

A Survey of Machine Learning Approaches for Computer Architecture Performance Modeling

Authors TBD

Affiliations TBD

Abstract—Machine learning-based performance modeling has emerged as a powerful alternative to traditional analytical models and cycle-accurate simulators for predicting computer system behavior. This survey comprehensively analyzes ML techniques for performance prediction across CPUs, GPUs, accelerators, and distributed systems, covering over 60 papers from architecture and ML venues published between 2016–2025. We propose an eight-dimension taxonomy organizing approaches by modeling technique, target hardware, workload types, prediction targets, accuracy metrics, input requirements, evaluation scope, and reproducibility. Our analysis reveals that specialized models achieve remarkable accuracy—below 5% error for narrow domains—while general-purpose models trade accuracy for broader applicability. Transfer learning and meta-learning techniques increasingly enable adaptation to new hardware with minimal profiling, addressing the challenge of hardware diversity. We identify key open challenges including benchmark diversity, cross-platform generalization, and integration with compiler and architecture exploration workflows. Hybrid approaches combining analytical structure with learned components represent a promising direction, offering both interpretability and accuracy. This survey provides practitioners guidance for selecting appropriate techniques and researchers a roadmap for advancing the field.

Index Terms—machine learning, performance modeling, computer architecture, neural networks, survey

I. INTRODUCTION

Performance modeling is fundamental to computer architecture research and development. Architects rely on accurate performance predictions to navigate vast design spaces, optimize hardware-software co-design, and make informed decisions about resource allocation. Traditional approaches—analytical models [1] and cycle-accurate simulators [2]—have served the community well, but face growing challenges as workloads and hardware become increasingly complex. Analytical models often oversimplify system behavior, while simulators can require hours or days to evaluate a single design point, making exhaustive exploration impractical.

The rise of deep learning workloads has intensified these challenges. Modern neural networks exhibit diverse computational patterns—from dense matrix operations in transformers to sparse irregular accesses in graph neural networks—that stress traditional modeling assumptions. Simultaneously, hardware diversity has exploded: GPUs, TPUs, custom accelerators, and multi-device distributed systems each present unique performance characteristics that resist unified analytical treatment. This complexity has motivated a new generation of *machine learning-based* performance models that learn predictive functions directly from profiling data.

ML-based performance modeling has emerged as a compelling alternative. Learned models can capture complex, non-linear relationships between workload characteristics and hardware behavior that elude closed-form analysis. Recent work demonstrates remarkable accuracy: NeuSight [3] achieves 2.3% error predicting GPT-3 latency on H100 GPUs, while nn-Meter [4] reaches 99% accuracy for edge device latency prediction. Beyond accuracy, these approaches offer practical benefits: models trained on one platform can transfer to new hardware with minimal adaptation [5], and inference-time predictions complete in milliseconds rather than hours.

This survey provides a comprehensive analysis of ML-based performance modeling techniques for computer architecture. We make the following contributions:

- A **taxonomy** organizing approaches along eight dimensions: modeling technique, target hardware, workload types, prediction targets, accuracy metrics, input requirements, evaluation scope, and reproducibility.
- A **systematic survey** of over 60 papers from architecture venues (MICRO, ISCA, HPCA, ASPLOS) and ML venues (MLSys, NeurIPS, ICML) published between 2016–2025.
- A **comparative analysis** examining trade-offs between accuracy, training cost, generalization, and interpretability across approaches.
- An identification of **open challenges** including data scarcity, cross-platform generalization, and integration with design automation flows.

The remainder of this paper is organized as follows. Section II provides background on traditional performance modeling and relevant ML techniques. Section III presents our classification taxonomy. Section IV surveys approaches organized by target hardware platform. Section V offers comparative analysis across key dimensions. Section VI discusses open challenges and future directions. Section VII presents hands-on reproducibility evaluations of representative tools. Section VIII concludes.

II. BACKGROUND

A. Traditional Performance Modeling

Performance modeling has traditionally relied on two complementary approaches: analytical models and cycle-accurate simulation. This section reviews both paradigms and their limitations, motivating the emergence of ML-based alternatives.

1) *Analytical Models*: Analytical models express performance as closed-form functions of hardware and workload parameters. The roofline model [1] exemplifies this approach, bounding attainable performance by peak compute throughput and memory bandwidth. Given operational intensity I (FLOP/byte), the roofline predicts performance as $P = \min(\pi, \beta \cdot I)$, where π is peak FLOPS and β is memory bandwidth. Despite its simplicity, roofline reasoning guides optimization by revealing compute-bound versus memory-bound regimes.

For DNN accelerators, analytical cost models have become standard practice. Timeloop [6] models data movement across memory hierarchies for any given mapping (loop order and tiling), computing access counts and energy from architectural parameters. MAESTRO [7] provides a data-centric framework that derives performance from dataflow descriptions. Sparseloop [8] extends this methodology to sparse tensor operations, achieving 2000 \times speedup over RTL simulation while maintaining accuracy.

Analytical models offer several advantages: fast evaluation (microseconds per design point), interpretability (designers can trace predictions to specific terms), and extrapolation to unseen configurations. However, they require manual derivation for each target architecture, struggle to capture complex microarchitectural effects (contention, pipeline stalls, caching behavior), and may oversimplify non-linear interactions.

2) *Cycle-Accurate Simulation*: Cycle-accurate simulators model hardware at register-transfer level, faithfully reproducing timing behavior. General-purpose simulators like gem5 [2] support flexible configuration of CPU cores, caches, memory controllers, and interconnects. For GPUs, simulators such as GPGPU-Sim [9] and Accel-Sim [10] model SIMT execution, warp scheduling, and memory coalescing.

Cycle-accurate simulation achieves high fidelity—typically within 5–15% of real hardware [2]—and supports detailed microarchitectural studies. However, simulation speed presents a fundamental limitation: evaluating a single ResNet-50 inference may require hours, making design space exploration impractical. ASTRA-sim [11] addresses distributed training at scale through analytical abstractions, but even coarse-grained simulation struggles with the combinatorial explosion of modern ML workloads and hardware configurations.

3) *The Modeling Gap*: Neither approach fully addresses modern performance modeling needs. Analytical models are fast but imprecise for complex microarchitectures. Simulators are accurate but too slow for iterative design. This tension has intensified as ML workloads diversify (from CNNs to transformers to mixture-of-experts models) and hardware specializes (GPUs, TPUs, custom accelerators). ML-based performance models offer a middle path: learning complex relationships from profiling data while enabling millisecond-scale inference.

B. Machine Learning Fundamentals

This section provides a brief primer on ML techniques frequently employed in performance modeling, establishing

terminology used throughout the survey.

1) *Classical Machine Learning*: Linear regression and its regularized variants (ridge, LASSO) remain widely used for performance prediction due to their simplicity and interpretability. Given feature vector \mathbf{x} (e.g., operator parameters, hardware counters), linear models predict $\hat{y} = \mathbf{w}^\top \mathbf{x} + b$. While unable to capture non-linear relationships, linear models provide baselines and feature importance rankings.

Tree-based ensembles—random forests and gradient boosted trees (XGBoost, LightGBM)—handle non-linearities through recursive partitioning. These methods dominate when training data is limited (<10K samples) and features are well-engineered, often outperforming deep learning in low-data regimes [4].

2) *Deep Learning*: Multi-layer perceptrons (MLPs) learn hierarchical feature representations through stacked non-linear transformations: $\mathbf{h}_{i+1} = \sigma(\mathbf{W}_i \mathbf{h}_i + \mathbf{b}_i)$. MLPs require minimal feature engineering but need sufficient training data and careful regularization to avoid overfitting.

Recurrent neural networks (RNNs) and their gated variants (LSTM, GRU) process sequential inputs, making them suitable for modeling operator sequences in neural network execution graphs. However, sequential processing limits parallelization and can miss long-range dependencies.

3) *Graph Neural Networks*: Graph neural networks (GNNs) operate on graph-structured data through message passing. For a node v with features \mathbf{h}_v , GNNs iteratively update representations by aggregating information from neighbors $\mathcal{N}(v)$:

$$\mathbf{h}_v^{(k+1)} = \phi \left(\mathbf{h}_v^{(k)}, \bigoplus_{u \in \mathcal{N}(v)} \psi(\mathbf{h}_u^{(k)}, \mathbf{e}_{uv}) \right) \quad (1)$$

where ϕ and ψ are learnable functions and \bigoplus is a permutation-invariant aggregation (sum, mean, max).

GNNs are particularly appealing for performance modeling because DNN computation graphs have natural graph structure. Nodes represent operators with features (type, parameters), edges represent data dependencies with features (tensor shapes, datatypes). GNNs can learn to propagate performance-relevant information along these dependencies [12].

4) *Attention and Transformers*: Attention mechanisms compute weighted combinations over input elements, with weights determined by learned compatibility functions. Self-attention allows each position to attend to all other positions:

$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax} \left(\frac{\mathbf{Q}\mathbf{K}^\top}{\sqrt{d_k}} \right) \mathbf{V} \quad (2)$$

Transformers stack self-attention with feedforward networks, enabling long-range dependency modeling without sequential processing. Recent performance models leverage transformer architectures to capture complex inter-operator interactions across entire computation graphs.

