**⑆ ChatGPT**

# Course Plan: Accelerated Parallel Computing with NVIDIA GPUs

**Overview:** This course explores how modern GPUs – especially NVIDIA's – enabled a revolution in parallel computing. GPUs like the 1999 GeForce 256 (Nvidia's first GPU) "outperformed existing products by a wide margin" [1] , and today GPUs dominate AI and HPC. For example, a GPU-based deep learning training ran ~3× faster than on CPU [2] , and a GPU array library (CuPy) finished a large vector math task 500× faster than NumPy on CPU [3] . We'll build on such examples (often using free Colab GPU runtimes) to teach fundamentals (parallel hardware, programming models, memory, sync, interconnect) with engaging stories (NVIDIA's history and innovations) and practical demos. Key libraries include NVIDIA CUDA, cuBLAS/cuFFT, Thrust, as well as Python tools like CuPy (GPU NumPy) [4] , Numba, PyCUDA and ML frameworks (PyTorch/ TensorFlow) for GPU acceleration. Real-world impacts abound: medical image processing that once took minutes on CPU now takes seconds with GPUs [5] , and large protein analyses run hundreds of times faster on GPUs [6] , powering advances in healthcare, defense, science and beyond.

## Week 1: Introduction to Parallel Architectures

*Story:* The origin of GPUs – three engineers at a Denny's in 1993 – and the GeForce 256 (1999) as the first true GPU [1] . GPUs turned parallel processing from niche into a trillion-dollar industry. For example, training a CNN on GPU vs CPU showed ~18 min → ~6 min speed-up [2] , and GPU array code (CuPy) gave ~532× speed-up on large vector math [3] . We introduce parallelism fundamentals and how NVIDIA's vision (compute-focused GPUs) launched modern AI.

- **Subtopics:** Flynn's taxonomy (SISD, SIMD, MISD, MIMD); comparing CPU vs GPU hardware (few fast cores vs thousands of parallel ALUs); history of GPU architectures (GeForce256 → Fermi → Pascal → Ampere → Hopper → Blackwell).
- **Key Concepts:** Amdahl's and Gustafson's laws (limits of scalability); parallel speedup/efficiency; SIMD (CPU) vs SIMT (GPU threads in warps).
- **Hands-on:** Launch a basic CUDA (or CuPy/Numba) vector-add kernel on Colab to measure parallel execution time versus serial CPU.
- **NVIDIA Focus:** CUDA (released 2006) for general-purpose GPU computing [7] ; milestone facts (e.g. Nvidia's >80% share of AI GPU market). Story note: GPUs processed GPUs' first "Big Bang" in AI ~2009 when deep nets were trained 10× faster [7] .

# Week 2: Parallel Programming Models and GPU Execution

*Story:* The "CUDA revolution" began in 2006 – democratizing GPU supercomputing on a desktop. Early adopters saw 100× speedups for physics/drug simulations by moving code to GPUs. We'll cover CPU vs GPU programming models and how GPUs execute kernels.

- **Subtopics:** Parallel programming models (shared memory, message passing/MPI, OpenMP, CUDA programming); GPU execution hierarchy (CUDA: grids of thread blocks, each block of threads, warps); Streaming Multiprocessors (SMs), warps, and SIMT scheduling.
- **Key Concepts:** Thread scheduling/occupancy (warp divergence cost); memory access patterns (coalescing for performance); CUDA streams and concurrency; debugging and error handling (race conditions, etc.).
- **Libraries/Tools:** CUDA Toolkit, Thrust (C++ STL for GPU), NVIDIA HPC SDK; Python tools like PyCUDA and Numba; CuPy for Python GPU arrays [4]. Example: CuPy uses CUDA under the hood so NumPy code runs on GPU (order-of-magnitude speed-ups) [4] [3].
- **Hands-on:** Write a CUDA kernel or use Numba to add vectors, or simply convert a NumPy loop to CuPy and compare runtimes. Demonstrate dynamic parallelism (Kepler feature) by spawning sub-kernels (e.g. fractal generator).
- **NVIDIA Focus:** Scalable CUDA model (one code runs on 1 vs many GPUs). Fun fact: Kepler (2012) introduced **dynamic parallelism** (kernels launching kernels) – think of a GPU raytracer spawning more tasks on-the-fly. Also note libraries like cuBLAS, cuFFT, cuRAND that accelerate common math kernels on GPU.

