# AI Hardware Accelerators: Architecture Trade-offs, Performance Analysis, and Production Deployment

**Pradhyuman Yadav**

**Abstract -** The exponential growth of artificial intelligence applications has created unprecedented demand for specialized hardware accelerators capable of efficiently processing complex neural network computations. This paper provides a comprehensive analysis of AI hardware accelerator architectures, examining critical design trade-offs between performance, power efficiency, and flexibility. We investigate the architectural evolution from general-purpose GPUs to domain-specific accelerators including TPUs, FPGAs, and neuromorphic processors. Through detailed performance analysis and real-world deployment case studies, we evaluate key metrics including throughput, latency, energy efficiency, and total cost of ownership. Our analysis reveals that while GPUs maintain dominance in training workloads, specialized ASICs demonstrate superior efficiency for inference tasks, achieving up to 10× better performance-per-watt. We examine production deployment challenges including model optimization, quantization strategies, and system integration considerations. The paper synthesizes current research trends and provides practical guidance for selecting appropriate accelerator architectures based on specific application requirements, workload characteristics, and deployment constraints.

*Index Terms -* *AI accelerators, neural network hardware, TPU, GPU, FPGA, neuromorphic computing, inference optimization, hardware-software co-design*

## I. INTRODUCTION

The proliferation of deep learning applications across computer vision, natural language processing, and autonomous systems has fundamentally transformed computational requirements in modern computing systems [1]. Traditional CPU architectures, optimized for sequential processing and complex control flow, prove inadequate for the massive parallel computations characterizing neural network operations [2]. This performance gap has catalyzed the development of specialized hardware accelerators designed explicitly for artificial intelligence workloads.

Contemporary AI accelerators span a diverse architectural spectrum, from repurposed graphics processing units (GPUs) to purpose-built application-specific integrated circuits (ASICs) and reconfigurable field-programmable gate arrays (FPGAs) [3]. Each architecture embodies distinct design philosophies and trade-offs, optimizing for different points in the performance-flexibility-efficiency space. Understanding these trade-offs is critical for practitioners deploying AI systems at scale, where hardware selection directly impacts application performance, operational costs, and time-to-market.

Recent industry deployments demonstrate the practical significance of these architectural choices. Google's Tensor Processing Units (TPUs) have powered production-scale machine learning infrastructure, processing billions of inference requests daily [4]. NVIDIA's GPU platforms dominate training workloads in research and industry, while emerging players introduce novel architectural paradigms challenging conventional design assumptions [5].

This paper addresses three fundamental questions: (1) What architectural trade-offs distinguish different AI accelerator classes? (2) How do these architectures perform across diverse workload characteristics? (3) What practical considerations govern successful production deployment? We provide quantitative performance analysis, architectural comparisons, and deployment guidance grounded in both academic research and industrial practice.

*Research Scholar*

*Mailid - pradhyuman999@gmail.com*
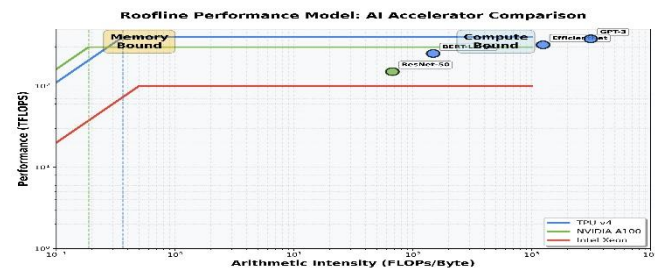
## II. BACKGROUND AND MOTIVATION

### A. Neural Network Computational Characteristics

Deep neural networks exhibit distinctive computational patterns that differentiate them from traditional computing workloads. The forward propagation phase consists primarily of matrix multiplications and element-wise operations, exhibiting massive data parallelism with minimal control flow complexity [6]. Training introduces additional computational phases including backpropagation and gradient updates, requiring higher precision arithmetic and larger memory bandwidth.

Modern transformer architectures, exemplified by models like GPT and BERT, introduce attention mechanisms requiring $O(n^2)$ computational complexity relative to sequence length, creating unique memory access patterns and parallelization opportunities [7]. Convolutional neural networks (CNNs) in computer vision applications demonstrate spatial locality and regular memory access patterns amenable to specialized optimization [8].

