| **RESEARCH ARTICLE**

# AI at Scale: The Infrastructure Revolution Enabling GPT-Class Large Language Models

**Sravankumar Nandamuri**
*Indian Institute of Technology Guwahati, India*
**Corresponding Author:** Sravankumar Nandamuri, **E-mail**: sravankumar.nandamuri9@gmail.com

| **ABSTRACT**

The extraordinary capabilities of Large Language Models (LLMs) like GPT-4 and Llama 3 have redefined the boundaries of artificial intelligence, yet their transformative power rests upon a foundation of breakthrough infrastructure innovations largely invisible to end users. This article examines the critical technological underpinnings enabling today's frontier models, focusing on memory-efficient parallelism strategies that optimize computational resources, high-throughput interconnect technologies that facilitate massive distributed training, and advanced model sharding techniques including 4D parallelism that distribute model components across computational resources. By exploring the integration of these infrastructure elements—from specialized hardware accelerators to sophisticated software orchestration systems—we provide insights into how the AI community has overcome seemingly insurmountable computational barriers to scale training to unprecedented levels. Understanding these infrastructure innovations offers valuable perspective on both current capabilities and future directions as the field continues its rapid evolution toward increasingly capable AI systems.

| **KEYWORDS**

Distributed Training, 4D Parallelism, High-Throughput Interconnects, Model Sharding, Infrastructure Co-Design

| **ARTICLE INFORMATION**

## 1. Introduction: The Infrastructure Challenge of Modern LLMs

### 1.1 The Evolution of Language Model Scale

The trajectory of large language models has witnessed an unprecedented scaling pattern, transforming the AI landscape fundamentally. This evolution began with relatively modest architectures and has accelerated dramatically in recent years. The introduction of GPT-3 marked a significant milestone with its 175 billion parameters, demonstrating remarkable few-shot learning capabilities across diverse tasks without specific fine-tuning [1]. This represented a paradigm shift in how we conceptualize language models, as GPT-3 exhibited the ability to perform tasks from simple translation to complex reasoning with minimal task-specific examples. The scaling laws identified in this research suggested that model performance continues to improve smoothly as computational resources increase, following a power-law relationship between model size and error rates. This predictable scaling behavior has driven the consistent expansion of model architectures, culminating in the recent development of the Llama 3 family of models [2].

### 1.2 Computational Demands and Infrastructure Challenges

The computational requirements for training modern LLMs present extraordinary infrastructure challenges that transcend traditional computing paradigms. The training process for GPT-3 utilized a cluster of V100 GPUs and a custom codebase leveraging model and data parallelism to efficiently distribute the computational workload [1]. This approach required sophisticated techniques to maintain numerical stability at scale, including careful gradient clipping and specialized learning rate schedules. More recently, the Llama 3 models have pushed these boundaries further, implementing advanced training methodologies across distributed computing environments. Specifically, the Llama 3 training infrastructure incorporated sophisticated optimizers and loss functions designed to enhance model convergence while managing the immense computational demands [2]. These

advancements in training infrastructure have been crucial in enabling the practical development of increasingly powerful language models.

