

1. Mechanism of Communication

PCIe (Peripheral Component Interconnect Express)

- **Wires and Protocol:**
 - **Parallel Communication:** Uses **lanes** (x1, x2, x4, x8, x16) for data transfer.
 - **Lines:** Each lane consists of **two differential pairs** (TX and RX) and a separate clock signal.
 - **Protocol:** PCIe uses a **packet-based protocol** with **transaction layers**, **data link layers**, and **physical layers**.
 - **Clock:** Separate clock signal (synchronous).
- **Throughput:**
 - PCIe Gen 3: **8 GT/s per lane** (\approx 985 MB/s per lane).
 - PCIe Gen 4: **16 GT/s per lane** (\approx 1.97 GB/s per lane).
 - PCIe Gen 5: **32 GT/s per lane** (\approx 3.94 GB/s per lane).
 - Overhead: Includes packet headers, error checking, and protocol overhead.
- **Controller/Device Config:**
 - **Multi-Controller:** PCIe supports **multiple endpoints** (devices) connected to a **root complex** (controller).
 - **Arbitration:** Uses **hardware-based arbitration** to manage access to the bus.

SATA (Serial Advanced Technology Attachment)

- **Wires and Protocol:**
 - **Serial Communication:** 1 bit at a time.
 - **Lines:** SATA uses **7 wires** (2 differential pairs for data TX/RX, 3 ground wires).
 - **Protocol:** SATA uses a **packet-based protocol** with separate control and data frames.
 - **Clock:** Embedded clocking (no separate clock line; clock is recovered from the data stream).
- **Throughput:**
 - SATA III: **6 Gbps** (\approx 600 MB/s).
 - Overhead: Includes control frames, error checking, and protocol overhead.
- **Controller/Device Config:**
 - **Single Controller:** SATA is a **single-master protocol**; the host (CPU) is the controller, and storage devices are slaves.
 - **No Arbitration:** Only one host communicates with the storage device at a time.

NVMe (Non-Volatile Memory Express)

- **Wires and Protocol:**
 - **Parallel Communication:** Uses PCIe lanes (x2, x4, x8, etc.).
 - **Lines:** PCIe uses **differential pairs** for TX/RX (e.g., x4 uses 8 data lines + clock and control lines).
 - **Protocol:** NVMe is a **queue-based protocol** with up to **64,000 queues**, each holding **64,000 commands**.
 - **Clock:** Separate clock signal (PCIe is synchronous).
- **Throughput:**
 - PCIe Gen 3 \times 4: **4 GB/s** (32 Gbps).
 - PCIe Gen 4 \times 4: **8 GB/s** (64 Gbps).
 - Overhead: Minimal due to efficient queueing and parallelism.
- **Controller/Device Config:**
 - **Single Controller:** NVMe is also **single-master**, with the host (CPU) as the controller.
 - **No Arbitration:** Only one host communicates with the NVMe SSD at a time.

NVLINK

- **Wires and Protocol:**
 - **Parallel Communication:** Uses **multiple high-speed lanes** (e.g., 8 or 12 lanes per link).
 - **Lines:** Differential pairs for TX/RX, with separate clock and control lines.
 - **Protocol:** Proprietary protocol optimized for **GPU-to-GPU** and **GPU-to-CPU** communication.
 - **Clock:** Separate clock signal (synchronous).
- **Throughput:**
 - NVLINK 3.0: **600 GB/s** bidirectional per link.

- Overhead: Minimal due to direct memory access and low-latency design.
 - **Controller/Device Config:**
 - **Multi-Controller:** NVLINK supports **multi-master communication** (e.g., multiple GPUs can communicate directly).
 - **Arbitration:** Uses **hardware-based arbitration** to manage access to the bus.
-

2. Types of Communication

PCIe

- **Parallel:** Uses multiple lanes (x1, x2, x4, etc.).
- **Off-Chip:** Connects devices to the motherboard.
- **Synchronous:** Uses a separate clock signal.

SATA

- **Serial:** 1 bit at a time.
- **Off-Chip:** Connects storage devices to the motherboard.
- **Synchronous:** Uses embedded clocking.