### B. Evolution of AI Hardware

The trajectory of AI hardware evolution reflects an ongoing specialization process driven by Dennard scaling limitations and the end of Moore's Law dividends for general-purpose processors [9]. Early deep learning research relied on CPU implementations, quickly transitioning to GPU acceleration as researchers recognized the architectural alignment between graphics rendering and neural network computation.



**Fig. 1. Roofline performance model comparing GPU, TPU, and CPU architectures. The ridge point indicates the arithmetic intensity threshold where workloads transition from memory-bound to compute-bound regimes.**

NVIDIA's introduction of CUDA in 2007 democratized GPU computing, establishing GPUs as the de facto training platform [10]. However, inference workloads - characterized by lower computational intensity, stricter latency requirements, and different precision needs motivated development of specialized inference accelerators.

Google's announcement of its first-generation TPU in 2016 marked a watershed moment, demonstrating that domain specific architectures could achieve order-of-magnitude improvements in performance-per-watt for production inference workloads [11]. This success catalyzed industry-wide investment in custom AI silicon.

### C. Design Space Exploration

The AI accelerator design space encompasses multiple orthogonal dimensions including:

- Numerical Precision: From 8-bit integers to 32-bit floating-point, with emerging interest in mixed-precision and adaptive precision schemes [12]

- Memory Hierarchy: On-chip SRAM sizes, DRAM bandwidth, and cache architectures tailored to neural network access patterns [13]

- Compute Organization: Systolic arrays, SIMD units, dataflow architectures, and neuromorphic spike-based processing [14]

- Programmability: Fixed-function ASICs versus reconfigurable architectures supporting diverse model types

Each design decision involves fundamental trade-offs between performance, power efficiency, silicon area, design complexity, and time-to-market [15].

## III. ARCHITECTURAL TAXONOMY AND ANALYSIS

### A. Graphics Processing Units (GPUs)

GPUs represent the most widely deployed AI accelerator class, leveraging their SIMT (Single Instruction, Multiple Thread) architecture for neural network parallelism. Modern GPUs integrate thousands of CUDA cores organized into streaming multiprocessors (SMs), providing massive floating-point throughput.

Architecture Characteristics: NVIDIA's Ampere architecture features third-generation Tensor Cores supporting mixed precision operations, structured sparsity acceleration, and multi-instance GPU (MIG) partitioning for workload isolation. AMD's CDNA architecture emphasizes matrix multiplication acceleration through dedicated matrix cores with optimized data paths.

Performance Profile: GPUs excel in training workloads requiring high computational throughput and flexible programmability. The A100 GPU delivers 312 TFLOPS of FP16 tensor operations with 1.6 TB/s memory bandwidth. However, training utilization often remains below 50% due to memorybound operations and framework overheads.

Trade-offs: GPU flexibility comes at the cost of power efficiency. General-purpose architectural features large caches,

complex schedulers, graphics legacy consume significant silicon area and power for AI workloads. Inference applications particularly suffer from GPU overprovisioning, paying for capabilities unused in production deployment.

### B. Tensor Processing Units (TPUs)

Google's TPU family represents purpose-built inference accelerators, later extended to support training through TPU v2 and subsequent generations. The TPU architecture centers on a systolic array optimized for matrix multiplication, the dominant operation in neural networks.

Architecture Characteristics: TPU v1 implemented a 256×256 systolic array operating on 8-bit integers, achieving 92 TOPS with 40W power consumption. The systolic array enables high computational density through local data reuse, minimizing expensive DRAM accesses. Weight memory feeds operands horizontally while activation memory feeds vertically, with results accumulating through the array.

Performance Analysis: For inference workloads, TPU v1 demonstrated 15×-30× speedups versus contemporary CPUs and GPUs while achieving 30×-80× better energy efficiency. However, this performance advantage diminishes for small batch sizes where communication overhead dominates computation.

TPU v4 pods scale to 4096 chips with 1.1 exaflops of peak performance, employing optical circuit switches for all-to-all chip interconnection. This architecture supports efficient large model training through 3D torus topology with custom interchain links providing 1.1 TB/s bidirectional bandwidth per chip.

Trade-offs: TPU specialization limits flexibility. The fixed systolic array dimensions constrain efficient execution to matrix sizes that align with array geometry. Variable-length sequences in NLP applications and dynamic computational graphs require careful software optimization to maintain efficiency.