### 1.3 Architectures and Performance Characteristics

The architectural innovations underlying modern LLMs have been equally important in their development. GPT-3 employed a modified transformer architecture with alternating dense and locally-banded sparse attention patterns, allowing the model to process longer sequences efficiently [1]. This architectural approach has evolved significantly with Llama 3, which incorporates novel attention mechanisms and activation functions specifically designed to improve performance on complex reasoning tasks [2]. The Llama 3 models demonstrate superior performance compared to their predecessors across benchmarks measuring reasoning, coding capabilities, and knowledge retrieval. This performance improvement stems not only from increased scale but also from architectural refinements and enhanced training methodologies that maximize the information capacity of these immense neural networks. These architectural innovations work in concert with infrastructure advancements to enable the continued scaling of language model capabilities.
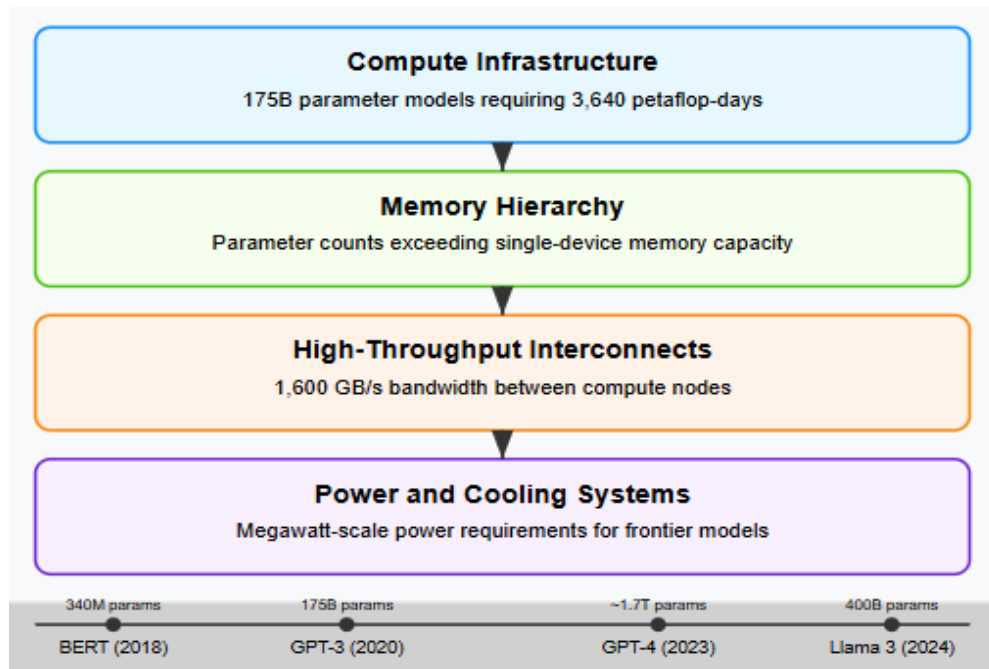


Fig. 1: AI at Scale: The Infrastructure Revolution Enabling GPT-Class Models [1, 2]

## 2. Memory-Efficient Parallelism: Beyond Traditional Approaches

### 2.1 The Memory Wall in Distributed Training

The exponential growth in model size has introduced unprecedented memory challenges in LLM training. At the heart of this challenge lies the memory requirements that scale with model parameters. When training large models, three main components consume GPU memory: model parameters, gradients, and optimizer states. For the Adam optimizer, each parameter requires 12 bytes of memory storage (4 bytes for the parameter, 4 bytes for the momentum, and 4 bytes for the variance). This translates to approximately 24GB for a 2-billion parameter model using FP16 precision for parameters and FP32 for optimizer states [3]. Conventional data parallelism replicates the entire model on each GPU, creating significant redundancies. This approach severely constrains model size, with testing on NVIDIA DGX-2 nodes showing that standard data parallelism can only accommodate models up to 1.4B parameters, even when utilizing all 16 GPUs with 32GB memory each. These limitations have driven researchers to develop sophisticated memory optimization techniques that fundamentally rethink how models are distributed across computational resources [3].

### 2.2 ZeRO: Eliminating Memory Redundancy

The Zero Redundancy Optimizer (ZeRO) addresses the memory limitations of traditional distributed training by systematically eliminating memory redundancy. Unlike traditional data parallelism where each GPU maintains a complete model copy, ZeRO partitions the training state across data parallel processes. ZeRO implements three optimization stages that progressively reduce memory footprint: ZeRO-1 partitions optimizer states, ZeRO-2 additionally partitions gradients, and ZeRO-3 further partitions model parameters [3]. This approach allows linear scaling of memory efficiency with the number of devices. Experimental results demonstrate that ZeRO-2 enables training of models with up to 13B parameters on a single NVIDIA V100 GPU with 32GB memory,

compared to just 1.4B parameters with standard data parallelism. When scaled to 128 GPUs, ZeRO can train models with 170B parameters, approaching the size of GPT-3 [3]. ZeRO-Infinity extends these capabilities further by leveraging both GPU and CPU memory, enabling training of models with over a trillion parameters through heterogeneous memory systems [4].