NVMe

- **Parallel:** Uses multiple PCIe lanes (e.g., x4, x8).
- **Off-Chip:** Connects SSDs to the motherboard.
- **Synchronous:** Uses a separate clock signal.

NVLINK

- **Parallel:** Uses multiple high-speed lanes.
 - **Off-Chip:** Connects GPUs and CPUs on a PCB or across systems.
 - **Synchronous:** Uses a separate clock signal.
-

3. Bus Control Methods

PCIe

- **Packet-Based:**
 - Uses **TLP (Transaction Layer Packets)** for data transfer.
 - Supports **ACK/NACK** for reliable communication.

SATA

- **Strobe-Based:**
 - The host (controller) sends a command, and the device responds within a fixed time.
 - No explicit ACK; relies on timing and error checking.

NVMe

- **Queue-Based:**
 - The host submits commands to queues, and the device processes them asynchronously.
 - Uses **completion queues** to signal when commands are done.

NVLINK

- **Handshake-Based:**
 - Uses explicit ACK signals for reliable communication.
 - Supports **memory pooling**, where GPUs can directly access each other's memory.
-

4. Burst Transfers

PCIe

- **Efficient Burst Transfers:**
 - PCIe is optimized for large, contiguous data transfers.
 - Uses **TLPs** to handle non-contiguous data.

SATA

- **Limited Burst Transfers:**
 - SATA supports small bursts but is not optimized for large contiguous transfers.
 - Overhead from control frames reduces efficiency.

NVMe

- **Efficient Burst Transfers:**
 - NVMe is optimized for large, contiguous data transfers (e.g., reading/writing large blocks of data).
 - Uses **scatter-gather lists** to handle non-contiguous data efficiently.

NVLINK

- **High-Speed Burst Transfers:**
 - NVLINK is designed for massive data transfers between GPUs and CPUs.
 - Supports **coherent memory access**, allowing GPUs to access each other's memory directly.
-

5. Bus Arbitration

PCIe

- **Priority-Based Arbitration:**
 - Uses hardware arbitration to manage access to the bus.
 - Supports **round-robin** and **fixed-priority** schemes.

SATA

- **No Arbitration:**
 - SATA is a single-master protocol; only the host communicates with the storage device.

NVMe

- **No Arbitration:**
 - NVMe is also single-master; only the host communicates with the SSD.

NVLINK

- **Priority-Based Arbitration:**
 - Uses hardware arbitration to manage access to the bus.
 - Supports **round-robin** and **fixed-priority** schemes.
-

6. Memory-Mapped I/O

PCIe

- **Memory-Mapped I/O:**
 - PCIe uses **memory-mapped registers** for communication.
 - Devices can access host memory directly.

SATA

- **Not Memory-Mapped:**
 - SATA uses a **command-based interface**; the host sends commands to the device.

NVMe

- **Memory-Mapped I/O:**
 - NVMe uses **memory-mapped registers** for command submission and completion.
 - The host writes commands to registers, and the device processes them.

NVLINK

- **Memory-Mapped I/O:**
 - NVLINK enables **direct memory access** between GPUs and CPUs.
 - GPUs can access each other’s memory as if it were their own.

Summary Table

| Feature | PCIe | SATA | NVMe | NVLINK |
|--------------------|---------------------------|---------------------|-------------------------------|-----------------------|
| Communication Type | Parallel | Serial | Parallel (PCIe) | Parallel (NVLINK) |
| Throughput | Up to 63 GB/s (Gen 5 x16) | 6 Gbps (≈ 600 MB/s) | Up to 64 Gbps (PCIe Gen 4 x4) | 600 GB/s (NVLINK 3.0) |
| Bus Control | Packet-Based | Strobe-Based | Queue-Based | Handshake-Based |
| Burst Transfers | Efficient | Limited | Efficient | High-Speed |
| Arbitration | Priority-Based | None | None | Priority-Based |
| Memory-Mapped I/O | Yes | No | Yes | Yes |

Key Takeaways

- **PCIe** is the foundational technology that enables high-speed communication for both **NVMe** (storage) and **NVLINK** (GPU/CPU communication).
- **SATA** is a simpler, slower protocol designed for legacy storage devices.
- **NVMe** leverages PCIe to provide high-speed storage access, while **NVLINK** extends PCIe-like capabilities for GPU-to-GPU and GPU-to-CPU communication.