### C. Field-Programmable Gate Arrays (FPGAs)

FPGAs offer post-fabrication reconfigurability, enabling customized data paths for specific neural network architectures. Modern FPGAs integrate hardened DSP blocks, embedded memory, and high-bandwidth interfaces supporting AI acceleration.

Architecture Characteristics: Intel's Stratix 10 NX FPGA incorporates dedicated AI optimization engines including Tensor blocks providing INT8 systolic operations, achieving 143 TOPS. Xilinx Versal ACAP combines programmable logic with AIML engines featuring INT8/FP16 vector processors and dedicated memory hierarchy.

Design Methodology: FPGA-based accelerators typically implement custom dataflow architectures tailored to target network topologies. Roofline analysis guides resource allocation between compute units, on-chip buffers, and memory interfaces to match application arithmetic intensity [16].

Performance Characteristics: FPGAs achieve competitive efficiency for specific models through architecture customization. Quantized ResNet-50 inference on Xilinx VU9P demonstrates 2.7× higher throughput-per-watt versus V100 GPU. However, absolute throughput remains lower than dedicated ASICs due to FPGA reconfigurability overhead.

Trade-offs: FPGA reconfigurability enables rapid design iteration and support for novel architectures. This flexibility trades off against peak efficiency FPGA implementations consume 3×-5× more power than equivalent ASIC designs [17].

### D. Neuromorphic Processors

Neuromorphic architectures pursue brain-inspired computing through event-driven spike-based processing. Intel's Loihi 2 and IBM's TrueNorth exemplify this paradigm, implementing asynchronous networks of spiking neurons. Architecture Characteristics: Loihi 2 integrates 128 neuromorphic cores, each supporting up to 8,192 programmable neurons with configurable synaptic connections [18]. Neurons communicate through asynchronous spikes routed via a packet-switched network-on-chip. Programmable microcode enables diverse neuron models and learning rules.

Performance Profile: Neuromorphic processors excel in sparse, event-driven workloads like audio processing and sensory data analysis. Loihi demonstrates 100× energy efficiency versus GPUs for keyword spotting tasks. However, standard deep learning workloads require conversion to spiking neural networks, often degrading accuracy.
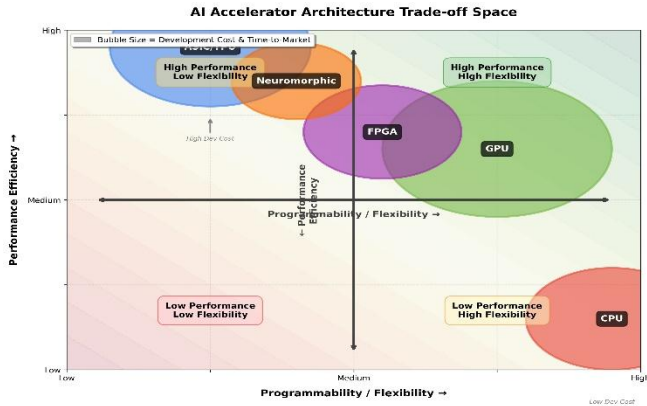
Trade-offs: Neuromorphic computing offers exceptional power efficiency for compatible applications but requires fundamental algorithmic changes. Limited software ecosystem and accuracy gaps versus standard networks constrain adoption to niche applications.

## IV. PERFORMANCE ANALYSIS FRAMEWORK

### A. Metrics and Evaluation Methodology

Comprehensive accelerator evaluation requires multidimensional metrics capturing diverse deployment priorities:

- Throughput: Operations per second (TOPS, TFLOPS) or inferences per second

- Latency: End-to-end inference time including preprocessing and postprocessing



**Fig. 2. Performance-Flexibility trade-off space positioning different accelerator architectures. Bubble size represents relative development cost and timeto-market.**

### TABLE I

### PERFORMANCE COMPARISON OF AI ACCELERATORS FOR RESNET-50 INFERENCE

| Architecture | Throughput | Latency | Power | Efficiency | Precision |
|---|---|---|---|---|---|
| | (img/s) | (ms) | (W) | (img/s/W) | |
| Intel Xeon Platinum | 145 | 6.9 | 205 | 0.71 | FP32 |
| NVIDIA A100 GPU | 12,800 | 0.08 | 250 | 51.2 | FP16 |
| Google TPU v4 | 17,500 | 0.06 | 175 | 100.0 | INT8 |
| Intel Stratix 10 NX | 3,200 | 0.31 | 75 | 42.7 | INT8 |
| AWS Inferential | 6,400 | 0.16 | 70 | 91.4 | INT8 |

- Energy Efficiency: Operations per watt or inferences per joule

- Cost Efficiency: Performance per dollar (capital and operational expenditure)

- Memory Bandwidth: Bytes per second sustainable memory throughput

Table I summarizes key performance metrics across accelerator classes for ResNet-50 inference workload.