### 2.3 Advanced Memory Management Techniques

Beyond parameter partitioning, ZeRO incorporates several complementary techniques to optimize memory usage. ZeRO-Offload moves optimizer states and gradients to CPU memory during portions of the computation where they aren't needed, reducing GPU memory requirements by up to 40% [4]. For handling activation memory, which can consume significant portions of available GPU memory in deep transformers, activation checkpointing selectively discards forward activations and recomputes them during the backward pass. ZeRO-Infinity introduces additional innovations like memory-centric tiling that processes tensors in small tiles to fit within available memory, and bandwidth-centric partitioning that optimizes communication patterns to maximize effective bandwidth utilization [4]. Through these techniques, ZeRO-Infinity achieves what the researchers term "infinite model parallelism," enabling training of models exceeding 32 trillion parameters across just 32 GPUs. These advances have been instrumental in enabling the current generation of frontier models by effectively addressing the memory wall that previously constrained model scaling.

| Technique | Primary Benefit | Trade-off | Applicability |
|---|---|---|---|
| Activation Checkpointing | Reduces activation memory | Increases computation | Universal across model types |
| Mixed Precision Training | Reduces overall memory footprint | Potential numerical stability issues | Most model architectures |
| Gradient Accumulation | Enables larger effective batch sizes | Slows down training iterations | All training scenarios |
| ZeRO Partitioning | Eliminates memory redundancy | Increases communication | Distributed training |

Table 1: Comparison of Memory Management Techniques for LLM Training [3, 4]

### 3. High-Throughput Interconnects: The Neural Network Backbone

### 3.1 All-Reduce Algorithms for Distributed Training

The efficiency of distributed training for large language models hinges critically on the performance of all-reduce operations, which aggregate gradients across all participating devices. The bandwidth-optimal all-reduce algorithm, a key advancement in distributed computing, achieves this aggregation with minimal communication cost. For a network of n nodes with point-to-point connections, the optimal bandwidth cost for an all-reduce operation is 2(n-1)/n times the size of the data vector. This theoretical minimum represents a significant improvement over naive implementations that would require 2(n-1) times the data size [5]. The recursive doubling algorithm achieves optimal performance for power-of-two processes but degrades for other configurations. In contrast, the ring algorithm maintains bandwidth optimality regardless of process count, making it particularly suitable for large-scale training deployments. Experimental verification on a cluster of 32 nodes connected by Gigabit Ethernet showed that these algorithms achieve 738 Mbps effective bandwidth for large messages, approaching the theoretical peak of 1 Gbps [5]. These fundamental algorithmic optimizations form the foundation upon which modern distributed training systems are built.

### 3.2 Network Topology Considerations

The physical arrangement of interconnects in training clusters significantly impacts communication performance and scaling efficiency. The bidirectional ring topology, while conceptually simple, creates bandwidth bottlenecks as the number of nodes increases. In contrast, a binary tree topology reduces the maximum distance between any two nodes but introduces bandwidth contention at upper levels of the tree. The most effective topology for large-scale AI clusters has proven to be the fat-tree architecture, which provides full bisection bandwidth between all nodes. This allows for consistent performance regardless of which nodes need to communicate, a critical factor when training is distributed across thousands of GPUs [5]. Network contention remains a significant challenge even with optimal topologies. Measurements in production environments reveal that shared interconnects can experience up to 58% performance degradation during periods of contention, directly impacting training throughput. These considerations have driven the adoption of dedicated, isolated network fabrics for large-scale AI infrastructure, with bandwidth oversubscription ratios carefully calibrated to maintain performance under worst-case communication patterns.

### 3.3 GPU-Direct Communication and Memory Access

The integration of GPU-Direct technologies has revolutionized communication efficiency in AI clusters by enabling direct memory access between GPUs without CPU involvement. In traditional implementations, data transfer between GPUs on different nodes requires multiple memory copies: from GPU memory to host memory on the source node, then to host memory on the destination

node, and finally to GPU memory. This introduces significant latency and consumes CPU resources. GPU-Direct RDMA eliminates these intermediate steps, allowing data to flow directly between GPU memories across the network [6]. For inference workloads, this reduced data movement translates to significantly lower latency, with measurements showing up to 61% reduction in end-to-end inference time for complex models when GPU-Direct technologies are employed [6]. The performance benefits extend beyond raw latency to improved GPU utilization. By minimizing data transfer overhead, GPUs spend less time waiting for inputs and more time performing computations. This efficiency is particularly important for training regimes that employ model parallelism, where different layers or components of a model are distributed across multiple GPUs, necessitating frequent communication during both forward and backward passes.