5) *Transfer Learning*: Transfer learning adapts models trained on one domain (source) to perform well on another (target). In performance modeling, this enables training on

easily-profiled hardware and transferring to new platforms with limited data. Common approaches include fine-tuning (adjusting pre-trained weights with target data), domain adaptation (learning domain-invariant representations), and meta-learning (learning to adapt quickly from few examples) [5].

C. Problem Formulation

We now formally define the performance modeling problem and establish the evaluation framework used throughout this survey.

1) *Inputs and Outputs*: Performance modeling maps workload and hardware descriptions to performance metrics. Formally, given workload specification \mathcal{W} and hardware configuration \mathcal{H} , a performance model f predicts metric y :

$$\hat{y} = f(\mathcal{W}, \mathcal{H}; \theta) \quad (3)$$

where θ represents model parameters (weights for ML models, equations for analytical models).

Workload representations vary by granularity and abstraction:

- *Operator-level*: Individual layer parameters (kernel size, channels, batch size)
- *Graph-level*: Full computation graph with node and edge features
- *IR-level*: Intermediate representations from compilers (TVM [13], XLA)
- *Trace-level*: Execution traces capturing runtime behavior

Hardware representations similarly span multiple levels:

- *Specification*: Static parameters (core count, memory size, bandwidth)
- *Counter-based*: Runtime performance counters (cache misses, stalls)
- *Embedding*: Learned dense representations of hardware platforms

2) *Prediction Targets*: Performance models target various metrics depending on application requirements:

Latency measures execution time, typically end-to-end inference time or per-layer latency. Latency prediction is critical for real-time applications with strict deadlines and for optimizing user-facing services.

Throughput captures sustained processing rate: samples per second for inference, tokens per second for language models, or images per second for training. Throughput optimization maximizes hardware utilization for batch processing.

Energy encompasses power consumption (Watts) and energy per operation (Joules/inference). Energy prediction is essential for mobile deployment, data center cost optimization, and sustainability considerations.

Memory includes peak memory footprint (for feasibility checking), memory bandwidth utilization, and memory access patterns.

Multi-objective formulations jointly predict multiple metrics, enabling Pareto-optimal design selection balancing latency, energy, and accuracy.

3) *Accuracy Metrics*: The field employs several accuracy metrics, each with distinct interpretations:

Mean Absolute Percentage Error (MAPE) measures average relative deviation:

$$\text{MAPE} = \frac{100\%}{n} \sum_{i=1}^n \left| \frac{y_i - \hat{y}_i}{y_i} \right| \quad (4)$$

MAPE is scale-invariant and interpretable (5% MAPE means predictions typically differ by 5% from ground truth).

Root Mean Square Error (RMSE) penalizes large errors more heavily:

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2} \quad (5)$$

Correlation coefficients (Pearson, Spearman) measure how well predictions track relative ordering—important when models guide design space exploration.

Ranking accuracy directly evaluates whether models correctly order configurations, often measured via Kendall’s τ or top- k accuracy.

4) *Hardware Targets*: Modern performance modeling spans diverse hardware platforms:

CPUs remain important for general-purpose inference and training of smaller models. CPU modeling must account for complex cache hierarchies, branch prediction, out-of-order execution, and SIMD vectorization.

GPUs dominate ML training and large-scale inference. GPU modeling addresses SIMT execution, warp scheduling, memory coalescing, and multi-GPU scaling.

TPUs and custom accelerators employ specialized dataflows for matrix operations. Modeling these devices requires understanding systolic arrays, on-chip memory hierarchies, and dataflow mappings.

Edge devices (mobile SoCs, embedded NPUs) impose strict power and memory constraints. Edge modeling emphasizes latency under thermal throttling and memory-limited execution.

Distributed systems scale training across multiple devices and nodes. Distributed modeling must capture communication overhead, synchronization barriers, and pipeline parallelism.

This diversity of targets, workloads, and metrics motivates our comprehensive taxonomy in Section III.

III. TAXONOMY

We organize the surveyed literature along three primary dimensions: the hardware target being modeled, the machine learning techniques employed, and the input representations used. Figure 1 illustrates how these dimensions intersect to characterize different performance modeling approaches. This taxonomy extends existing classifications [6], [7] by incorporating the emerging diversity of ML-based methods and their distinctive design choices.

Our classification scheme serves two purposes. First, it provides a systematic framework for understanding the design space of ML-based performance models—researchers can

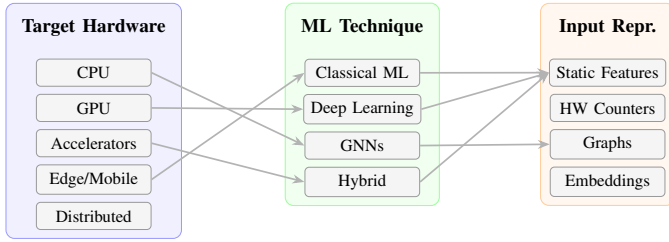


Fig. 1. Three-dimensional taxonomy for ML-based performance modeling. Arrows indicate common pairings observed in the literature (e.g., GPU models often use deep learning with static features).

identify which combinations of targets, techniques, and representations have been explored versus those that remain open. Second, it enables practitioners to select appropriate methods for their use cases by matching problem characteristics (target hardware, available data, accuracy requirements) to model capabilities.

A. By Modeling Target

The choice of hardware target fundamentally shapes model design, as different platforms exhibit distinct performance characteristics and modeling challenges.

1) *CPU Performance Modeling*: CPUs present complex modeling challenges due to deep out-of-order pipelines, sophisticated cache hierarchies, and branch prediction. ML models for CPU performance must capture instruction-level parallelism, cache behavior, and memory access patterns. Traditional approaches relied on microbenchmark-based linear regression [2], while recent work employs graph neural networks to model basic block throughput [12]. CPU modeling remains challenging due to the diversity of microarchitectures and the difficulty of capturing dynamic effects like branch misprediction and cache contention.

2) *GPU Performance Modeling*: GPUs dominate modern ML training and inference, making accurate GPU performance prediction critical. GPU modeling must account for SIMT execution, warp scheduling, memory coalescing, and memory bandwidth limitations. Early approaches used analytical roofline models [1], but these struggle with the complex memory hierarchies and occupancy effects of modern GPUs.

ML-based GPU models have achieved remarkable accuracy. NeuSight [3] introduces tile-based prediction that mirrors CUDA’s execution model, achieving 2.3% error on GPT-3 inference across H100, A100, and V100 GPUs. Habitat [14] pioneered runtime-based cross-GPU prediction using wave scaling analysis. These approaches demonstrate that learned models can capture GPU performance characteristics that elude analytical treatment.

3) *DNN Accelerator Modeling*: Custom DNN accelerators—including TPUs, NPUs, and systolic array designs—employ specialized dataflows optimized for matrix operations. Modeling these devices requires understanding the interaction between dataflow, memory hierarchy, and tensor tiling.

Analytical frameworks like Timeloop [6] and MAESTRO [7] provide systematic approaches for accelerator design space exploration. Timeloop models data movement and compute utilization for any valid mapping of operations to hardware, achieving 5–10% accuracy versus RTL simulation at $2000\times$ speedup. MAESTRO offers a data-centric perspective using intuitive dataflow directives. Sparseloop [8] extends these frameworks to sparse tensor operations, critical for efficient transformer inference.

ML-based approaches complement analytical models by learning residual corrections or capturing effects not modeled analytically. ArchGym [15] demonstrates that ML surrogate models can achieve 0.61% RMSE while providing $2000\times$ speedup over simulation, enabling rapid design space exploration for accelerator development.

4) *Edge and Mobile Device Modeling*: Edge devices impose strict power, memory, and latency constraints, making accurate prediction essential for deploying ML models on mobile phones, IoT devices, and embedded systems. The diversity of edge hardware—spanning mobile CPUs, mobile GPUs, NPUs, and DSPs—creates significant challenges for cross-platform prediction.

nn-Meter [4] addresses this challenge through kernel-level prediction with adaptive sampling, achieving 99% accuracy across mobile CPUs, GPUs, and Intel VPUs. LitePred [5] extends this work with transfer learning, achieving 99.3% accuracy across 85 edge platforms with less than one hour of adaptation per new device. These results demonstrate that ML models can effectively generalize across the heterogeneous edge hardware landscape.

5) *Distributed System Modeling*: Multi-GPU and multi-node systems introduce communication overhead, synchronization barriers, and parallelism strategy choices that fundamentally change performance characteristics. Distributed training performance depends on the interplay between compute, memory bandwidth, and network communication.

ASTRA-sim [11] provides end-to-end distributed training simulation, modeling collective communication algorithms, network topology, and compute-communication overlap. VIDUR [16] focuses specifically on LLM inference serving, capturing the unique characteristics of prefill and decode phases, KV cache management, and request scheduling. These simulation frameworks achieve 5–15% accuracy versus real clusters while enabling exploration of parallelization strategies at scale.

B. By ML Technique

The choice of ML technique reflects trade-offs between accuracy, data efficiency, interpretability, and generalization capability.

1) *Classical Machine Learning*: Tree-based ensembles—random forests and gradient boosted trees (XGBoost, LightGBM)—remain highly effective for performance modeling, particularly in low-data regimes. These methods handle non-linear relationships through recursive partitioning, provide

feature importance rankings for interpretability, and require minimal hyperparameter tuning.