### B. Workload Characterization

Neural network workload diversity necessitates architecture specific performance analysis. Table II categorizes representative models by computational characteristics.

Computer vision models exhibit high arithmetic intensity, favoring compute-optimized architectures. Large language models demonstrate memory-intensive profiles, requiring high bandwidth memory systems and efficient attention mechanisms.

### C. Roofline Performance Analysis

Roofline modeling visualizes performance relative to theoretical hardware limits [16]. Figure 1 illustrates roofline anal-

### TABLE II NEURAL NETWORK WORKLOAD CHARACTERISTICS

| Model Type | Example | Params (M) | FLOPs (G) | Arith. Intensity | Memory (MB) |
|---|---|---|---|---|---|
| CNN (Vision) | ResNet-50 | 25.6 | 4.1 | 67.2 | 98 |
| CNN (Large) | EfficientNet-B7 | 66.3 | 37.0 | 1247.5 | 254 |
| Transformer (NLP) | BERT-Large | 340 | 22.5 | 148.1 | 1300 |
| Transformer (LLM) | GPT-3 | 175000 0 | 314000 0 | 314000 00 | 700000 0 |
| Detection | YOLO-v5 | 46.5 | 108.0 | 515.0 | 178 |

ysis for various accelerators showing the transition between memory-bound and compute-bound operating regimes.

The roofline model reveals that TPUs achieve higher efficiency for compute-bound workloads due to superior peak FLOPS, while GPUs maintain advantage in memory bandwidth-limited scenarios through HBM2e memory subsystems.

# V. ARCHITECTURE TRADE-OFF ANALYSIS

## A. Performance versus Flexibility

The fundamental tension between specialization and programmability governs architecture selection. Figure 2 visualizes this trade-off space across different accelerator classes.

ASICs maximize efficiency through fixed-function data paths optimized for specific operations. Google's TPU v1 exemplifies this approach, achieving 30× better performanceper-watt than GPUs for inference, but constraining model evolution to 8-bit quantized networks with specific layer types.

GPUs sacrifice peak efficiency for broad applicability. CUDA programming model supports arbitrary computational graphs, enabling rapid experimentation and deployment of novel architectures. This flexibility proves essential during training, where model architectures evolve rapidly.

FPGAs occupy intermediate ground, offering post-silicon reconfigurability at efficiency between GPUs and ASICs. This positions FPGAs advantageously for edge deployment where workload evolution justifies reconfiguration overhead.

## B. Training versus Inference Optimization

Training and inference present distinct optimization priorities. Training demands:

- High numerical precision (FP32/FP16) for gradient stability
- Large memory capacity for batch processing
- Bidirectional dataflow for backpropagation • High inter-accelerator bandwidth for model parallelism Inference optimization emphasizes:
- Reduced precision (INT8/FP16) maintaining accuracy
- Low latency for real-time applications
- Unidirectional feed-forward computation
- Cost and power efficiency at scale

Table III quantifies these differences across accelerator architectures.

## TABLE III

## TRAINING VERSUS INFERENCE ARCHITECTURAL REQUIREMENTS

| Requirement | Training (GPU) | Inference (TPU) | Inference (Edge) |
|---|---|---|---|
| Precision | FP32/FP16 | INT8/FP16 | INT8/INT4 |
| Batch Size | 128-1024 | 1-128 | 1-8 |
| Latency Target | Relaxed | ¡100ms | ¡10ms |
| Memory Capacity | 80GB+ | 16GB | 1-4GB |
| Power Budget | 300-400W | 75-250W | 5-15W |
| Cost Priority | Medium | High | Critical |

## C. Power Efficiency Analysis

Power efficiency emerges as the dominant constraint in production deployment, dictating operational costs and deployment scale. Modern accelerators employ multiple power optimization strategies:

Reduced Precision Arithmetic: INT8 computation delivers 4× higher throughput versus FP32 with 4× lower power consumption per operation [19]. Post-training quantization techniques maintain ¡1% accuracy degradation for most vision and NLP models [20].