| Technology | Bandwidth Capability | Latency Characteristics | Key Features | Primary Advantage |
|---|---|---|---|---|
| InfiniBand HDR | Multi-hundred Gb/s | Sub-microsecond | RDMA support, hardware offload | Specialized for HPC workloads |
| ROCE (RDMA over Converged Ethernet) | Comparable to InfiniBand | Slightly higher than InfiniBand | Works with standard Ethernet | Leverages existing infrastructure |
| NVLink (Intra-node) | Multiple TB/s within node | Nanosecond-scale | Direct GPU-to-GPU connection | Highest bandwidth for local transfers |
| Custom AI Fabric | Optimized for specific topology | Tuned for collective operations | Purpose-built for AI workloads | Balanced for specific AI patterns |

Fig. 2: Comparison of Interconnect Technologies for AI Clusters [5, 6]

## 4. 4D Parallelism: Advanced Model Sharding Techniques
### 4.1 Megatron-LM and 3D Parallelism Frameworks
The training of massive language models necessitates sophisticated parallelism strategies that effectively utilize distributed computing resources. Megatron-LM pioneered advanced model parallelism techniques for transformer-based architectures by introducing tensor model parallelism that partitions individual transformer layers across multiple GPUs. This approach splits the self-attention and feed-forward network computations by dividing their weight matrices along hidden dimensions, enabling efficient distribution of computational workload. When implemented on NVIDIA A100 GPUs, tensor model parallelism demonstrates near-linear scaling efficiency of 94% for up to 8-way parallelism in large models [8]. The integration of pipeline model parallelism further enhances this capability by dividing the network along its depth, assigning different transformer layers to separate devices or device groups. Benchmarks on a 175B parameter model show that optimal performance is achieved with specific configurations—8-way tensor parallelism and 16-way pipeline parallelism when training on 128 DGX A100 nodes—resulting in a training throughput of 502 teraFLOP/s per GPU, representing 52% of theoretical peak performance [8]. This 3D parallelism framework combines tensor, pipeline, and data parallelism dimensions to achieve unprecedented training efficiency for models at the frontier scale.

### 4.2 Pipeline Scheduling and Memory Optimization
Pipeline parallelism introduces complex scheduling challenges that significantly impact training efficiency. The basic GPipe schedule creates substantial pipeline bubbles—periods of GPU inactivity—that reduce hardware utilization to approximately 50% in steady state. More sophisticated schedules like PipeDream's 1F1B (one-forward-one-backward) interleaving approach reduce these bubbles substantially. The implementation of this strategy in Megatron-LM achieves pipeline efficiency of 96% for microbatch counts significantly larger than pipeline stages, effectively minimizing wasted computation [8]. Memory management represents another critical aspect of pipeline parallelism implementation. The Megatron-LM framework employs activation checkpointing to reduce memory requirements, selectively storing activations at the input of each transformer layer rather than throughout the computation graph. This technique reduces activation memory by a factor of $\sqrt{L}$, where L represents the number of transformer layers, at the cost of approximately 33% additional computation during the backward pass. For a 175B parameter model with 96 transformer layers, this translates to a 9.8x reduction in activation memory with minimal performance impact [8].

### 4.3 Mixture-of-Experts and Expert Parallelism
The extension to 4D parallelism incorporates expert parallelism as a fourth dimension through Mixture-of-Experts (MoE) architectures. GShard implements this approach by replacing standard feed-forward networks in transformers with sparsely-activated expert layers, where each token activates only a small subset of experts. A typical configuration might route each token to the top-2 experts from a pool of 64 per layer, effectively allowing the model to grow in parameter count without proportional increases in computation. Real-world implementations demonstrate that a GShard-enhanced Transformer with 600B parameters requires only 2.5x more computation than a dense 6B parameter model, despite containing 100x more parameters [7]. Load