Classical ML models dominate when training data is limited (<10K samples) or when features are well-engineered. nn-Meter [4] demonstrates that random forests achieve competitive accuracy with careful kernel-level feature engineering. The ALCOP framework combines XGBoost with analytical pre-training, using analytical model predictions as features to accelerate autotuning convergence.

2) *Deep Learning*: Multi-layer perceptrons (MLPs) learn hierarchical feature representations without manual feature engineering. MLPs are widely used as the prediction head in more complex architectures and as standalone models when sufficient training data is available. NeuSight [3] uses MLPs to predict tile-level GPU utilization, learning complex interactions between tile parameters and hardware characteristics.

Recurrent neural networks (RNNs and LSTMs) process sequential inputs, making them suitable for modeling operator sequences in neural network execution. However, sequential processing limits parallelization, and attention-based architectures increasingly replace RNNs for sequence modeling tasks.

3) *Graph Neural Networks*: Graph neural networks (GNNs) have emerged as particularly effective for performance modeling because computational graphs have natural graph structure. Nodes represent operators with features (type, parameters, shapes), edges represent data dependencies with features (tensor dimensions, datatypes). GNNs propagate performance-relevant information along these dependencies through message passing.

GRANITE [12] applies GNNs to basic block throughput estimation, learning to predict CPU performance from instruction dependency graphs. For DNN workloads, GNN-based models capture inter-operator interactions that flat feature representations miss. The graph structure also enables natural handling of variable-size networks without padding or truncation.

4) *Hybrid Analytical+ML Models*: Hybrid approaches combine physics-based analytical models with learned components, achieving both interpretability and high accuracy. The analytical component provides a strong prior based on hardware characteristics, while the ML component learns residual corrections and complex interactions.

This design philosophy has produced state-of-the-art results. Analytical pre-training initializes ML models with reasonable predictions, reducing data requirements and improving convergence. Physics-informed architectures incorporate analytical insights into model structure—NeuSight’s tile-based prediction mirrors CUDA’s execution model, providing inductive bias that improves generalization. Residual learning trains ML models to predict the error of analytical models, combining analytical interpretability with ML’s ability to capture unmodeled effects.

The latency predictor study [17] demonstrates that hybrid approaches with transfer learning achieve 22.5% average improvement over baselines, with up to 87.6% improvement on challenging cross-platform prediction tasks.

C. By Input Representation

Input representation determines what information the model can access and how effectively it can learn performance-relevant patterns.

1) *Static Features*: Static features derive from workload and hardware specifications without runtime measurement. For DNN workloads, these include layer parameters (kernel size, channels, stride, batch size), tensor dimensions, and operator types. Hardware specifications include core counts, memory sizes, bandwidth, and clock frequencies.

Static features enable prediction without profiling, supporting use cases like neural architecture search where thousands of candidate networks must be evaluated. Feature engineering plays a critical role: effective representations capture computation-to-communication ratios, memory footprint estimates, and parallelization potential.

2) *Hardware Counters*: Performance counters provide runtime measurements of hardware behavior: cache miss rates, memory bandwidth utilization, instruction throughput, and stall cycles. Counter-based models can capture dynamic effects invisible to static analysis, including contention, thermal throttling, and runtime scheduling decisions.

The primary limitation is that counter-based models require hardware execution, limiting their applicability for design space exploration or new architecture evaluation. However, for optimizing existing deployments or debugging performance anomalies, counter-based models provide valuable insights that static approaches cannot match.

3) *Graph Representations*: Graph representations encode computational graphs with nodes representing operators and edges representing data dependencies. Node features capture operator characteristics (type, parameters), while edge features encode tensor properties (shape, datatype, memory format).

Graph representations provide several advantages over flat feature vectors: they naturally handle variable-size networks, preserve structural information about operator interactions, and enable permutation-invariant predictions. GNNs operating on these representations can learn which subgraph patterns indicate performance bottlenecks.

4) *Learned Embeddings*: Learned embeddings compress high-dimensional or categorical information into dense vector representations. Hardware embeddings represent diverse devices as points in a learned feature space, enabling transfer learning across platforms. Operator embeddings capture semantic similarities between operator types that may share performance characteristics.

HELP formulates hardware prediction as meta-learning, learning hardware embeddings that represent devices as black-box functions. With just 10 measurement samples on a new device, HELP achieves accurate predictions by positioning the device appropriately in the learned embedding space. This approach is particularly valuable for the fragmented edge hardware landscape, where collecting exhaustive training data for each device is impractical.

Table I summarizes representative papers across our taxonomy dimensions, illustrating the diversity of approaches and their key characteristics.

IV. SURVEY OF APPROACHES

This section surveys ML-based performance modeling approaches organized by target hardware platform. For each category, we examine the modeling challenges specific to that platform, describe representative techniques, and synthesize key findings across the literature. Table II provides a comprehensive comparison of the surveyed approaches.

A. CPU Performance Modeling

CPU performance modeling for ML workloads presents unique challenges due to complex microarchitectural effects including out-of-order execution, branch prediction, and deep cache hierarchies. While GPUs have received more attention for DNN training, CPUs remain important for inference—particularly on edge devices and for operators that map poorly to SIMT execution.

1) *Traditional CPU Performance Modeling*: Traditional CPU modeling relies on cycle-accurate simulation through frameworks like gem5 [2]. The gem5 simulator provides multiple fidelity levels: fast functional simulation for correctness validation, and detailed out-of-order models achieving 10–20% accuracy versus real hardware. For ML workloads, gem5 extensions such as gem5-Aladdin and SMAUG enable accelerator integration studies.

However, cycle-accurate simulation suffers from fundamental speed limitations—simulating even modest DNN inference requires hours, making design space exploration impractical. This limitation has motivated ML-based alternatives that learn to predict performance from static program features.

2) *ML-Based Basic Block Modeling*: GRANITE [12] represents the state of the art in ML-based CPU performance modeling. The key insight is that basic block throughput—the steady-state execution rate of a loop body—can be predicted from the instruction dependency graph without simulation. GRANITE encodes basic blocks as directed graphs where nodes represent instructions with features (opcode, operand types) and edges capture data dependencies.

A graph neural network processes this representation through message passing layers:

$$\mathbf{h}_i^{(k+1)} = \text{MLP} \left(\mathbf{h}_i^{(k)} + \sum_{j \in \mathcal{N}(i)} \mathbf{h}_j^{(k)} \right) \quad (6)$$

where $\mathbf{h}_i^{(k)}$ represents instruction i ’s embedding at layer k . After several message passing rounds, a global pooling operation aggregates instruction embeddings into a single block representation, which a final MLP maps to throughput prediction.

GRANITE achieves 0.97 Kendall’s τ correlation with ground-truth measurements on x86 basic blocks, significantly outperforming prior analytical models like IACA and llvm-mca. Critically, the learned model generalizes across

microarchitectures—a model trained on Skylake transfers to Haswell with only modest accuracy degradation.

3) *Challenges and Opportunities*: Despite GRANITE’s success, several challenges remain for CPU performance modeling. First, DNN operators often involve memory-bound execution where cache behavior dominates—GRANITE focuses on compute-bound basic blocks and does not model memory hierarchy effects. Second, modern CPUs feature increasingly complex prefetchers and branch predictors whose behavior is difficult to capture in static features. Third, CPU-based DNN inference often involves highly optimized library code (Intel MKL, ARM Compute Library) whose performance depends on runtime scheduling decisions.

Hybrid approaches combining coarse-grained simulation with learned correction factors represent a promising direction. Rather than simulating every cycle, these methods use fast simulation to establish approximate behavior, then train ML models to predict residual errors, potentially achieving simulation accuracy at reduced cost.

B. GPU Performance Modeling

GPGPU-Sim [9] pioneered detailed GPU simulation, modeling SIMT cores, warp scheduling, memory coalescing, and cache hierarchies. Accel-Sim [10] extended this foundation with trace-driven simulation and improved correlation with modern GPUs (Turing, Ampere), achieving 0.90–0.97 IPC correlation.

1) *Cycle-Accurate GPU Simulation*: GPGPU-Sim [9] pioneered detailed GPU simulation, modeling SIMT cores, warp scheduling, memory coalescing, and cache hierarchies. Accel-Sim [10] extended this foundation with trace-driven simulation and improved correlation with modern GPUs (Turing, Ampere), achieving 0.90–0.97 IPC correlation.

These simulators provide high fidelity—essential for microarchitectural studies—but suffer from 1000–10000 \times slowdown versus real GPU execution. Simulating a single ResNet-50 inference can require hours, making design space exploration impractical. This has motivated the development of ML-based predictors that achieve comparable accuracy at dramatically reduced cost.

2) *Learned GPU Performance Models*: Habitat [14] introduced *wave scaling* for cross-GPU prediction. The key insight is that GPU execution time can be decomposed into compute and memory components that scale differently across devices:

$$T_{\text{target}} = T_{\text{compute}} \cdot \frac{P_{\text{source}}}{P_{\text{target}}} + T_{\text{memory}} \cdot \frac{B_{\text{source}}}{B_{\text{target}}} \quad (7)$$

where P denotes peak compute throughput and B memory bandwidth. By profiling on a source GPU and measuring how kernels respond to artificially reduced parallelism (“wave scaling”), Habitat estimates the compute and memory fractions, enabling prediction on unseen target GPUs.

