

CUDA

A100 tensorcore

<https://zhuanlan.zhihu.com/p/620257581>

GPGPU v2.01.pdf

<https://zhuanlan.zhihu.com/p/166180054>

https://www.tinyedi.com/cuda_learning/#pipeline

<https://docs.nvidia.com/cuda/cuda-c-programming-guide/index.html>

<https://zhuanlan.zhihu.com/p/486224812>

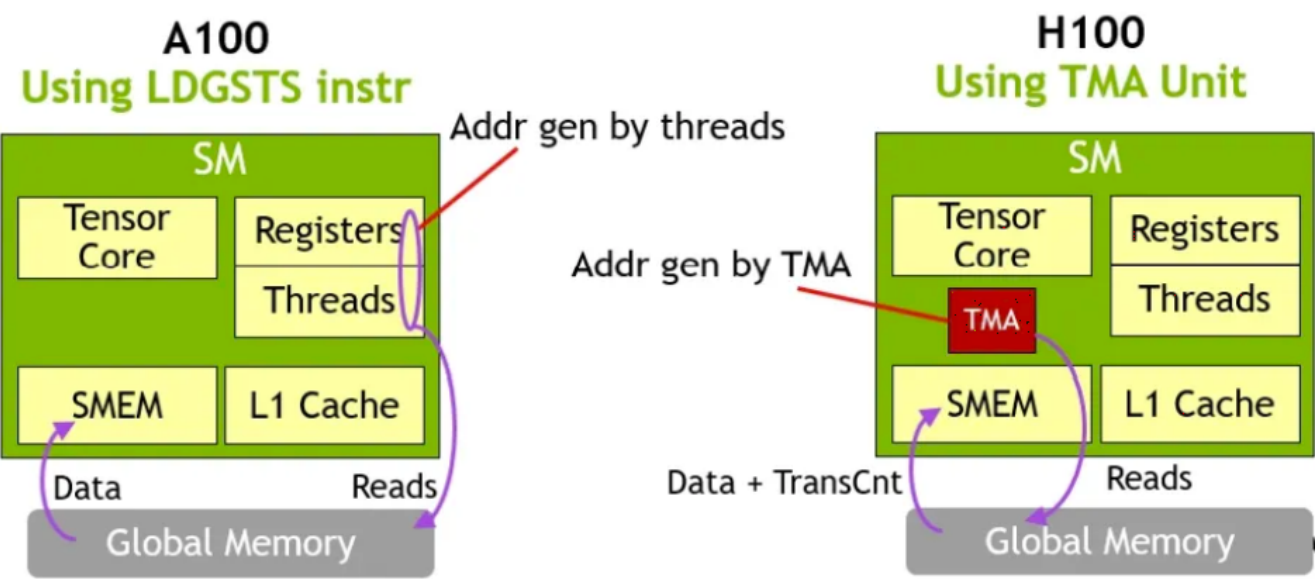
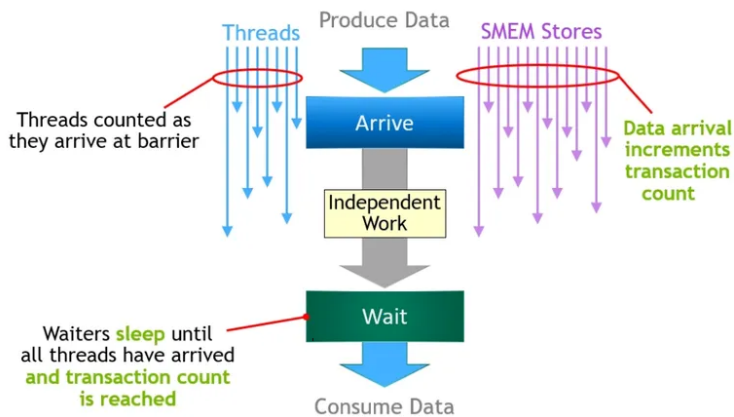