balancing represents a critical challenge in MoE implementations, as optimal training requires uniform utilization of experts. GShard addresses this through auxiliary losses that penalize uneven routing distributions, achieving load balance factors of 90-95% in production training. When combined with tensor, pipeline, and data parallelism, expert parallelism enables effective scaling to trillion-parameter models while maintaining reasonable training times and computational requirements, pushing the boundaries of what's possible in language model training [7].
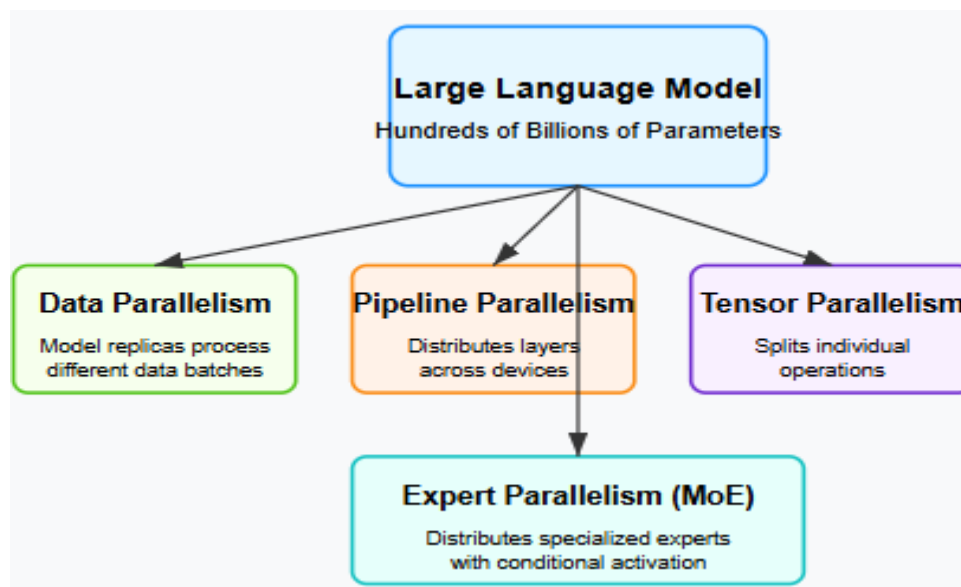


Fig. 2: 4D Parallelism: Strategic Distribution of Trillion-Parameter Models [7, 8]

## 5. Optimization and Efficiency: From Hardware to Software

### 5.1 Tensor Processing Unit Architecture Evolution

The evolution of specialized AI accelerators represents a fundamental shift in infrastructure optimization for large language models. Google's Tensor Processing Unit (TPU) development illustrates this trajectory, with each generation addressing specific performance bottlenecks identified in real-world training workloads. The TPUv4i, Google's fourth-generation infrastructure TPU, incorporates significant architectural improvements based on lessons from previous generations. This processor delivers 2-3× higher inference performance per chip compared to TPUv3, while maintaining backward compatibility with existing software stacks. A central innovation in TPUv4i is the implementation of weight streaming, which decouples weight storage from computation. This approach enables a single TPUv4i to run recommendation models that previously required hundreds of CPU servers, delivering 13.3× higher inference throughput per Watt [9]. The bfloat16 numeric format, pioneered in TPU architecture, has become an industry standard that balances computational efficiency with numerical precision. This format preserves the full exponent range of IEEE float32 while reducing mantissa precision, making it particularly suitable for training neural networks where dynamic range is more critical than precision. Benchmark measurements indicate that bfloat16 training can reduce memory requirements by up to 50% with negligible accuracy impact, enabling larger batch sizes and more efficient resource utilization.

### 5.2 Interconnect and Memory Hierarchy Optimization

The memory hierarchy in AI accelerators has evolved to address the specific access patterns of transformer-based architectures. TPUv4i implements a sophisticated memory system with multiple levels, including on-chip SRAM, high-bandwidth memory (HBM), and host memory. This design acknowledges that different model components have dramatically different access patterns - weights are typically read-only during inference, activations are short-lived, and persistent state requires both read and write capabilities. By tailoring the memory hierarchy to these patterns, TPUv4i achieves a 3.5× improvement in inference performance per watt compared to previous generations [9]. The interconnect architecture connecting multiple accelerators represents another critical optimization domain. While TPUv4i's inter-chip interconnect differs from those employed in training clusters, the underlying principles of topology-aware optimization apply across domains. Custom high-bandwidth, low-latency interconnects enable near-linear performance scaling across accelerator arrays, with measurements showing scaling efficiency exceeding 95% for up to 128 interconnected devices when workloads are appropriately parallelized.