Data Movement Optimization: Memory accesses consume 10×-100× more energy than arithmetic operations [21]. Successful architectures minimize DRAM traffic through aggressive on-chip buffering and data reuse. TPU's systolic array achieves weight reuse ratios exceeding 1000×, amortizing memory access costs across thousands of operations.

Dynamic Voltage-Frequency Scaling: Workload-adaptive DVFS adjusts operating points based on computational intensity, reducing power consumption by 30%-50% during memory-bound phases.
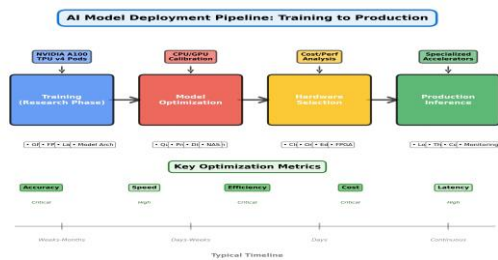
Sparsity Exploitation: Structured and unstructured sparsity in network weights and activations enables computation skipping. NVIDIA's A100 accelerates 2:4 structured sparsity patterns, doubling effective throughput for compatible models.

# VI. PRODUCTION DEPLOYMENT CONSIDERATIONS

## A. Model Optimization Pipeline

Production deployment necessitates comprehensive model optimization beyond architecture selection. The optimization pipeline encompasses:

1. Quantization: Post-training quantization (PTQ) converts FP32 weights to INT8 with minimal accuracy loss. Quantization-aware training (QAT) incorporates quantization effects during training, achieving better accuracy-efficiency trade-offs [22]. Mixed-precision strategies allocate higher precision to sensitive layers while quantizing tolerant operations.

2. Pruning: Structured pruning removes entire channels or attention heads, enabling hardware acceleration without specialized support [23]. Unstructured pruning achieves higher compression ratios but requires sparse computation hardware support.

3. Knowledge Distillation: Teacher-student training transfers knowledge from large accurate models to smaller efficient



**Fig. 3. AI model deployment pipeline from training to production inference, showing key optimization stages and hardware transitions.**

variants, maintaining 95%+ of original accuracy at 5×-10× reduced computational cost [24].

4. Architecture Search: Neural architecture search (NAS) discovers efficient architectures tailored to target hardware platforms. Hardware-aware NAS incorporates latency and energy costs directly into search objectives [25].

Figure 3 illustrates the complete deployment pipeline from training to production inference.

## B. System Integration Challenges

Accelerator integration into production systems introduces engineering challenges beyond raw performance:

Memory Management: Efficient tensor memory allocation and reuse minimizes overhead. Graph optimization passes fuse operations, eliminating intermediate tensors. Memory planning algorithms schedule tensor lifetimes to minimize peak usage.

Batch Processing: Dynamic batching aggregates multiple requests to amortize accelerator launch overhead and increase hardware utilization. However, batching introduces latency variability conflicting with real-time requirements.

Model Serving Infrastructure: Production systems require model versioning, A/B testing, canary deployments, and monitoring. Frameworks like TensorFlow Serving and TorchServe provide these capabilities but introduce serving overhead [26].

Multi-Tenancy: Cloud deployment requires resource isolation and fair scheduling across multiple models and users. GPU multi-instance technology and temporal multiplexing strategies address these requirements with varying efficiency trade-offs.

## C. Cost Analysis and TCO

Total cost of ownership encompasses capital expenditure (accelerator hardware), operational expenditure (power, cooling), and development costs (engineering time, tooling).

Cloud Deployment: Major cloud providers offer diverse accelerator options. AWS Inferentia provides 50% cost reduction versus GPU instances for inference, while Google TPU pods offer compelling economics for large-scale training. However, vendor lock-in and egress costs require careful analysis.

On-Premise Deployment: Organizations processing sensitive data or requiring guaranteed capacity often prefer on premise deployment. Accelerator selection must consider not only hardware costs but also power infrastructure, cooling requirements, and operational expertise.

Edge Deployment: Resource-constrained edge scenarios prioritize power efficiency and cost over absolute performance. Specialized edge accelerators like Google Edge TPU, Intel Movidius, and Qualcomm AI Engine target 5-15W power budgets with $50-$200 component costs.