Habitat achieves 11.8% average error predicting training iteration time across GPU generations (V100 to A100). However, it requires actual GPU execution for wave scaling measurements and cannot predict performance for unseen models.

TABLE I
REPRESENTATIVE PAPERS CLASSIFIED BY OUR TAXONOMY DIMENSIONS. ACCURACY REPORTED AS MAPE OR CORRELATION WHERE AVAILABLE.

Paper	Target	Technique	Input	Accuracy	Key Contribution
NeuSight [3]	GPU	Hybrid	Static	2.3%	Tile-based prediction
nn-Meter [4]	Edge	Classical ML	Static	<5%	Kernel detection
LitePred [5]	Edge	Transfer	Static	0.7%	85-platform transfer
GRANITE [12]	CPU	GNN	Graph	0.97 corr	Basic block modeling
Timeloop [6]	Accelerator	Analytical	Static	5–10%	Loop-nest DSE
ASTRA-sim [11]	Distributed	Simulation	Traces	5–15%	Collective modeling
ArchGym [15]	Accelerator	Hybrid	Static	0.61% RMSE	ML-aided DSE

TABLE II
SUMMARY OF SURVEYED ML-BASED PERFORMANCE MODELING APPROACHES, ORGANIZED BY TARGET HARDWARE PLATFORM.

Paper	Platform	ML Technique	Prediction Target	Error	Key Innovation
<i>CPU Performance Modeling</i>					
GRANITE [12]	CPU	GNN	Basic block throughput	0.97 corr	Instruction graph encoding
gem5+ML [2]	CPU	Hybrid	Execution time	10–20%	Simulation + learning
<i>GPU Performance Modeling</i>					
NeuSight [3]	GPU	Hybrid MLP	Kernel/E2E latency	2.3%	Tile-based prediction
Habitat [14]	GPU	MLP	Training time	11.8%	Wave scaling analysis
Accel-Sim [10]	GPU	Simulation	Cycle-accurate	10–20%	SASS trace-driven
<i>DNN Accelerator Modeling</i>					
Timeloop [6]	NPU	Analytical	Latency/Energy	5–10%	Loop-nest DSE
MAESTRO [7]	NPU	Analytical	Latency/Energy	5–15%	Data-centric directives
Sparseloop [8]	NPU	Analytical	Sparse tensors	5–10%	Compression modeling
ArchGym [15]	Multi	RL+Surrogate	Multi-objective	0.61%	ML-aided DSE
<i>Edge Device Modeling</i>					
nn-Meter [4]	Edge	RF ensemble	Latency	<1%	Kernel detection
LitePred [5]	Edge	VAE+MLP	Latency	0.7%	85-platform transfer
HELP [18]	Multi	Meta-learning	Latency	1.9%	10-sample adaptation
<i>Distributed and LLM Systems</i>					
ASTRA-sim [11]	Distributed	Simulation	Training time	5–15%	Collective modeling
VIDUR [16]	GPU cluster	Simulation	LLM serving	<5%	Prefill/decode phases

NeuSight [3] addresses these limitations through *tile-based prediction*. The key innovation is decomposing GPU kernel execution into tiles—the basic scheduling unit in CUDA—and predicting per-tile behavior:

$$T_{\text{kernel}} = \max_{w \in \text{waves}} \sum_{t \in w} (T_{\text{compute}}^{(t)} + T_{\text{memory}}^{(t)}) \quad (8)$$

This formulation mirrors actual GPU execution semantics: tiles are scheduled in waves, and kernel time is dominated by the slowest wave. NeuSight uses MLPs to predict tile-level compute and memory times from static features (tile dimensions, register usage, shared memory allocation).

By capturing the wave-level structure, NeuSight achieves remarkable accuracy: 2.3% error on GPT-3 inference across H100, A100, and V100 GPUs. This represents a $50\times$ reduction in error compared to prior approaches like Habitat (121.4% \rightarrow 2.3% on H100 for GPT-3). NeuSight’s physics-informed architecture—encoding GPU execution semantics into the model structure—provides strong inductive bias that enables generalization to unseen models and GPUs.

3) *Compiler Cost Models for GPUs*: The TVM [13] and Ansor [19] systems use learned cost models to guide tensor program optimization. Rather than executing every candidate program, XGBoost or MLP models predict execution time from program features (loop bounds, vectorization widths, memory access patterns).

Ansor’s hierarchical search combines sketch generation, random annotation, and evolutionary refinement, using the cost model to prune the search space. With 10K profiled samples, Ansor achieves approximately 15% MAPE on GPU kernel prediction. The TenSet dataset provides 52 million program performance records across CPUs and GPUs, enabling pre-trained cost models that accelerate autotuning convergence by $10\times$.

4) *LLM Inference Prediction*: Large language model inference presents unique GPU modeling challenges. LLM execution exhibits distinct *prefill* (compute-bound, parallel prompt processing) and *decode* (memory-bound, sequential token generation) phases with fundamentally different performance characteristics.

VIDUR [16] provides discrete-event simulation for LLM

serving systems. Rather than modeling GPU microarchitecture, VIDUR simulates request scheduling, KV cache management, and batching decisions—the system-level factors that dominate serving performance. VIDUR achieves <5% error on end-to-end serving metrics including time-to-first-token and request latency.

Roofline-LLM extends traditional roofline analysis to LLM inference by decomposing transformer execution into compute-bound (prefill attention, FFN) and memory-bound (decode attention, KV cache access) components. Combined with learned correction factors, this hybrid approach achieves 87% reduction in MSE compared to pure roofline predictions.

C. Accelerator Performance Modeling

DNN accelerators—including TPUs, NPUs, and custom ASIC designs—employ specialized dataflows and memory hierarchies optimized for tensor operations. Modeling these devices requires understanding the interaction between dataflow choices, memory hierarchy utilization, and workload characteristics.

1) *Analytical Accelerator Modeling*: Timeloop [6] provides the foundational framework for DNN accelerator design space exploration. The key insight is that accelerator performance can be accurately predicted from loop-nest representations of tensor computations. For a given architecture specification and mapping (loop order, tiling, spatial distribution), Timeloop analytically computes:

- **Data reuse** at each memory level: how many times each tensor element is accessed from each buffer
- **Latency**: compute cycles plus memory stall cycles based on bandwidth constraints
- **Energy**: access counts multiplied by per-access energy at each memory level

Timeloop decouples architecture specification (PEs, buffer sizes, bandwidth) from mapping decisions, enabling systematic exploration of dataflow choices. The framework achieves 5–10% accuracy versus RTL simulation while providing 2000× speedup, making million-point design sweeps tractable.

MAESTRO [7] offers a complementary *data-centric* perspective. Rather than loop-nest transformations, MAESTRO models performance through data movement analysis using compact dataflow directives. This representation is more intuitive—designers specify how tensors flow through the architecture rather than manipulating loop indices—while achieving comparable accuracy.

Sparseloop [8] extends analytical modeling to sparse tensor accelerators. The key challenge is that sparse execution time depends on runtime sparsity patterns, not just static tensor dimensions. Sparseloop models compression formats (CSR, bitmap, RLE), gating logic, and sparse-dense conversion overhead, enabling accurate prediction for pruned neural networks and sparse attention patterns.

2) *ML-Augmented Accelerator Design*: ArchGym [15] demonstrates how ML-based surrogate models can accelerate accelerator design. The framework connects ML optimization algorithms (reinforcement learning, Bayesian optimization,

evolutionary strategies) to hardware simulators through a unified interface.

A key finding is the *hyperparameter lottery*: ML algorithms show high variance across hyperparameter choices, with optimal settings differing substantially between target designs. ArchGym addresses this through systematic hyperparameter sweeps enabled by fast surrogate models. Trained surrogate models achieve 0.61% RMSE while providing 2000× speedup over simulation, enabling exploration of hyperparameter configurations that would be intractable with direct simulation.

3) *FPGA and Emerging Accelerator Modeling*: FPGA-based accelerators present additional modeling challenges due to the flexibility of reconfigurable fabric and the complexity of HLS-generated datapaths. Recent work applies transfer learning to FPGA design space exploration: models trained on one design can adapt to new architectures with limited additional profiling.

Emerging accelerators—including processing-in-memory (PIM), neuromorphic, and analog compute-in-memory designs—remain underexplored. These platforms exhibit fundamentally different performance characteristics (energy-dominated by activations, analog noise effects, sparse event-driven computation) that existing frameworks do not address. Developing unified modeling approaches for this diverse hardware landscape represents an important open challenge.

D. Memory System Modeling

Memory system behavior increasingly dominates ML workload performance. Large language models may require hundreds of gigabytes for weights and KV cache, while training workloads stress memory bandwidth through gradient communication. Accurate memory modeling is essential for understanding performance across the modern hardware landscape.

1) *Cache and Memory Hierarchy Modeling*: Traditional memory system modeling relies on cache simulation within frameworks like gem5 [2] and GPGPU-Sim [9]. These simulators model replacement policies, bank conflicts, memory coalescing (for GPUs), and DRAM controller behavior with high fidelity.