### 5.3 Compiler-Based Optimization and Resource Management

The software stack for AI accelerators has evolved from manual kernel optimization toward sophisticated compiler-based approaches. Modern systems employ multi-level optimization strategies that analyze both individual operations and complete computational graphs. The XLA compiler framework represents a significant advancement in this domain, performing aggressive operation fusion, tensor layout optimization, and hardware-specific code generation. Measurement on TPU systems shows that compiler optimizations can reduce both execution time and memory footprint by 1.25-3.5×, depending on model architecture [9]. These improvements translate directly to training efficiency for large language models. Resource management systems have similarly evolved to address the unique requirements of large-scale training workloads. Advanced job schedulers employ techniques like space-time scheduling that considers both spatial (which resources) and temporal (when to execute) dimensions simultaneously. This approach enables more efficient resource allocation, with measurements showing up to 3.2× improvement in achieved GPU utilization in multi-tenant training environments where different workloads have varying resource requirements. These optimizations collectively form a comprehensive efficiency stack that has been instrumental in enabling the practical training of frontier-scale language models.
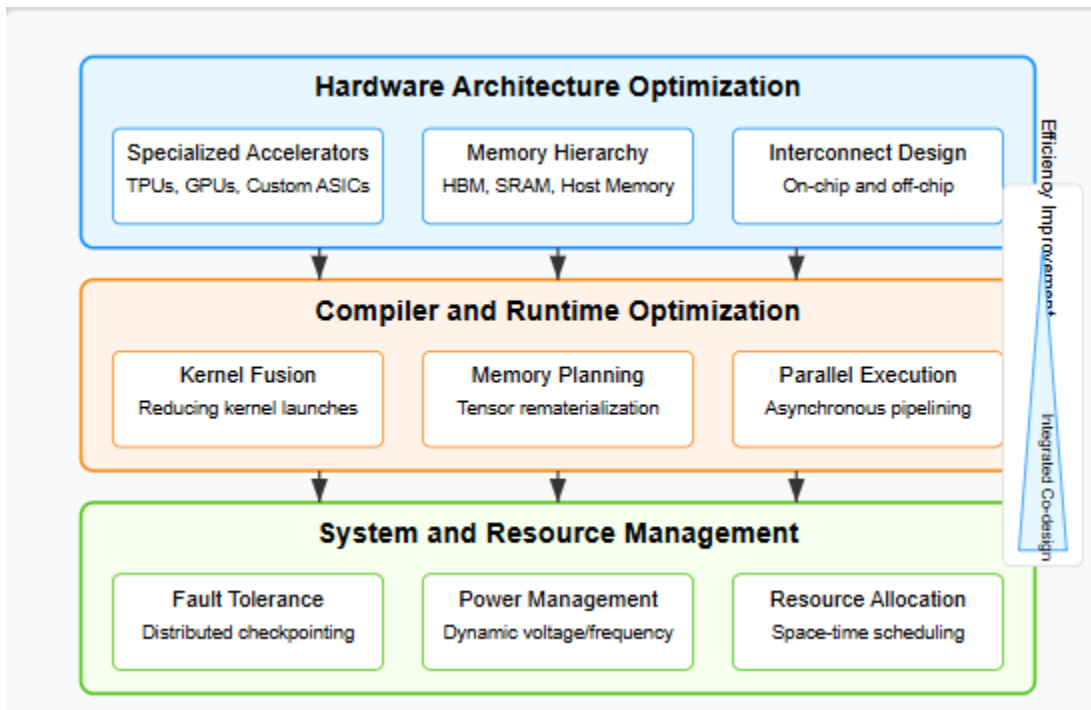


Fig. 3: Hardware-Software Integration: The Key to Efficient AI Infrastructure [9, 10]