# Week 3: Memory Hierarchy – Registers, Caches and Global Memory

*Story:* Real-time graphics pushed NVIDIA to innovate memory (rendering huge game worlds). Modern GPUs pack vast bandwidth – for example, the NVIDIA Hopper H100 GPU's HBM3 memory delivers ~3 TB/s bandwidth [8] (orders of magnitude above typical CPU). We'll explore GPU memory levels and optimizing data access.

- **Subtopics:** GPU memory types – registers (per-thread), shared memory/L1 (per-SM), L2 caches, and global DRAM (HBM). Virtual memory on GPU (paging, Translation Lookaside Buffers, UVA unified address space).
- **Key Concepts:** Latency vs bandwidth; memory coalescing, caching strategies, prefetching; cache hit/miss rates and replacement (GPUs often use simple write-through L1, e.g. "GPU-VI" 2-state protocol 【37†】 to reduce miss stalls). CUDA Unified Memory (UVA) allows the system to migrate data.
- **Hands-on:** Benchmark memory bandwidth on Colab: e.g., copy vs compute overhead. Use NVIDIA Nsight Compute (or simple CUDA timers) to profile global vs shared memory access.
- **NVIDIA Focus:** Hopper's HBM3 (3 TB/s) [8] and large shared/L2 caches (e.g. 60 MB L2 in Hopper). GPUs invert cache ratios: sometimes L1 is small relative to registers to serve many threads. Introduce tools: NVIDIA Nsight and the CUDA profiler to visualize memory bottlenecks.

## Week 4: Cache Coherence Basics (Snooping and Directories)

*Story:* Ensuring consistency across many caches was a key challenge in early multi-GPU systems (e.g. SLI for gaming). NVIDIA's research led to relaxed coherence for GPUs to reduce overhead.

- **Subtopics:** Cache coherence problem and states (Valid/Invalid basics). Snooping protocols (bus-based broadcasting), vs directory-based protocols. Overview of the VI (Valid-Invalid) 2-state protocol.
- **Key Concepts:** Broadcast vs point-to-point coherence traffic; invalidation overhead; coherence analysis. (E.g., GPUs often use simpler schemes like write-through L1 and rely on software coherence.)
- **Hands-on:** Simulate a simple coherence scenario: e.g., two CUDA threads on different SMs reading/writing the same data, and observe consistency with `__threadfence()` or atomics.
- **NVIDIA Focus:** Volta's NVLink allowed hardware-supported multi-GPU coherence (e.g., Nvidia's NVLink-C2C provides coherent shared memory at ~900 GB/s [9] ). Note: Relaxed coherence in GPU hardware can cut overhead ~30–50% compared to full MESI.

## Week 5: Cache Coherence Protocols – MSI, MESI, Dragon

*Story:* Classic protocols (MSI/MESI) were invented in the 1980s, but GPUs adapted them. The Dragon protocol (update-based) resembles how GPUs minimize invalidations during gaming workloads.