For DNN workloads, memory access patterns are often highly regular—streaming through weight and activation tensors—making analytical prediction feasible. Timeloop [6] models memory hierarchy through data reuse analysis: given a tiling and loop order, the framework computes exact access counts at each memory level. This analytical approach achieves high accuracy for regular workloads but may miss dynamic effects like cache contention in multi-tenant scenarios.

2) *KV Cache for LLM Inference*: KV cache management has emerged as the dominant memory challenge for LLM serving. The attention mechanism requires storing key-value tensors for all previously generated tokens, with memory growing linearly with sequence length and batch size. For long-context models serving concurrent requests, KV cache can consume hundreds of gigabytes.

vLLM’s PagedAttention introduces virtual memory concepts to KV cache management. By storing KV blocks in non-contiguous physical memory with page tables for address translation, PagedAttention achieves near-zero memory waste from fragmentation. This system-level optimization yields 2–4 \times throughput improvement over prior approaches.

VIDUR [16] models KV cache behavior at the serving system level, simulating allocation, eviction, and paging decisions that affect request latency. More recent work explores KV cache compression through quantization (Oaken), sparsity (ALISA), and adaptive token selection (MorphKV), with potential memory savings exceeding 50%. Accurate performance models for these compression techniques—predicting the latency-accuracy tradeoff for different compression levels—remain an open challenge.

3) *Distributed Memory and Communication*: Multi-GPU and multi-node training introduces communication overhead that can dominate performance at scale. ASTRA-sim [11] provides end-to-end simulation of distributed training, modeling collective communication algorithms (ring, tree, halving-doubling all-reduce), network topology, and compute-communication overlap.

The simulation decomposes collective operations into point-to-point messages, tracks network contention, and models the interaction between computation and communication phases. ASTRA-sim achieves 5–15% error versus real multi-GPU clusters, enabling exploration of parallelization strategies (data parallel, model parallel, pipeline parallel) before expensive hardware experiments.

A key insight from distributed training modeling is that communication overhead depends strongly on message granularity and overlap opportunities. Chunked gradient communication, where gradients are transmitted in pieces overlapped with backward pass computation, can hide communication latency. Accurate modeling of this overlap—which depends on operator ordering, chunk sizes, and network bandwidth—is essential for predicting distributed training performance.

E. Cross-Platform and Transfer Learning

The proliferation of hardware platforms—from edge devices to datacenter GPUs to custom accelerators—creates demand for performance models that generalize across configurations. Training separate models for each target device is impractical given the diversity of the hardware landscape. Transfer learning and meta-learning approaches address this challenge by learning shared representations that adapt efficiently to new platforms.

1) *Hardware-Adaptive Latency Prediction*: HELP [18] formulates cross-hardware prediction as meta-learning. The key insight is that hardware platforms can be treated as “tasks” in meta-learning: each device provides a small sample of profiled networks, and the goal is rapid adaptation to new devices.

HELP learns:

- **Architecture encoder**: A GNN that embeds neural network architectures into a fixed-dimensional space

- **Hardware encoder**: A learned function that represents devices from their profiled samples
- **Predictor**: An MLP that maps (architecture, hardware) pairs to latency

Using MAML-style meta-learning, HELP achieves 93.2% accuracy with just 10 profiled samples on new devices, reaching 98.1% with 100 samples. This sample efficiency is critical for the fragmented edge hardware landscape where collecting exhaustive training data for each device type is impractical.

2) *Transfer Learning at Scale*: LitePred [5] scales cross-platform prediction to 85 edge devices—the most comprehensive evaluation to date. The framework introduces a VAE-based data sampler that intelligently selects which architectures to profile on new devices. Rather than random sampling, the VAE identifies architectures that are most informative for learning the device’s performance characteristics.

With less than one hour of profiling on a new device, LitePred achieves 99.3% accuracy on held-out architectures. This combines pre-trained representations from source platforms with efficient adaptation, demonstrating that the cross-platform transfer learning problem is tractable even at scale.

The latency predictors study [17] provides a systematic comparison of transfer learning approaches for NAS. Key findings include:

- End-to-end training on pooled multi-platform data outperforms sequential fine-tuning
- Transfer learning provides 22.5% average improvement over training from scratch
- Benefits are largest for challenging cross-platform transfers (up to 87.6% improvement)

3) *Hybrid Analytical-ML Transfer*: Hybrid approaches combine analytical models with learned components to improve transfer efficiency. SynPerf decomposes GPU kernel execution into pipeline demands (compute, memory, cache) using analytical models, then trains MLPs to capture cross-pipeline interactions. The analytical decomposition provides physics-based structure that transfers across GPUs, while the learned component captures device-specific effects.

This hybrid architecture achieves 6.1% kernel-level error and has been applied to guide Triton kernel optimization, demonstrating 1.7 \times speedup on generated kernels. The combination of interpretable analytical structure with learned flexibility represents a promising direction for transferable performance modeling.

4) *Open Challenges in Transfer Learning*: Despite progress, several challenges remain. First, most transfer learning work focuses on CNN architectures; transformers and mixture-of-experts models remain underexplored. Second, transfer across *workload types* (not just hardware) is challenging—models trained on vision networks may not transfer to language models or graph neural networks. Third, continual learning for performance models—adapting to hardware and software evolution over time—is largely unexplored.

Foundation models for performance prediction represent an emerging opportunity. Pre-trained on large-scale profiling

datasets spanning diverse architectures and hardware, such models could provide strong initialization for any new prediction task. The TenSet dataset with 52 million records represents a step in this direction, but comprehensive datasets covering the full range of modern workloads and hardware remain to be developed

V. COMPARISON AND ANALYSIS

Having surveyed the landscape of ML-based performance modeling approaches, we now provide a comparative analysis across key dimensions, including commonly used analytical and simulation-based baselines. This analysis synthesizes trade-offs that practitioners face when selecting or developing performance models, examining accuracy, training cost, generalization, and interpretability. Table III provides a comprehensive comparison across these dimensions.

A. Accuracy vs. Training Cost

A fundamental trade-off exists between prediction accuracy and the cost of data collection and model training. We analyze this trade-off across the surveyed approaches, identifying regimes where different techniques excel.

1) *Data Collection Overhead*: The cost of obtaining training data varies dramatically across approaches. *Profiling-based methods* require executing workloads on target hardware, with costs ranging from minutes (single operators) to hours (full model sweeps). nn-Meter [4] requires approximately 1,000 profiled samples per kernel type per device, translating to several hours of automated measurement. LitePred [5] reduces this to approximately 100 samples for new devices through intelligent VAE-based sampling.

Simulation-based training uses cycle-accurate or analytical simulators as ground truth. ArchGym [15] trains surrogate models on Timeloop [6] outputs, avoiding real hardware entirely but requiring validated simulator configurations. This approach achieves 0.61% RMSE while providing 2000 \times speedup over direct simulation.

Transfer learning amortizes data collection across platforms. HELP [18] demonstrates that meta-learning enables 93.2% accuracy with just 10 samples on new devices, reaching 98.1% with 100 samples. This sample efficiency is critical for the fragmented edge hardware landscape.

2) *Model Training Cost*: Training complexity varies from minutes for classical ML to days for large-scale pre-training. Tree-based ensembles (random forests, XGBoost) train in minutes on modest datasets and require minimal hyperparameter tuning. Deep learning models require careful architecture design, regularization, and often GPU training, but can achieve higher accuracy on large datasets.

The TenSet dataset [20] with 52 million tensor program performance records enables pre-trained cost models that accelerate autotuning convergence by 10 \times . However, creating such datasets requires substantial infrastructure investment.

3) *Accuracy Stratification*: We observe three accuracy tiers across the surveyed approaches:

Tier 1 (<5% error): Specialized models achieving near-perfect accuracy on narrow domains. nn-Meter achieves <1% error on edge device latency through kernel-level decomposition. NeuSight reaches 2.3% error on GPU inference through physics-informed tile-based prediction. LitePred achieves 0.7% error across 85 edge platforms through extensive transfer learning.

Tier 2 (5–15% error): General-purpose models with broader applicability. Habitat achieves 11.8% error on cross-GPU prediction using wave scaling. Analytical frameworks like Timeloop and MAESTRO typically achieve 5–15% error versus RTL simulation.

Tier 3 (15–25% error): Compiler cost models optimized for ranking rather than absolute accuracy. TVM’s AutoTVM [13] achieves approximately 20% MAPE, sufficient for guiding autotuning search. These models prioritize speed and online adaptation over absolute precision.

The key insight is that accuracy requirements depend on the use case: neural architecture search may tolerate 10–15% error if rankings are preserved, while hardware cost estimation for procurement decisions demands <5% accuracy.

B. Generalization Capabilities

Generalization—the ability to predict accurately on unseen workloads, configurations, or hardware—is perhaps the most critical capability for practical deployment. We analyze generalization along three axes: workload generalization, hardware generalization, and temporal generalization.

1) *Workload Generalization*: Models must handle neural network architectures not seen during training. GNN-based approaches offer natural workload generalization because the graph structure captures compositional relationships. GRANITE [12] generalizes across basic blocks by learning instruction-level patterns that compose into block-level predictions.

However, generalization often fails across workload types. Models trained on CNNs may not transfer to transformers due to fundamentally different computational patterns. NeuSight [3] addresses this by training on diverse operator types (GEMM, attention, convolution) and learning GPU execution semantics that generalize across operations.