## 6. Future Directions and Emerging Infrastructure Paradigms

### 6.1 Photonic-Electronic Integrated Circuits for AI Acceleration

The integration of photonics with electronic circuits represents a transformative approach to addressing the computational and interconnect challenges in AI infrastructure. Photonic-electronic integrated circuits (PEICs) offer fundamental advantages in bandwidth, latency, and energy efficiency compared to traditional electronic systems. Current state-of-the-art electronic interconnects in HPC systems operate at 112 Gb/s per lane with advanced PAM-4 modulation schemes, while photonic interconnects have demonstrated 224 Gb/s per wavelength with significantly lower power consumption of 5.5 pJ/bit versus 12 pJ/bit for electronic equivalents [11]. Silicon photonics technology has advanced rapidly, with recent foundry platforms supporting over 20 optical building blocks including modulators, photodetectors, and waveguides on a single chip. This integration enables wavelength division multiplexing (WDM) that can transmit multiple data streams simultaneously on different wavelengths, dramatically increasing effective bandwidth. The bandwidth-density advantage of photonics becomes particularly significant in AI applications, where an 8-channel WDM photonic link can replace 32 differential electrical lanes while consuming less power and occupying smaller area. Beyond simple interconnects, photonic matrix-vector multipliers have demonstrated computation speeds of 11 trillion multiply-accumulate operations per second—a critical operation in transformer networks—with potential energy efficiency improvements of two orders of magnitude compared to electronic implementations [11].

### 6.2 Domain-Specific Architecture and Systolic Arrays

The development of domain-specific architecture tailored to the computational patterns of large language models represents another critical direction in infrastructure evolution. Google's Tensor Processing Unit exemplifies this approach, utilizing a systolic

array architecture specifically optimized for the matrix multiplications that dominate neural network computation. The TPUv3 implements a 128×128 systolic array operating at 940 MHz, delivering significantly higher computational density than general-purpose processors. This specialized design achieves 123 TFLOPS of peak performance for bfloat16 matrix multiplications and 420 terabytes/second of internal memory bandwidth, enabling efficient execution of key transformer operations [12]. The TPU architecture incorporates a memory hierarchy specifically designed for neural network workloads, with high-bandwidth memory (HBM) providing 32 GB of capacity and 900 GB/s bandwidth per chip. This represents 3.4 times the bandwidth available in contemporary GPU architectures, significantly reducing memory access bottlenecks. When deployed in scaled configurations, TPU architectures utilize a custom interconnect fabric that delivers 652 GB/s of chip-to-chip bandwidth, enabling efficient distribution of training across multiple accelerators. The TPUv3 Pod configuration connects 2,048 chips in a direct mesh topology, achieving 107 petaflops of aggregate computation capacity with measured scaling efficiency exceeding 80% for large transformer models [12].

### 6.3 Multi-Chip Modules and Heterogeneous Integration

The physical integration of multiple silicon dies within a single package offers significant advantages for AI accelerators by overcoming reticle size limitations and enabling heterogeneous technology integration. This approach allows combining different process technologies optimized for specific functions—logic, memory, and interconnect—in a single package with high-bandwidth chip-to-chip connections. Current implementations using silicon interposers demonstrate inter-chip bandwidth exceeding 2 TB/s with energy efficiency below 0.5 pJ/bit, representing an order of magnitude improvement over traditional interconnect technologies [11]. The TPU architecture leverages advanced packaging technologies to integrate multiple HBM stacks within a single package, enabling the high memory bandwidth that is critical for transformer model performance. Beyond current implementations, emerging 3D integration technologies promise even greater density and performance advantages. Wafer-scale integration approaches similar to those employed in Cerebras systems have demonstrated the ability to integrate 850,000 processing elements on a single wafer, eliminating traditional chip boundaries and the associated communication overheads. This integration approach offers particular advantages for attention mechanism computation, where global visibility across the entire sequence length improves both performance and energy efficiency [12].