- **Subtopics:** MSI/MESI states and transitions (Modified, Shared, Exclusive, Invalid). Dragon protocol (write-update variant). Comparison of **update vs invalidate** strategies and bandwidth trade-offs.
- **Key Concepts:** Cache miss types (compulsory, capacity, conflict, coherence miss). Correctness: ensuring no stale data (race conditions, memory fences). Protocol verification basics (e.g. serialization).
- **Hands-on:** Use the CUDA profiler (nvprof/nvidia-smi) to analyze cache hit rates and L2 traffic on a simple CUDA program. Try toggling unified memory or atomic operations to see coherence effects.
- **NVIDIA Focus:** Tesla V100/A100 GPUs use a 2-state GPU-VI variant (write-through L1, making coherence simpler). Interesting: GPUs often favor "update" for some data (like textures) to reduce frame glitches, akin to Dragon's idea. (No direct cite; conceptual link.)

## Week 6: Bus-Based Multiprocessors – Atomic Bus and Snoop Design

*Story:* Early multi-GPU work (e.g. SLI) required atomic buses to avoid tearing between screens. Today, atomic memory ops are crucial for parallel algorithms (like reductions).

- **Subtopics:** Bus arbitration and tagging for snooping; single vs multi-level cache hierarchies. Atomic bus operations (read-modify-write on bus). Differences between simple bus vs split-transaction bus.
- **Key Concepts:** Bus contention and latency; hierarchical snoop (multi-level caches with coherence); atomic primitives (locks vs lock-free, compare-and-swap). CUDA supports atomic intrinsics (`atomicAdd`, etc.) across threads and even GPUs.
- **Hands-on:** Demo CUDA atomic operations: e.g., parallel reduction using `atomicAdd` vs per-block partial sums. Time how atomics scale with contention.

- **NVIDIA Focus:** Modern NVLink has native atomic support across GPUs, enabling GPU clusters (e.g. graph processing) to scale lock-free. Example: NVIDIA's Hopper NVSwitch uses atomics for all-reduce with 900 GB/s multi-GPU fabric [10] .

## Week 7: Split-Transaction Buses and Concurrency

*Story:* Split transactions allow higher throughput (like NVLink's pipelined links). Just as NVSwitch splits flows in Summit supercomputer, split buses let requests and responses overlap.

- **Subtopics:** Mechanics of split-transaction bus (address and data phases); handling pending transactions and flow-control. Interaction with multi-level cache (responses may come out-of-order).
- **Key Concepts:** Overlapping transactions improves bandwidth; buffering and retries to handle conflicts. Modeling performance (effective latency with pipelining).
- **Hands-on:** Use CUDA streams to simulate overlapping data transfers (e.g. `cudaMemcpyAsync` with streams) to see throughput gains. Implement a toy request/response protocol in CUDA to see pending-transaction management.
- **NVIDIA Focus:** NVLink 3/4.0 splits data lanes across GPUs for up to ~1.8 TB/s bidirectional per link. In practice, DGX systems with NVSwitch treat GPU communication like a fat-tree network (all-to-all) to maximize throughput [11] .

## Week 8: Scalable Shared Memory – Directory Coherence Protocols

*Story:* As GPU clusters grew, directory protocols (SGI Origin-like) were needed instead of broadcasts. NVIDIA's multi-GPU DGX A100 (2020) uses directory-like tracking over NVSwitch.

- **Subtopics:** Directory-based coherence (centralized memory-based vs distributed cache-based). Sharer lists and states (Modified, Shared). Directory operations: lookups, updates, invalidations.
- **Key Concepts:** Scalability of directory (storage overhead per line); correctness issues (deadlock, livelock avoidance). Hybrid protocols (snoop+directory).
- **Hands-on:** Experiment with CUDA Unified Memory across multiple GPUs in Colab (with NCCL or cuBlas multi-GPU). Observe when pages migrate vs stay resident (e.g. `cudaMemcpyPeer` ).
- **NVIDIA Focus:** NVLink-C2C (Grace Hopper) offers hardware coherence with a directory-like fabric: GPUs can transparently share data with CPU/NIC, greatly reducing manual traffic [12] [13] . For example, Grace Hopper can access all 150 TB of memory across 256 GPUs coherently [14] .