2) *Hardware Generalization*: Cross-hardware prediction remains challenging due to microarchitectural diversity. Three approaches have shown promise:

Meta-learning treats hardware platforms as tasks. HELP [18] learns hardware embeddings that position devices in a shared latent space, enabling few-shot adaptation to new platforms.

Feature-based transfer uses hardware specifications as input features. LitePred [5] learns relationships between hardware characteristics (compute capability, memory bandwidth) and performance, enabling zero-shot prediction (92.1% accuracy) on entirely new devices.

TABLE III
COMPARATIVE ANALYSIS OF REPRESENTATIVE PERFORMANCE MODELS—INCLUDING ML-BASED AND ANALYTICAL/SIMULATION APPROACHES—ACROSS KEY DIMENSIONS. THE ACCURACY COLUMN REPORTS THE METRIC AND VALUE AS GIVEN IN EACH ORIGINAL WORK (E.G., MAPE, RMSE, KENDALL’S τ , RANGES).

Model	Accuracy (as reported)	Training Data	Adaptation Cost	Generalization	Interpretability	Infer
<i>Classical ML</i>						
nn-Meter [4]	<1% MAPE	1K/kernel	Hours/device	Device-specific	Medium	
XGBoost (TVM) [13]	20% MAPE	10K+	Online	Operator-level	Medium	
<i>Deep Learning</i>						
NeuSight [3]	2.3% MAPE	100K+	Pre-trained	Cross-GPU	Low	
Habitat [14]	11.8% MAPE	Online profiling runs	None (requires GPU)	Cross-GPU	Medium	Per-ke
<i>Graph Neural Networks</i>						
GRANITE [12]	0.97 τ	10K+	Hours	Cross- μ arch	Low	
HELP [18]	1.9% MAPE	Meta-training	10 samples	Cross-platform	Low	
<i>Transfer Learning</i>						
LitePred [5]	0.7% MAPE	85 platforms	100 samples	85+ devices	Low	
<i>Hybrid Analytical+ML</i>						
Timeloop [6]	5–10%	Arch spec	None	Any accelerator	High	
ArchGym [15]	0.61% RMSE	Simulation	Surrogate training	Architecture-specific	Medium	
VIDUR [16]	<5%	Kernel profiles	Per-model	LLM-specific	High	

Analytical decomposition factors predictions into hardware-dependent and hardware-independent components. Habitat [14] decomposes execution into compute and memory components that scale with known hardware parameters, achieving cross-GPU prediction without retraining.

3) *Temporal Generalization*: An underexplored dimension is generalization across time—as software stacks evolve (new compiler versions, framework updates, driver changes), performance characteristics shift. Models trained on older configurations may degrade on current systems.

Continual learning approaches that adapt to evolving hardware-software stacks represent an important open direction. The TenSet dataset’s versioned releases provide a starting point for studying temporal generalization in compiler cost models.

C. Interpretability

Interpretability—understanding *why* a model makes particular predictions—is valuable for debugging, optimization guidance, and building practitioner trust. We categorize approaches by their interpretability characteristics.

1) *Analytical Models: High Interpretability*: Analytical frameworks like Timeloop [6] and MAESTRO [7] provide full interpretability. Predictions decompose into explicit terms: data movement at each memory level, compute utilization, bandwidth constraints. Practitioners can trace high-latency predictions to specific bottlenecks (e.g., “DRAM bandwidth limits this mapping”).

This interpretability enables *actionable insights*: if the model predicts memory-bound execution, the designer knows to explore mappings with better data reuse. The roofline model [1] exemplifies this—identifying compute-bound versus memory-bound regimes immediately suggests optimization directions.

2) *Classical ML: Medium Interpretability*: Tree-based ensembles provide feature importance rankings, indicating which input features most influence predictions. nn-Meter’s kernel-level decomposition enables interpretability: practitioners can identify which kernels dominate latency and focus optimization efforts accordingly.

However, feature importance does not explain *how* features interact. A model may indicate that “kernel size” is important without revealing whether large or small kernels are faster for a given hardware platform.

3) *Deep Learning: Low Interpretability*: Deep neural networks, including GNNs and transformers, function as black boxes. While techniques like attention visualization and gradient-based attribution provide some insight, they rarely yield actionable optimization guidance.

NeuSight [3] partially addresses this through physics-informed architecture: by decomposing predictions into compute and memory components that mirror GPU execution, the model structure itself provides interpretability even though individual weight values remain opaque.

4) *Hybrid Approaches: Balanced Interpretability*: Hybrid analytical+ML models offer a middle ground. The analytical component provides interpretable baselines, while the ML component captures residual effects. When predictions diverge from analytical expectations, practitioners know the difference stems from effects not captured in the analytical model (contention, cache effects, scheduling decisions).

VIDUR [16] exemplifies this for LLM serving: discrete-event simulation provides interpretable system-level behavior, while learned kernel-time predictors capture GPU execution details. The simulation structure enables “what-if” analysis (e.g., “how would P99 latency change with larger batch sizes?”) that pure ML models cannot support.

5) *The Interpretability-Accuracy Trade-off*: A general trade-off exists between interpretability and accuracy. Analytical models sacrifice accuracy for transparency; deep learning models sacrifice transparency for accuracy. For production deployment, hybrid approaches that combine interpretable structure with learned components increasingly represent the best of both worlds.

VI. OPEN CHALLENGES AND FUTURE DIRECTIONS

Despite remarkable progress, significant challenges remain in ML-based performance modeling. This section identifies key open problems and promising research directions that will shape the field’s evolution.

A. Data Availability and Quality

The effectiveness of ML-based performance models fundamentally depends on training data quality and availability. Several challenges persist in this dimension.

1) *Benchmark Diversity*: Existing datasets predominantly cover CNN architectures optimized for image classification. TenSet [20] provides 52 million tensor program records but focuses on operators from ResNet, MobileNet, and similar architectures. Modern workloads—transformers, mixture-of-experts models, graph neural networks, diffusion models—remain underrepresented.

The rapid evolution of model architectures exacerbates this gap. Models trained on 2022-era workloads may poorly predict performance of 2025 architectures featuring sparse attention, conditional computation, or novel activation functions. Continuously updated, community-maintained benchmark suites could address this challenge.

2) *Hardware Coverage*: Hardware diversity creates data collection bottlenecks. LitePred [5] covers 85 edge devices, but the mobile hardware landscape spans hundreds of distinct SoC configurations. Data center hardware (H100, TPU v5, custom accelerators) often has restricted access, limiting public dataset creation.

Simulation-based data generation offers a partial solution: ArchGym [15] trains on Timeloop outputs, avoiding hardware access requirements. However, simulation accuracy itself requires validation against real hardware, creating a chicken-and-egg problem.

3) *Measurement Noise and Reproducibility*: Performance measurements exhibit variance from thermal throttling, OS scheduling, memory allocation, and caching effects. Industrial-strength profiling requires careful warm-up periods, multiple runs, and statistical aggregation. Many published models train on single-run measurements, potentially learning noise rather than signal.

Standardized measurement protocols—specifying warm-up iterations, cooling periods, statistical aggregation methods—would improve cross-study comparability and model reliability.

B. Model Generalization

Generalization remains the central challenge: models that excel on training distributions often fail on realistic deployment scenarios.

1) *Cross-Workload Generalization*: Models struggle to generalize across workload types. A predictor trained on CNNs may fail on transformers due to different computational patterns: CNNs are compute-dominated by convolutions with high data reuse, while transformers feature attention mechanisms with sequence-length-dependent memory access patterns.

Promising directions include workload-agnostic representations (learning from computation graphs rather than architecture-specific features) and multi-task learning across workload families.

2) *Cross-Hardware Generalization*: Hardware generalization faces fundamental obstacles. Different hardware families (CPUs, GPUs, TPUs, FPGAs) employ distinct execution models, memory hierarchies, and parallelism patterns. Even within GPU families, architectural changes (Volta to Ampere to Hopper) introduce new features (tensor cores, TMA, FP8) that alter performance characteristics.

Transfer learning approaches [5], [18] show promise for related hardware, but truly cross-family prediction (e.g., GPU to TPU) remains elusive. Hardware-agnostic intermediate representations that capture essential computational patterns while abstracting platform details could enable broader transfer.

3) *Distribution Shift*: Performance models face distribution shift as software stacks evolve. Compiler optimizations, framework updates, and driver changes alter the workload-to-hardware mapping, invalidating models trained on older configurations.

Online adaptation and continual learning techniques could address distribution shift, but few studies systematically evaluate temporal generalization. Developing benchmarks that explicitly measure robustness to software evolution would accelerate progress.

C. Integration with Design Flows

For ML-based performance models to impact practice, they must integrate seamlessly with existing design and optimization workflows.

1) *Compiler Integration*: Compiler autotuning represents a natural application: ML models guide the search for optimal tensor program configurations. TVM [13] and Ansor [19] demonstrate this integration, but challenges remain.

Cost model accuracy directly affects autotuning efficiency. Mispredictions cause the search to explore suboptimal regions, wasting compilation time. Uncertainty quantification—knowing when predictions are unreliable—could enable more efficient exploration-exploitation trade-offs. Recent uncertainty-aware cost models can provide calibrated uncertainty estimates, but such techniques are not yet standard.