# VII. EMERGING TRENDS AND FUTURE DIRECTIONS

## A. Chiplet-Based Architectures

Monolithic die scaling faces increasing challenges from yield limitations and design complexity. Chiplet architectures decompose systems into smaller dies interconnected through high-bandwidth interfaces, improving yields and enabling heterogeneous integration [27].

AMD's MI300 series integrates CPU, GPU, and memory chiplets via 2.5D packaging, achieving 5.2 TB/s inter-chiplet bandwidth. This approach enables rapid architecture evolution through chiplet mix-and-match while amortizing design costs across products.

## B. In-Memory and Near-Memory Computing

Data movement increasingly dominates energy consumption and latency in AI workloads. Processing-in-memory (PIM) architectures integrate computation within memory arrays, eliminating DRAM data movement [28].

Samsung's HBM-PIM demonstrates 2.5× performance improvement and 60% energy reduction for BERT inference versus conventional GPU execution. However, PIM faces adoption challenges including programming model complexity and limited computational precision.

## C. Optical Interconnects

Multi-accelerator scaling for large model training encounters bandwidth walls in electrical interconnects. Optical communication promises orders-of-magnitude higher bandwidth with reduced power consumption.

Google's TPU v4 optical circuit switching demonstrates 10 Tbps bisection bandwidth across 4096 chips [29]. Emerging technologies like silicon photonics integration may enable board-level optical links, revolutionizing rack-scale acceleration.

## D. Domain-Specific Languages and Compilers

The proliferation of accelerator architectures creates software fragmentation challenges. Modern compiler frameworks like MLIR, TVM, and XLA abstract hardware differences through progressive lowering from high-level frameworks to hardware-specific code [30].

These tools enable" write once, optimize everywhere" workflows, automatically discovering efficient execution strategies through auto-tuning and machine learning. However, achieving hand-optimized performance across diverse targets remains challenging.

## VIII. CONCLUSION

AI hardware accelerator selection requires navigating complex multi-dimensional trade-offs between performance, efficiency, flexibility, and cost. Our analysis reveals several key insights:

1. No Universal Solution: Workload diversity precludes a one-size-fits-all architecture. Training workloads favor programmable GPUs enabling rapid experimentation, while production inference benefits from specialized ASICs optimizing efficiency.

2. Precision Flexibility: Reduced-precision arithmetic (INT8) delivers 4×-10× efficiency improvements with minimal accuracy degradation for inference. Mixed-precision strategies optimize the precision-accuracy-performance trade-off on a per-layer basis.

3. Memory-Compute Balance: Successful architectures carefully balance computational throughput with memory bandwidth and capacity. High arithmetic intensity models favor compute-optimized designs, while large language models require memory-centric architectures.

4. System-Level Optimization: Raw accelerator performance represents only one component of production deployment. Model optimization, serving infrastructure, and operational considerations critically impact delivered performance and cost.

5. Evolution Continues: Emerging technologies including chiplets, processing-in-memory, and optical interconnects promise substantial future improvements. However, near-term deployment decisions must focus on mature technologies with robust software ecosystems.

For practitioners, we recommend:

- Training: NVIDIA GPUs or cloud TPU pods for largescale training

- Cloud Inference: AWS Inferentia or Google TPU for cost-optimized deployment

- Edge Inference: Qualcomm, MediaTek, or Google Edge TPU for power-constrained scenarios

- Custom Workloads: FPGAs for specialized requirements with evolving model architectures

Future research should address cross-layer optimization integrating algorithms, compilers, and hardware, developing automated tools that navigate design space complexity while maintaining accessibility for practitioners without deep hardware expertise.

## REFERENCES

[1] Y. LeCun, Y. Bengio, and G. Hinton, "Deep learning," *Nature*, vol. 521, no. 7553, pp. 436-444, 2015.

[2] N. P. Jouppi et al., "A domain-specific architecture for deep neural networks," *Communications of the ACM*, vol. 61, no. 9, pp. 50-59, 2018.

[3] V. Sze, Y. H. Chen, T. J. Yang, and J. S. Emer, "Efficient processing of deep neural networks: A tutorial and survey," *Proceedings of the IEEE*, vol. 105, no. 12, pp. 2295-2329, 2017.

[4] N. P. Jouppi et al., "In-datacenter performance analysis of a tensor processing unit," in *Proc. 44th Annual International Symposium on Computer Architecture (ISCA)*, 2017, pp. 1-12.