## Week 9: Case Study: Legacy NUMA Systems vs GPU NUMA

*Story:* Compare classic NUMA machines (SGI Origin, Sequent NUMA-Q) to modern "NUMA-aware GPUs." NVIDIA's Grace CPU+GPU platform uses NVLink directories for NUMA coherence.

- **Subtopics:** SGI Origin (hub-based cache-coherent NUMA) vs Sequent NUMA-Q (IQ-Link interconnect); protocols used (SCI, cache-only memory). Latency/bandwidth in NUMA. Comparison of programming on GPU versus NUMA CPU.
- **Key Concepts:** NUMA effects (memory affinity, page migration); correctness (data race avoidance). Verification glimpses for NUMA protocols.

- **Hands-on:** Benchmark a NUMA scenario in software: e.g., write a simple CUDA code where threads on different GPUs process local vs remote data, and compare access times (if multi-GPU available). Or simulate NUMA in Python (psutil numa).
- **NVIDIA Focus:** Modern NVIDIA systems outperform old NUMA by large margins. For instance, an 8-GPU DGX A100 achieves ~5 PFLOPS (DGX A100 with eight A100s) [11] , far above legacy NUMA boxes. Grace Hopper Superchip builds a coherent NUMA with GPUs that can run quantum-material simulations 3× faster than Origin-era systems (story).

## Week 10: Memory Consistency Models (Sequential vs Relaxed)

*Story:* Strict "sequential consistency" is intuitive but slow; GPUs use relaxed models (like CUDA's memory model introduced in PTX 2008). This flexibility fueled NVIDIA's AI speedups.

- **Subtopics:** Sequential consistency definition (Lamport). Relaxed models: release consistency, weak ordering. CUDA memory model (scopes like `__threadfence_block` ) and C++11 atomics on GPU.
- **Key Concepts:** Synchronization (fences, `volatile` , atomics) to ensure ordering. Data races and their avoidance. Trade-offs: ease of reasoning vs performance.
- **Hands-on:** Write a CUDA or OpenMP example with a deliberate reorder bug (e.g. missing fence) and fix it. Use `__threadfence()` in CUDA to enforce order across threads.
- **NVIDIA Focus:** PTX allows aggressive reordering; correctness requires careful use of `__syncthreads` and atomics. This weak model lets Hopper train LLMs ~10× faster [15] (but bugs can occur if unsynced).

## Week 11: Synchronization in GPUs

*Story:* Synchronization primitives keep parallel worlds orderly. NVIDIA's `__syncthreads()` (2007) simplified block-level barriers; later cooperative groups (Pascal+) enabled grid-wide sync.

- **Subtopics:** LL/SC vs atomic locks; point-to-point (producer-consumer queues) vs barriers. CUDA barriers: warp-level ( `__syncwarp()` ), block-level ( `__syncthreads()` ), grid-level (cooperative groups or CUDA 9+ global barriers).
- **Key Concepts:** Deadlock/starvation issues; efficient barrier algorithms (tree barriers, combining networks). GPU-specific: __barrier_sync, `cudaDeviceSynchronize()` .
- **Hands-on:** Implement a custom barrier using atomics (for a grid) and compare to `__syncthreads()` . Use CUDA Cooperative Groups to sync multiple blocks.
- **NVIDIA Focus:** Pascal introduced Cooperative Groups for flexible synchronization, enabling GPU-wide collective operations. For example, inter-block global barriers can double performance of some multi-block workloads (Nvidia whitepaper).

# Week 12: Interconnects – Topologies, Routing, Flow Control

*Story:* Interconnects are the "highways" of modern AI. NVLink (2014) began NVIDIA's high-speed links; NVSwitch (2018) created a fat-tree for DGX; NVLink 4.0 (Hopper) gives 1.8 TB/s per GPU. These topologies let supercomputers like OLCF Summit/Frontier share data at enormous rates.