2) *Architecture Exploration*: Hardware design space exploration requires evaluating millions of configurations. ML

surrogate models can accelerate this process, as demonstrated by ArchGym [15], but integration challenges persist.

The design space is often too large for exhaustive surrogate training. Active learning strategies that intelligently select which configurations to simulate could improve sample efficiency. Additionally, surrogate models must provide reliable uncertainty estimates to avoid overconfident predictions that mislead designers.

3) *Serving System Optimization*: LLM serving systems require real-time performance prediction for scheduling decisions. VIDUR [16] provides offline simulation, but online serving requires predictions within microseconds.

Lightweight models suitable for real-time inference, combined with periodic retraining on observed performance, could enable adaptive serving optimization. The challenge is maintaining accuracy while meeting strict latency requirements.

D. Emerging Opportunities

Several emerging trends create new opportunities for ML-based performance modeling.

1) *Foundation Models for Performance Prediction*: The success of foundation models in NLP and vision suggests potential for performance modeling. A foundation model pre-trained on diverse workload-hardware-performance tuples could provide strong initialization for downstream tasks.

Key requirements include: (1) massive, diverse training datasets spanning hardware platforms and workload types; (2) representation learning that captures transferable performance patterns; and (3) efficient adaptation mechanisms for new domains. The TenSet dataset [20] provides a starting point, but orders-of-magnitude more data may be needed for true foundation model capabilities.

2) *LLM-Assisted Performance Analysis*: Large language models offer new modalities for performance understanding. Recent work explores using code LLMs to extract performance-relevant features, explain performance anomalies, and suggest optimizations.

Challenges include hallucination (LLMs may generate plausible-sounding but incorrect performance estimates) and limited numerical reasoning. Hybrid approaches combining LLM-based code understanding with principled performance models could leverage the strengths of both paradigms.

3) *Hardware-Software Co-Design*: ML performance models can enable tighter hardware-software co-design by rapidly evaluating how software changes affect hardware utilization and vice versa.

ArchGym [15] demonstrates ML-guided accelerator design, but joint optimization of hardware architecture, compiler mappings, and model structure remains underexplored. Differentiable performance models could enable gradient-based co-optimization across the full stack.

4) *Emerging Hardware Paradigms*: New computing paradigms—processing-in-memory, neuromorphic computing, analog accelerators, quantum-classical hybrids—require new performance modeling approaches. These platforms exhibit

TABLE IV
SUMMARY OF REPRODUCIBILITY EVALUATION FOR REPRESENTATIVE PERFORMANCE MODELING TOOLS. SCORES REFLECT PRACTICAL REPRODUCIBILITY ON A 10-POINT SCALE.

Tool	Target	Score	Key Issue
Timeloop	Accelerators	9/10	Docker recommended
FlashAttention	GPUs	9/10	GPU required
ASTRA-sim	Distributed	8/10	Complex build
VIDUR	LLM serving	7/10	Python 3.10 only
nn-Meter	Edge devices	5/10	Pickle incompatibility

fundamentally different performance characteristics (energy-dominated costs, stochastic execution, analog noise) that existing frameworks do not address.

Early-stage performance modeling for emerging hardware could accelerate adoption by enabling software optimization before hardware availability. Transfer learning from related platforms (e.g., digital accelerators to analog) represents a promising direction.

5) *Multi-Objective and Constraint-Aware Prediction*: Practical deployment involves multiple objectives: latency, throughput, energy, memory, cost. Most current models predict single metrics independently. Joint multi-objective prediction could enable Pareto-optimal design selection.

Additionally, constraint-aware prediction—determining whether a workload *fits* on target hardware given memory limits—is often more valuable than precise latency estimates. Models that directly predict feasibility and constraint violations could better support deployment decisions.

VII. EXPERIMENTAL EVALUATION

To validate the practical applicability of surveyed performance modeling tools, we conducted hands-on reproducibility evaluations of five representative systems spanning different hardware targets and modeling approaches. This section presents our methodology, tool-by-tool findings, and synthesizes key lessons for practitioners.

A. Evaluation Methodology

We evaluated each tool along four dimensions:

Setup Complexity. We assessed installation difficulty, dependency management, documentation quality, and time to first result. Tools were tested in clean environments following their documented procedures.

Reproducibility. We verified whether example configurations produce consistent results and whether reference outputs are provided for validation.

Practical Usability. We examined API design, configuration flexibility, output interpretability, and integration with existing workflows.

Accuracy Validation. Where possible, we compared tool outputs against published accuracy claims or reference implementations.

Table IV summarizes our findings across all evaluated tools.

B. Tool-by-Tool Results

1) *Timeloop: DNN Accelerator Modeling*: Timeloop [6] provides analytical performance and energy modeling for DNN accelerators through loop-nest analysis.

Setup. Docker-based installation succeeds in 10–15 minutes with pre-built images for both x86 and ARM platforms. Native installation requires 1–2 hours due to complex dependencies (Barvinok, NTL libraries).

Reproducibility. Excellent—reference outputs are provided for all example architectures (Eyeriss, Simba), and results are deterministic. Tutorials with Jupyter notebooks enable systematic learning.

Key Finding. Energy breakdown analysis reveals DRAM dominates (>60%) for typical configurations, validating the importance of dataflow optimization. The mapper may not find globally optimal solutions but provides interpretable trade-off analysis.

2) *ASTRA-sim: Distributed Training Simulation*: ASTRA-sim [11] simulates distributed DNN training with configurable network backends.

Setup. Docker recommended due to Protobuf version sensitivity. Native build requires 1–2 hours with careful dependency management.

Reproducibility. Good—validated configurations for HGX-H100 and DGX-V100 are included. However, reference timing outputs are not provided, requiring trust in published accuracy claims.

Key Finding. The Chakra trace format has a learning curve, but enables detailed collective communication modeling. Multiple network backends (analytical, NS-3) allow accuracy-speed trade-offs.

3) *VIDUR: LLM Inference Simulation*: VIDUR [16] provides discrete-event simulation for LLM serving systems.

Setup. Python-only installation, but strict Python 3.10 requirement creates compatibility issues—Python 3.14 fails due to argparse API changes.

Reproducibility. Good for supported configurations. Pre-profiled data covers A100, H100, and A40 GPUs for Llama-family models. Adding new models requires GPU access for profiling.

Key Finding. Rich scheduler implementations (vLLM, Orca, Sarathi) enable direct algorithm comparison. Metrics include time-to-first-token, time-per-output-token, and memory utilization—essential for SLO-driven capacity planning.

4) *nn-Meter: Edge Device Latency Prediction*: nn-Meter [4] predicts DNN latency on edge devices through kernel-level decomposition.

Setup. Simple pip installation, but critical compatibility issue: pre-trained predictors fail to load with current scikit-learn versions due to pickle format changes.

Reproducibility. Poor in current state—the core functionality (pre-trained predictors) is broken without pinning scikit-learn to version 1.0.x.

Key Finding. This case highlights a critical reproducibility anti-pattern: ML models serialized with pickle are fragile

TABLE V
REPRODUCIBILITY BEST PRACTICES DERIVED FROM TOOL EVALUATION.

Practice	Rationale
Provide Docker images	Isolates dependencies
Document Python version	Prevents API incompatibilities
Include reference outputs	Enables result verification
Use portable model formats	Avoids pickle versioning issues
Pin dependency versions	Ensures reproducible environments

across library versions. Researchers should prefer version-agnostic serialization formats (ONNX, SavedModel) or pin exact dependency versions.

5) *FlashAttention: GPU Attention Optimization*: FlashAttention provides IO-aware attention kernels achieving 2–4× speedup over standard implementations.

Setup. PyPI package available; compilation from source requires 3–5 minutes with ninja. GPU (Ampere or newer) required for both building and running.

Reproducibility. Excellent—widely adopted in major frameworks (HuggingFace, vLLM), comprehensive test suite, and benchmark scripts for validation.

Key Finding. FlashAttention demonstrates successful reproducibility through framework integration rather than standalone distribution. Native integration into PyTorch’s scaled_dot_product_attention ensures continued maintenance and compatibility.

C. Synthesis and Recommendations

Our evaluation reveals systematic patterns affecting reproducibility:

Containerization dramatically improves reproducibility. Tools providing Docker images (Timeloop, ASTRA-sim) achieve higher reproducibility scores by isolating complex dependency chains. Native builds consistently encounter platform-specific issues.

Python version sensitivity is a major concern. VIDUR requires Python 3.10 specifically; nn-Meter’s pickle files are incompatible with current scikit-learn. Projects should document version constraints prominently and consider providing locked dependency specifications.

Pre-trained models age poorly. nn-Meter’s reliance on pickled scikit-learn models created a time bomb. FlashAttention avoids this by focusing on kernel optimization rather than learned models. For projects distributing trained models, ONNX or similar portable formats are preferable.

Reference outputs enable validation. Timeloop’s inclusion of expected outputs for all examples enables immediate verification. ASTRA-sim and VIDUR lack this, requiring users to trust published accuracy claims.

Table V summarizes best practices derived from our evaluation.

VIII. CONCLUSION

This survey has provided a comprehensive analysis of machine learning approaches for computer architecture performance modeling. We have examined over 60 papers spanning

traditional analytical models, simulation-based approaches, and modern ML techniques including classical machine learning, deep learning, graph neural networks, and hybrid methods.