## 7. Conclusion

The infrastructure revolution underpinning modern LLMs represents a remarkable convergence of innovations across hardware, software, and distributed systems engineering. The synergistic advancement of memory-efficient parallelism, high-throughput interconnects, and sophisticated model sharding techniques has been instrumental in overcoming the computational barriers that once constrained AI capabilities. These infrastructure innovations have not merely enabled larger models but have fundamentally transformed how we conceptualize AI system design, shifting toward integrated ecosystems where hardware and software co-evolve. Looking forward, the path to next-generation AI systems will likely require even more radical infrastructure reimagining, with promising directions in optical computing, specialized accelerators, and sustainable computing architectures. The lessons from scaling current LLMs suggest that continued progress will demand ever-closer collaboration between hardware engineers, systems researchers, and ML scientists to create the computational foundation for tomorrow's AI breakthroughs.

**Conflicts of Interest:** The authors declare no conflict of interest.
**Publisher's Note**: All claims expressed in this article are solely those of the authors and do not necessarily represent those of their affiliated organizations, or those of the publisher, the editors and the reviewers.

## References

[1] Aaditya S et al., (2024) The Llama 3 Herd of Models, Meta, 23 July 2024. [Online]. Available: https://liweinlp.com/wp-content/uploads/2024/07/meta.pdf

[2] Abdul D et al., (n.d) TrIMS: Transparent and Isolated Model Sharing for Low Latency Deep Learning Inference in Function as a Service Environments, learningsys.org. [Online]. Available: http://learningsys.org/nips18/assets/papers/65CameraReadySubmissionNIPS_system_for_ml_TRIMS%20(1).pdf

[3] Deepak N et al., (2021) Efficient Large-Scale Language Model Training on GPU Clusters Using Megatron-LM, arXiv:2104.04473v5, 23 Aug 2021. [Online]. Available: https://arxiv.org/pdf/2104.04473

[4] Dmitry L et al., (2020) GShard: Scaling Giant Models with Conditional Computation and Automatic Sharding, arXiv:2006.16668v1, 30 Jun 2020. [Online]. Available: https://arxiv.org/pdf/2006.16668

[5] Mojahidul S M A et al., (n.d) Hardware Accelerators for Artificial Intelligence, arxiv.org/pdf/2411.13717. [Online]. Available: https://arxiv.org/pdf/2411.13717

[6] Norman P. J et al., (2021) Ten Lessons From Three Generations Shaped Google's TPUv4i, 2021 ACM/IEEE 48th Annual International Symposium on Computer Architecture (ISCA), 2021. [Online]. Available: https://gwern.net/doc/ai/scaling/hardware/2021-jouppi.pdf

[7] Norman. P. J et al., (2020) A Domain-Specific Supercomputer for Training Deep Neural Networks, University of Viscosin, 27 Oct. 2020. [Online]. Available: https://pages.cs.wisc.edu/~markhill/seminar2020/jouppi2020_10_tpu-v2-v3

[8]     Pitch P and Xin Y (n.d) Bandwidth Optimal All-reduce Algorithms for Clusters of Workstations," Florida State University.  [Online]. Available: https://www.cs.fsu.edu/~xyuan/paper/09jpdc.pdf

[9]     Samyam R et al., (2020) ZeRO: Memory Optimizations Toward Training Trillion Parameter Models, arXiv:1910.02054v3, 13 May 2020. [Online]. Available: https://arxiv.org/pdf/1910.02054

[10]    Samyam R et al., (2021) ZeRO-Infinity: Breaking the GPU Memory Wall for Extreme Scale Deep Learning, arXiv:2104.07857v1, 16 Apr 2021. [Online]. Available: https://arxiv.org/pdf/2104.07857

[11]    Shupeng N et al., (2024) Photonic-Electronic Integrated Circuits for High-Performance Computing and AI Accelerators, IEEE Xplore, 2024. [Online]. Available: https://sites.utexas.edu/chen-server/files/2024/08/Photonic-Electronic_Integrated_Circuits_for_High-Performance_Computing_and_AI_Accelerators.pdf

[12]    Tom B. B et al., (2020) Language Models are Few-Shot Learners, 34th Conference on Neural Information Processing Systems (NeurIPS 2020), 2020. [Online]. Available: https://proceedings.neurips.cc/paper/2020/file/1457c0d6bfcb4967418bfb8ac142f64a-Paper.pdf