- **Subtopics:** Network topologies (bus, ring, mesh, torus, fat-tree); on-chip interconnect (crossbar vs mesh inside GPUs). GPU interconnects: PCIe (legacy), NVLink, NVSwitch (in DGX) and NVLink-C2C (CPU-GPU). Routing strategies (deterministic, adaptive) and flow control (credit-based backpressure).
- **Key Concepts:** Bisection bandwidth and latency; congestion management (QoS, buffering). Communication primitives (NCCL for collectives on GPUs). GPU cluster fabric: InfiniBand, Ethernet + NVIDIA Magnum IO.
- **Hands-on:** Use NCCL (NVIDIA Collective Communications Library) in Colab (with multi-GPU runtime) to perform an all-reduce or broadcast; measure scaling. Or simulate simple mesh routing in code.
- **NVIDIA Focus:** NVLink 4.0 in Blackwell offers 900 GB/s bidirectional per GPU link. For example, DGX A100 interconnects 8 GPUs at 4.8 TB/s total via NVSwitch [11] . NVLink-C2C in Grace Hopper unifies CPU/GPU memory (900 GB/s) [12] , letting GPUs directly load/store CPU memory. These fabrics enable trillion-parameter AI models by scaling GPU memory and compute.

**Additional Notes:** Throughout the course, we will highlight real-world cases and libraries: e.g. RAPIDS (cuDF/cuML) for GPU data science, TensorRT for inference, and domain examples (GPUs in genomics, medical imaging, autonomous vehicles, etc.). Google Colab's free NVIDIA GPUs (T4/P100) will be used for demos. Students will also explore group projects (e.g. optimizing a CUDA kernel, designing a multi-GPU configuration) and guest lectures to connect theory with cutting-edge practice.

**Sources:** Content is drawn from NVIDIA developer blogs and technical papers [16] [2] [4] [5] [17] [18] [1] [7] [8] [11] [12] , illustrating GPU performance and architecture facts. (If any topic lacked a recent source, it is based on standard knowledge of GPU design.)

---

[1] [7] Nvidia - Wikipedia
https://en.wikipedia.org/wiki/Nvidia

[2] Exploring CPU vs GPU Speed in AI Training: A Demonstration with TensorFlow | Microsoft Community Hub
https://techcommunity.microsoft.com/blog/azurehighperformancecomputingblog/exploring-cpu-vs-gpu-speed-in-ai-training-a-demonstration-with-tensorflow/4014242

[3] [4] Leveraging the Power of GPUs with CuPy in Python - KDnuggets
https://www.kdnuggets.com/leveraging-the-power-of-gpus-with-cupy-in-python

[5] Medical | Industrial Introduction | GPGPU & Embedded Edge AI Computing Solutions || Aetina Corporation
https://www.aetina.com/application.php?t=45

[6] [16] Boost Alphafold2 Protein Structure Prediction with GPU-Accelerated MMseqs2 | NVIDIA Technical Blog
https://developer.nvidia.com/blog/boost-alphafold2-protein-structure-prediction-with-gpu-accelerated-mmseqs2/

[8] [10] [15] NVIDIA Hopper Architecture In-Depth | NVIDIA Technical Blog
https://developer.nvidia.com/blog/nvidia-hopper-architecture-in-depth/

[9] [12] [13] [14] NVIDIA Grace Hopper Superchip Architecture In-Depth | NVIDIA Technical Blog
https://developer.nvidia.com/blog/nvidia-grace-hopper-superchip-architecture-in-depth/

[11] images.nvidia.com
https://images.nvidia.com/aem-dam/en-zz/Solutions/data-center/nvidia-ampere-architecture-whitepaper.pdf

[17] Army Research Lab boosts computer processing using gaming technology | Article | The United States Army
https://www.army.mil/article/46837/army_research_lab_boosts_computer_processing_using_gaming_technology

[18] new.math.uiuc.edu
https://new.math.uiuc.edu/im2008/dakkak/proposal/papers/michalakes_lspp.pdf