A. Key Findings

Our analysis reveals several key findings that characterize the current state of the field:

ML-based models achieve remarkable accuracy. State-of-the-art approaches achieve prediction errors below 5% for their target domains. NeuSight [3] reaches 2.3% error on GPU inference through physics-informed tile-based prediction. LitePred [5] achieves 0.7% error across 85 edge platforms through transfer learning. These accuracy levels are sufficient for production deployment in neural architecture search, autotuning, and hardware-aware optimization.

Hybrid approaches dominate recent work. The most successful models combine analytical structure with learned components. Analytical decomposition provides interpretable baselines and physics-based inductive bias, while ML captures complex effects that elude closed-form analysis. This hybrid philosophy—exemplified by NeuSight’s tile-based prediction and VIDUR’s [16] simulation-based framework—consistently outperforms pure analytical or pure ML approaches.

Transfer learning is essential for scalability. The proliferation of hardware platforms makes per-device training impractical. Meta-learning (HELP [18]) and VAE-based sampling (LitePred [5]) enable adaptation to new devices with 10–100 samples, demonstrating that cross-platform generalization is tractable.

Kernel-level decomposition improves accuracy. nn-Meter’s [4] insight that end-to-end latency decomposes into kernel latencies has become standard practice. By modeling at the kernel level and capturing framework fusion behavior, models achieve compositional predictions that generalize across architectures.

LLM inference presents unique challenges. Large language model serving has distinct characteristics—autoregressive generation, KV cache growth, prefill-decode phase separation—that require specialized modeling. VIDUR [16] and similar frameworks provide discrete-event simulation capturing these dynamics with <5% accuracy.

B. Promising Research Directions

Looking forward, we identify the most promising directions for advancing the field:

Foundation models for performance prediction. Pre-trained models that transfer across workloads and hardware could dramatically reduce data requirements for new prediction tasks. Creating the large-scale, diverse datasets needed to train such models represents a key community challenge.

Uncertainty quantification. Knowing when predictions are reliable enables better decision-making in autotuning, design space exploration, and serving optimization. Calibrated uncertainty estimates remain underexplored despite their practical importance.

Temporal generalization. As software stacks evolve, performance models must adapt. Continual learning approaches and benchmarks measuring robustness to software evolution deserve increased attention.

Multi-objective prediction. Practical deployment involves latency, throughput, energy, memory, and cost trade-offs. Joint multi-objective prediction could enable Pareto-optimal design selection across these dimensions.

Emerging hardware support. Processing-in-memory, neuromorphic computing, and analog accelerators require new modeling paradigms. Early-stage performance modeling for emerging hardware could accelerate adoption.

C. Concluding Remarks

Machine learning has transformed performance modeling from an art requiring deep architectural intuition to an increasingly systematic discipline. The surveyed approaches demonstrate that learned models can capture complex hardware-software interactions while enabling millisecond-scale prediction. As ML workloads continue to grow in importance and hardware diversity expands, accurate, generalizable performance models will become ever more critical for efficient system design and deployment.

We hope this survey serves as both a comprehensive reference for practitioners selecting performance modeling approaches and a roadmap for researchers identifying impactful open problems. The field’s rapid progress suggests that the coming years will bring continued advances in accuracy, generalization, and practical deployment of ML-based performance models.

REFERENCES

- [1] 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.
- [2] N. Binkert, B. Beckmann, G. Black, S. K. Reinhardt, A. Saidi, A. Basu, J. Hestness, D. R. Hower, T. Krishna, S. Sardashti, R. Sen, K. Sewell, M. Shoaib, N. Vaish, M. D. Hill, and D. A. Wood, “The gem5 simulator,” *ACM SIGARCH Computer Architecture News*, vol. 39, no. 2, pp. 1–7, 2011.
- [3] S. Lee, A. Phanishayee, and D. Mahajan, “NeuSight: GPU performance forecasting via tile-based execution analysis,” in *Proceedings of the 30th ACM International Conference on Architectural Support for Programming Languages and Operating Systems (ASPLOS)*, 2025, pp. 1–15.
- [4] L. L. Zhang, S. Han, J. Wei, N. Zheng, T. Cao, Y. Yang, and Y. Liu, “nn-meter: Towards accurate latency prediction of deep-learning model inference on diverse edge devices,” in *Proceedings of the 19th Annual International Conference on Mobile Systems, Applications, and Services (MobiSys)*, 2021, pp. 81–93, best Paper Award.
- [5] Y. Feng, Z. Li, J. Yang, and Y. Liu, “LitePred: Transferable and scalable latency prediction for hardware-aware neural architecture search,” in *Proceedings of the 21st USENIX Symposium on Networked Systems Design and Implementation (NSDI)*, 2024, pp. 1–18.
- [6] A. Parashar, P. Raina, Y. S. Shao, Y.-H. Chen, V. A. Ying, A. Muber, R. Venkatesan, B. Khailany, S. W. Keckler, and J. Emer, “Timeloop: A systematic approach to DNN accelerator evaluation,” in *Proceedings of the IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, 2019, pp. 304–315.
- [7] H. Kwon, P. Chatarasi, M. Sarber, M. Pellauer, A. Parashar, and T. Krishna, “MAESTRO: A data-centric approach to understand reuse, performance, and hardware cost of DNN mappings,” in *Proceedings of the 52nd IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2019, pp. 1–14.

- [8] Y. N. Wu, J. Emer, and V. Sze, “Sparseloop: An analytical approach to sparse tensor accelerator modeling,” in *Proceedings of the 55th IEEE/ACM International Symposium on Microarchitecture (MICRO)*, 2022, pp. 1–15.
- [9] A. Bakhoda, G. L. Yuan, W. W. L. Fung, H. Wong, and T. M. Aamodt, “Analyzing CUDA workloads using a detailed GPU simulator,” in *Proceedings of the IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, 2009, pp. 163–174.
- [10] M. Khairy, A. Jain, T. M. A. Shen, and T. G. Rogers, “Accel-sim: An extensible simulation framework for validated GPU modeling,” in *Proceedings of the 47th International Symposium on Computer Architecture (ISCA)*, 2020, pp. 473–486.
- [11] W. Won, T. Heo, S. Rashidi, S. Talati, S. Srinivasan, and T. Krishna, “ASTRA-sim2.0: Modeling hierarchical networks and disaggregated systems for large-model training at scale,” in *Proceedings of the IEEE International Symposium on Performance Analysis of Systems and Software (ISPASS)*, 2023, pp. 283–294.
- [12] O. Sykora, A. Rucker, C. Mendis, R. Barik, P. M. Phothilimthana, and S. Amarasinghe, “Granite: A graph neural network model for basic block throughput estimation,” in *Proceedings of the IEEE International Symposium on Workload Characterization (IISWC)*, 2022, pp. 1–13.
- [13] T. Chen, T. Moreau, Z. Jiang, L. Zheng, E. Yan, M. Cowan, H. Shen, L. Wang, Y. Hu, L. Ceze, C. Guestrin, and A. Krishnamurthy, “TVM: An automated end-to-end optimizing compiler for deep learning,” in *Proceedings of the 13th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2018, pp. 578–594.
- [14] G. X. Yu, Y. Gao, P. Golber, and A. Cidon, “Habitat: A runtime-based computational performance predictor for deep neural network training,” in *Proceedings of the USENIX Annual Technical Conference (ATC)*, 2021, pp. 503–521.
- [15] S. Krishnan, A. Yazdanbakhsh, S. Prakash, N. P. Jouppi, J. Parmar, H. Kim, J. Laudon, and C. Narayanaswami, “ArchGym: An open-source gymnasium for machine learning assisted architecture design,” in *Proceedings of the 50th International Symposium on Computer Architecture (ISCA)*, 2023, pp. 1–16.
- [16] A. Agrawal, A. Panwar, J. Mohan, N. Kwatra, B. S. Gulavani, and R. Ramachandran, “VIDUR: A large-scale simulation framework for LLM inference,” in *Proceedings of Machine Learning and Systems (MLSys)*, 2024, pp. 1–15.
- [17] L. Dudziak, T. Chau, M. S. Abdelfattah, R. Lee, H. Kim, and N. D. Lane, “Latency predictors for neural architecture search,” in *Proceedings of Machine Learning and Systems (MLSys)*, 2024, pp. 1–14.
- [18] H. Lee, S. Lee, S. Chong, and S. J. Hwang, “HELP: Hardware-adaptive efficient latency prediction for NAS via meta-learning,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021, pp. 27 016–27 028.
- [19] L. Zheng, C. Jia, M. Sun, Z. Wu, C. H. Yu, A. Haj-Ali, Y. Wang, J. Yang, D. Zhuo, K. Sen, J. E. Gonzalez, and I. Stoica, “Ansor: Generating high-performance tensor programs for deep learning,” in *Proceedings of the 14th USENIX Symposium on Operating Systems Design and Implementation (OSDI)*, 2020, pp. 863–879.
- [20] L. Zheng, R. Liu, J. Shao, T. Chen, J. E. Gonzalez, I. Stoica, and Z. Zhang, “TenSet: A large-scale program performance dataset for learned tensor compilers,” in *Advances in Neural Information Processing Systems (NeurIPS)*, vol. 34, 2021, pp. 29 876–29 888.