ASIDE

Why **NVMe** is considered **single-master** even though **PCIe** supports **multi-master** communication.

TL;DR

- **PCIe** is a **multi-master protocol** that supports communication between multiple devices, including peer-to-peer transfers.
- **NVMe** is a **single-master protocol** that operates on top of PCIe and is designed for storage access, where the host (CPU) is the sole initiator of commands.
- NVMe leverages PCIe’s multi-master capabilities to coexist with other devices on the bus but does not implement multi-master or peer-to-peer communication itself.

1. PCIe as a Multi-Master Protocol

- **PCIe Architecture:**
 - PCIe is designed to support **multiple endpoints** (devices) connected to a **root complex** (host).
 - Each endpoint can act as a **master** (initiator) or **slave** (target) depending on the transaction.
 - PCIe supports **peer-to-peer communication**, where two endpoints can communicate directly without involving the root complex.
- **Multi-Master Capability:**
 - PCIe allows multiple devices to initiate transactions (e.g., GPUs, network cards, storage devices).
 - Arbitration is handled by the **root complex** or **switch** to manage access to the bus.

2. NVMe as a Single-Master Protocol

- **NVMe Architecture:**

- NVMe is a **higher-level protocol** that operates on top of PCIe.
- It defines how the host (CPU) communicates with NVMe SSDs.
- NVMe is designed specifically for **storage devices** and assumes a **single host** (CPU) controlling the SSD.
- **Single-Master Design:**
 - In an NVMe system, the **host (CPU)** is the **master**, and the **NVMe SSD** is the **slave**.
 - The host initiates all commands (e.g., read, write) and the SSD responds.
 - NVMe does not support **peer-to-peer communication** between SSDs or other devices.

Why NVMe is Single-Master

1.

Storage-Specific Design:

- NVMe is optimized for **storage access**, where the host (CPU) is the sole initiator of commands.
- There is no need for SSDs to communicate directly with each other or with other devices.

2.

Queue-Based Model:

- NVMe uses a **queue-based model** where the host submits commands to **submission queues** and the SSD processes them.
- The host is the only entity that can submit commands to the queues.

3.

No Peer-to-Peer Support:

- NVMe does not implement the mechanisms required for peer-to-peer communication (e.g., addressing other devices, arbitration between devices).
- All communication is between the host and the SSD.

How NVMe Leverages PCIe Multi-Master Capability

While NVMe itself is single-master, it relies on PCIe’s multi-master capabilities in the following ways:

1.

Multiple NVMe SSDs:

- A system can have multiple NVMe SSDs connected to the same PCIe bus.
- Each SSD is controlled by the host independently, but they share the PCIe bus.
- PCIe’s arbitration ensures fair access to the bus.

2.

Other PCIe Devices:

- In a system with NVMe SSDs and other PCIe devices (e.g., GPUs, network cards), PCIe’s multi-master capability allows these devices to coexist.
- The host can communicate with NVMe SSDs and other devices simultaneously.

3.

DMA (Direct Memory Access):

- NVMe SSDs use PCIe’s DMA capabilities to transfer data directly to/from host memory.
- This is managed by the host, not the SSD.

Key Differences Between PCIe and NVMe

| Feature | PCIe | NVMe |
|---------------------|----------------------------------|-----------------------------------|
| Scope | Physical and link layer protocol | Higher-level protocol for storage |
| Multi-Master | Yes | No (single-master) |
| Peer-to-Peer | Supported | Not supported |
| Communication Model | Host ↔ Device or Device ↔ Device | Host ↔ SSD only |
| Queues | Not applicable | Submission and completion queues |

Example Scenario

- A system has:
 - 1 CPU (host).
 - 2 NVMe SSDs.
 - 1 GPU.
- **PCIe Multi-Master:**
 - The CPU can communicate with both NVMe SSDs and the GPU.
 - The GPU can also communicate with the CPU or perform peer-to-peer transfers with other GPUs.
- **NVMe Single-Master:**
 - Each NVMe SSD is controlled exclusively by the CPU.
 - The SSDs cannot communicate directly with each other or with the GPU.