

# Report

## Analytical Report: Scalability, Computational Performance, and Modern Architectures

---

### 1. Scalability and Performance Limits

#### 1.1. Concept

Scalability is the ability of a system to handle increasing workloads while maintaining efficiency. Performance limits are physical and algorithmic barriers that prevent linear performance growth when adding resources.

#### 1.2. Key Limits and Analysis

- **Amdahl's Law:** Maximum speedup is limited by the sequential fraction of work.

*P = parallel fraction, N = number of processors.*

$$\text{Speedup} \leq \frac{1}{(1 - P) + \frac{P}{N}}$$

- **Communication & Synchronization Overhead:** Multiple cores must share data and synchronize states, which causes latency (e.g., cache coherence).
- **Memory Wall:** With many cores accessing RAM, bandwidth saturates. This bottleneck reduces scaling. High-bandwidth memory (HBM) is used in modern GPGPUs to mitigate this.
- **Load Imbalance:** Uneven task distribution leads to idle cores.
- **Resource Contention:** Shared resources like system bus, L3 cache, or GPU SMs create conflicts, reducing throughput.

## 2. My personal computer ability

**Model:** HP Victus 16-e0xxx

- **CPU:** AMD Ryzen 5 5600H (6C/12T, 3.3–4.2 GHz)
- **GPU:** NVIDIA GTX 1650 (4 GB GDDR6) + AMD Radeon iGPU
- **RAM:** 20 GB DDR4 (16 GB + 4 GB, Flex Mode)
- **Storage:** Samsung NVMe SSD 512 GB (~22 GB free on C:)

### 2.1. Analysis

- **Strengths:**

- Multicore CPU (Zen 3) delivers solid multithreading performance.
  - Adequate RAM and fast SSD ensure smooth performance.
- **Weaknesses:**
    - GPU lacks Tensor Cores, only 4 GB VRAM → unsuitable for large AI training.
    - Flex Mode reduces memory bandwidth beyond 8 GB.
    - Low SSD free space impacts I/O performance.

## 2.2. Recommendations

- Replace 4 GB RAM with another 16 GB stick → 32 GB full dual-channel.
- Free or upgrade SSD.
- Use cloud GPU services (AWS, GCP, Colab Pro) for deep learning.

## 3. Parallel Processing Architectures

### 3.1. Thread Scheduling and CPU Affinity

- **Thread Scheduling:** OS decides which thread runs on which core.

- **CPU Affinity:** Pinning a thread to specific cores improves cache locality but may cause load imbalance if misconfigured.

### 3.2. Multicore vs MIC vs GPGPU

	<b>Architecture Features</b>	<b>Advantages</b>	<b>Limitations</b>
<b>Multicore CPU</b>	4–64 complex cores, large cache, SIMD (AVX).	Low latency, flexible, handles complex branching.	Limited scaling, resource contention.
<b>MIC</b>	Many simple x86 cores, wide vector units.	Parallel-friendly, easy CPU code porting.	Largely obsolete, outperformed by GPUs.
<b>GPGPU</b>	Hundreds–thousands of simple SIMT cores, high-bandwidth memory.	Massive throughput, ideal for AI/data-parallel transfer workloads.	Poor with branching code, high CPU↔GPU overhead.

## **Multicore CPUs (Central Processing Units)**

- **Organization:** A multicore CPU has a relatively small number of powerful, complex cores (4–64). Each core offers high single-thread performance, large caches, and strong support for branch-heavy code.
- **Strengths:** Very flexible. Best choice for operating system tasks, office applications, and programs with complex, hard-to-parallelize logic.
- **Limitations:** Scalability is constrained by Amdahl's Law and shared memory contention.

**Example: Intel Core i7, AMD Ryzen 7** in desktops or laptops.

- **Cores:** 4–16 powerful, complex cores.
- **Independence:** Each core has private cache and can run tasks independently.
- **Purpose:** Optimized for multitasking (web browsing, office apps, compilation, gaming).
- **Programming:** Standard multithreading libraries (**OpenMP, Pthreads**).

## **MIC (Many Integrated Core)**

- **Organization:** Integrates dozens to hundreds of simpler cores on a single chip. Example: Intel Xeon Phi. These

cores are weaker than traditional CPU cores but effective for vectorized execution.

- **Strengths:** Good for highly parallel, vectorizable workloads. Easier to port CPU code compared to GPGPU in some cases.
- **Limitations:** Now rare. Outperformed by GPGPU in raw compute and memory bandwidth.

**Example: Intel Xeon Phi** (discontinued).

- **Cores:** 60–72 simpler x86 cores.
- **Independence:** Weaker than CPU cores but effective with wide vector units and shared memory.
- **Purpose:** High-performance computing (HPC), machine learning, big data workloads.
- **Programming:** x86-compatible, easier to port CPU code.

## **GPGPU (General-Purpose Graphics Processing Units)**

- **Organization:** Contains hundreds to thousands of simple cores optimized for SIMD (Single Instruction, Multiple Threads). Equipped with extremely high-bandwidth memory (e.g., HBM) to feed all cores simultaneously.

- **Strengths:** Extremely high throughput for data-parallel tasks. Dominant in AI/ML, graphics rendering, and scientific simulations.
- **Limitations:** Inefficient for branch-heavy or complex logic. Overhead exists when transferring data between system RAM and GPU VRAM.

### **Example: NVIDIA RTX 4090, NVIDIA A100.**

- **Cores:** Thousands of simple CUDA cores/stream processors.
- **Independence:** Operates under **SIMT** (Single Instruction, Multiple Threads).
- **Purpose:** Graphics, deep learning, crypto mining, scientific simulations.
- **Programming:** **CUDA** (NVIDIA) or **OpenCL**.

## **4. NVIDIA Research Frontiers**

- **AI Model Scaling:** *ZeRO-Infinity* leverages GPU+CPU+NVMe for massive models.
- **Agentic AI:** Research on *Small Language Models* for efficient agents.
- **Scaling Laws in AI:** Predictable performance improvements when increasing compute/data.

- **GPU Architecture Insights:** Papers like *Dissecting the Volta GPU Architecture* help optimize software.

- *Conclusion and Outlook:*

Scalability is constrained by theoretical laws, hardware design, and bottlenecks. Choosing the right architecture for the workload is crucial: CPUs for low-latency complex tasks, GPGPUs for large-scale parallelism.

For the analyzed personal system, the optimal path is:

- Maximize local resources for development and DevOps.
- Leverage cloud GPU infrastructure for large-scale AI experiments.

This hybrid approach aligns with modern computing trends in both HPC and AI.