcpyHostToDevice);
    scal_cuda_func(n, val_gpu, factor);
    cudaMemcpy(val_cpu, val_gpu, n * sizeof(float), cudaMemcpyDeviceToHost);
}
```



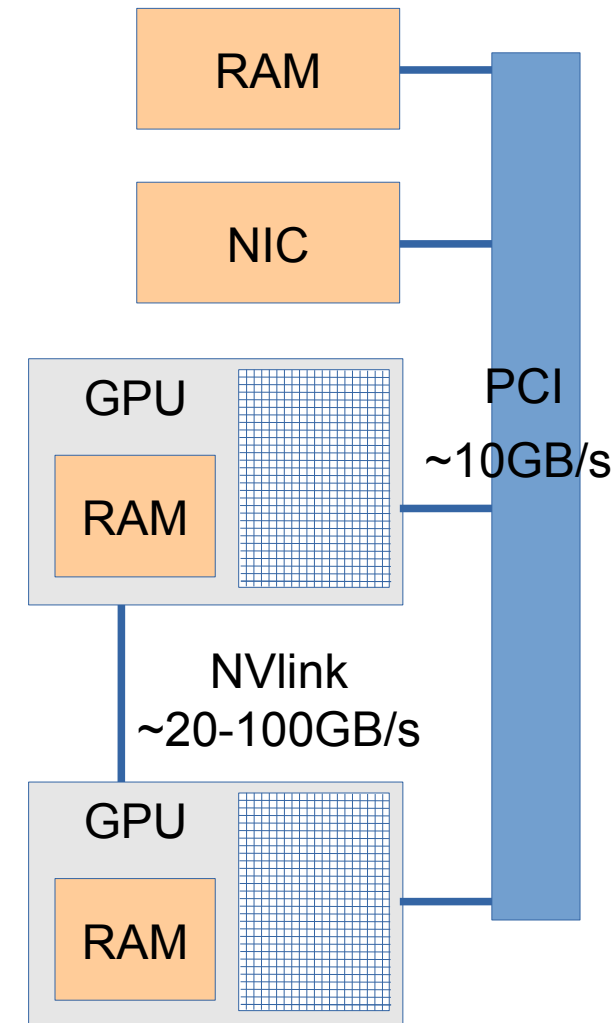
Or use cudaGraphs

- Express dependencies explicitly between cudaMemcpyAsync and kernel calls
- Let the GPU process the task graph

# Data management

## Topology matters

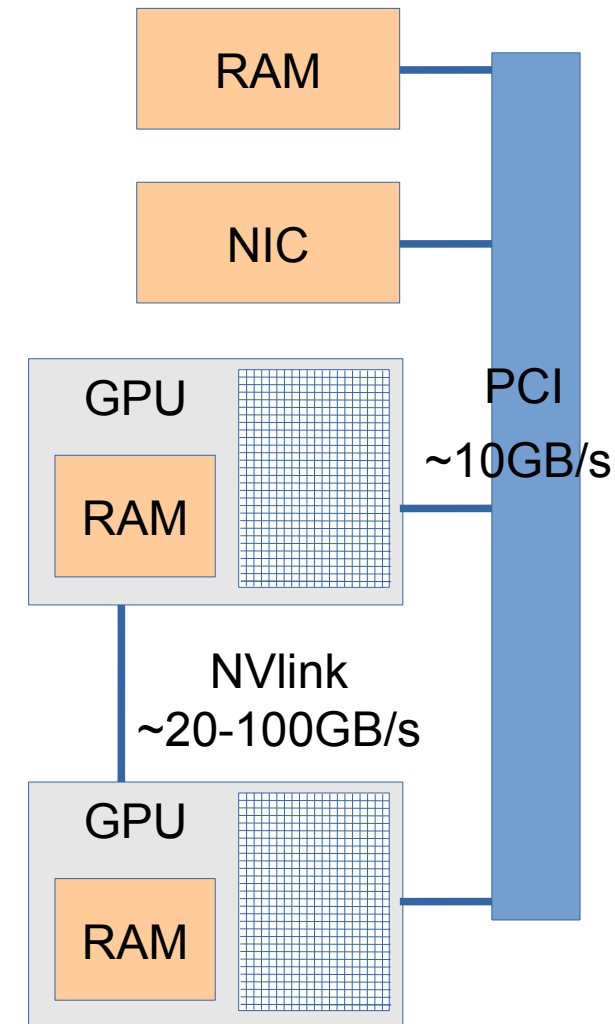
- NVlink provides fast GPU-GPU transfers
- GPU-Direct provides fast NIC-GPU transfers
  - With CUDA-aware MPI implementation
- Beware of PCI topology
- Beware of NUMA topology



# Data management

Or use unified memory

- Pass managed pointers to kernels
- CUDA runtime handles transfers on the fly
- But CUDA runtime is mostly blind
  - On-demand migration, introduces latency
  - With `cudaMemAdvise` / `cudaMemPrefetch` hints, can be efficient



## Summary

- HBM memory embedded in GPU
- Requires transfers
- To be pipelined with computation kernels

If you can avoid managing this it's best...

GPUs are beasts

- Require massive parallelism
  - With massive data accesses
- Optimizing is very involved
- Better delegate it to another software layer