Figure 19. Asynchronous Memory Copy with TMA on H100 vs LDGSTS on A100

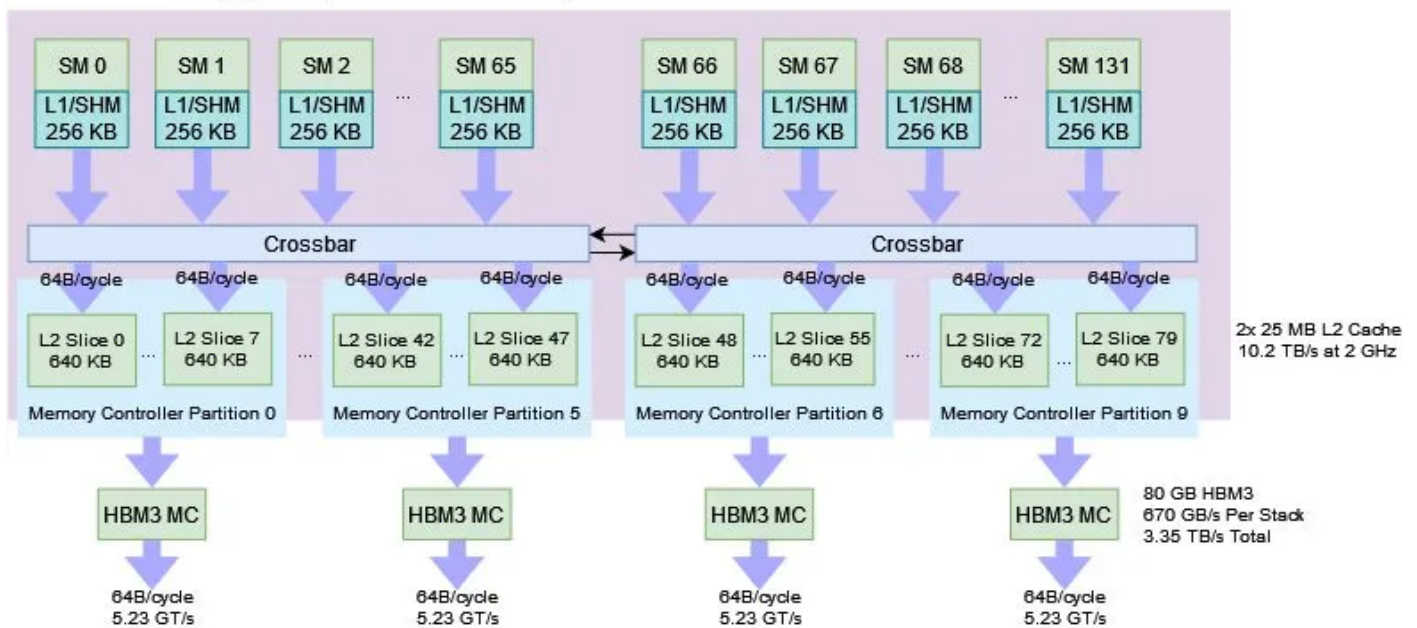
Async Transaction Barrier (New on H100)



Control code reuse read barrier write barrier wait barrier yield hint stall count
reuse GPR bank conflict
GPR barrier thread
yield stall count warp
memory memory bound
memory bound latency

stall□□□□ cycle□□□□□□□□□□

Nvidia Hopper (H100 SXM5)



CUDA

- `cudaLaunchKernel`
- `cudaLaunchKernelEx`
- `cudaLaunchKernelEx` maximize device utilization.
- `cudaLaunchKernelEx`
- `cudaLaunchKernelEx` whenever possible.
- `cudaLaunchKernelEx` threads within the same warp.

Adjust kernel launch configuration to

Minimize redundant accesses to global memory

Avoid long sequences of diverged execution by

TMA

TMA `cudaMemcpyAsync`

SM `cudaMemcpyAsync`

1D

`cudaMemcpyAsync` 5D `cudaMemcpyAsync`

`cudaMemcpyAsync`

add/min/max `cudaMemcpyAsync`

/

`cudaMemcpyAsync`

memcpy_async API `cudaMemcpyAsync`, `cudaMemcpyAsync` group, barrier, pipeline.

- group: `cudaMemcpyAsync` `cooperative_groups::wait(group)`, `cudaMemcpyAsync` `group.sync()`, `cudaMemcpyAsync` `group.sync()` `cooperative_groups::wait`, `cudaMemcpyAsync` `group.sync()`.
- barrier: `cudaMemcpyAsync` `barrier.arrive_and_wait()`, `cudaMemcpyAsync` `barrier` count, `cudaMemcpyAsync` count `cudaMemcpyAsync` `async_memcpy`.
- pipeline: `cudaMemcpyAsync` queue