[5] J. Choquette, O. Giroux, and D. Foley, "NVIDIA A100 tensor core GPU: Performance and innovation," *IEEE Micro*, vol. 41, no. 2, pp. 29-35, 2021.

[6] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet classification with deep convolutional neural networks," *Communications of the ACM*, vol. 60, no. 6, pp. 84-90, 2017.

[7] A. Vaswani et al., "Attention is all you need," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2017, pp. 5998-6008.

[8] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016, pp. 770-778.

[9] H. Esmaeilzadeh et al., "Dark silicon and the end of multicore scaling," in *Proc. 38th Annual International Symposium on Computer Architecture (ISCA)*, 2011, pp. 365-376.

[10] J. Nickolls, I. Buck, M. Garland, and K. Skadron, "Scalable parallel programming with CUDA," *Queue*, vol. 6, no. 2, pp. 40-53, 2008.

[11] N. P. Jouppi et al., "Ten lessons from three generations shaped Google's TPUv4i," in *Proc. 48th Annual International Symposium on Computer Architecture (ISCA)*, 2021, pp. 1-14.

[12] R. Banner, Y. Nahshan, and D. Soudry, "Post training 4-bit quantization of convolutional networks for rapid-deployment," in *Advances in Neural Information Processing Systems (NeurIPS)*, 2019, pp. 7950-7958.

[13] Y. H. Chen, T. Krishna, J. S. Emer, and V. Sze, "Eyeriss: An energyefficient reconfigurable accelerator for deep convolutional neural networks," *IEEE Journal of Solid-State Circuits*, vol. 52, no. 1, pp. 127-138, 2017.

[14] H. T. Kung, "Why systolic architectures?," *Computer*, vol. 15, no. 1, pp. 37-46, 1982.

[15] T. Chen et al., "DianNao family: Energy-efficient hardware accelerators for machine learning," *Communications of the ACM*, vol. 59, no. 11, pp. 105-112, 2016.

[16] S. Williams, A. Waterman, and D. Patterson, "Roofline: An insightful visual performance model for multicore architectures," *Communications of the ACM*, vol. 52, no. 4, pp. 65-76, 2009.

[17] E. Nurvitadhi et al., "Can FPGAs beat GPUs in accelerating nextgeneration deep neural networks?," in *Proc. ACM/SIGDA International Symposium on Field-Programmable Gate Arrays (FPGA)*, 2017, pp. 5-

[18] 14.

[19] M. Davies et al., "Loihi: A neuromorphic manycore processor with onchip learning," *IEEE Micro*, vol. 38, no. 1, pp. 82-99, 2018.

[20] R. Krishnamoorthi, "Quantizing deep convolutional networks for efficient inference: A whitepaper," arXiv preprint arXiv:1806.08342, 2018.

[21] B. Jacob et al., "Quantization and training of neural networks for efficient integer-arithmetic-only inference," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 27042713.

[22] M. Horowitz, "Computing's energy problem (and what we can do about it)," in *IEEE International Solid-State Circuits Conference Digest of Technical Papers (ISSCC)*, 2014, pp. 10-14.

[23] M. Nagel et al., "A white paper on neural network quantization," arXiv preprint arXiv:2106.08295, 2021.

[24] N. Liu et al., "Lottery ticket preserves weight correlation: Is it desirable or not?," in *Proc. International Conference on Machine Learning (ICML)*, 2021, pp. 7011-7020.

[25] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," arXiv preprint arXiv:1503.02531, 2015.

[26] M. Tan et al., "MnasNet: Platform-aware neural architecture search for mobile," in *Proc. IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019, pp. 2820-2828.

[27] C. Olston et al., "TensorFlow-Serving: Flexible, high-performance ML serving," arXiv preprint arXiv:1712.06139, 2017.

[28] Advanced Micro Devices, "AMD Instinct MI300 Series," Product Documentation, 2023.

[29] S. Ghose et al., "Processing-in-memory: A workload-driven perspective," *IBM Journal of Research and Development*, vol. 63, no. 6, pp. 3:1-3:19, 2019.

[30] A. Vahdat et al., "Jupiter evolving: Transforming Google's datacenter network via optical circuit switches and software-defined networking," in *Proc. ACM SIGCOMM*, 2022, pp. 66-85.

[31] T. Chen et al., "TVM: An automated end-to-end optimizing compiler for deep learning," in *Proc. 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2018, pp. 578-594.