

CALO-GNN: Calibrated-Uncertainty Graph Cost Models for Cross-Device TVM Meta-Schedule

Sanjay Kumar Patnala
Georgia Institute of Technology
Atlanta, USA
sanjaykumarpatnala@gmail.com

Abstract

Autotuning dominates the compilation budget for deep-learning models: profiling 5,000–30,000 candidate schedules *per operator* can burn tens of GPU-hours and delay deployment by days. TVM’s Meta-Schedule mitigates this cost with a learned latency predictor, yet its existing models output only point estimates, causing two chronic pathologies: (i) *over-exploitation*, where the tuner locks onto an early but mediocre schedule, and (ii) *over-exploration*, where it thrashes in poorly modelled regions. We introduce **CALO-GNN**, the first *evidential* graph-neural cost model for TVM. CALO-GNN provides single-pass predictions of latency *and* calibrated epistemic uncertainty, enabling a new uncertainty-decaying UCB rule (UEC-UCB). A two-stage transfer procedure reuses 4M historical schedules to warm-start new devices with just 2k measurements. Across seven fused operators and five heterogeneous accelerators—including NVIDIA H100 and a 32-core Xeon—CALO-GNN cuts end-to-end tuning time by **32.4%** and reaches within 5% of oracle performance **1.74×** faster than state-of-the-art baselines, all while respecting a strict 20 ms inference budget.

CCS Concepts

• **Software and its engineering** → **Just-in-time compilers**; • **Computing methodologies** → **Machine learning**; *Uncertainty quantification*.

Keywords

autotuning, graph neural networks, uncertainty calibration, compiler optimization, Bayesian optimization

ACM Reference Format:

Sanjay Kumar Patnala. 2025. CALO-GNN: Calibrated-Uncertainty Graph Cost Models for Cross-Device TVM Meta-Schedule. In *Proceedings of KDD 2025 Workshop on Inference Optimization for Generative AI (Inference Optimization for GenAI ’25)*. ACM, New York, NY, USA, 5 pages. <https://doi.org/10.1145/nnnnnnn.nnnnnnn>

1 Introduction

Modern compilers for machine-learning models face a daunting search problem: a single transformer layer, once lowered to low-level

TensorIR [2], can admit $O(10^4 - 10^5)$ legal *schedules*. Each schedule instantiates concrete loop tilings, thread bindings, vector widths, shared-memory layouts, and unroll factors—design choices that interact non-linearly with device caches and tensor shapes. Brute-force measurement is prohibitive: profiling 30k schedules on an H100 at 1ms per kernel already costs 8.3GPU-hours.

Meta-Schedule and its limitations. TVM Meta-Schedule [14] tackles this space with an evolutionary search guided by a latency *cost model*. The original XGBoost model and its GNN successor MetaTune [13] make point predictions; they cannot express *epistemic uncertainty*. Consequently, early “lucky” predictions dominate exploitation (*over-exploit*), while unexplored but risky regions consume many probes (*over-explore*). Our empirical results show that more than 40% of total measurements on depthwise convolutions are wasted in this manner.

Our thesis. Well-calibrated uncertainty is as valuable as accuracy in a compiler cost model. If the tuner knows *where* its predictions are uncertain, it can steer probes strategically—yet any viable solution must keep model inference under 20 ms to avoid dominating the inner loop.

Contributions.

- (1) **CALO-GNN**: a 480-line evidential GNN that returns mean latency and variance in one forward pass.
- (2) **UEC-UCB**: an uncertainty-decaying exploration rule derived from empirical variance scaling, eliminating exhaustive κ grid search.
- (3) **Two-stage transfer**: pre-train on 4M A100 schedules, then fine-tune 5epochs on 2k warm-up probes, trimming 38% of cold-start cost.
- (4) **Extensive validation**: seven operators, five devices, 35 workload/device pairs, plus ablation, calibration, and transfer studies.

2 Related Work

Compiler autotuning. TVM’s AutoTVM [3] and Ansor [19] automate search but rely on tabular Bayesian models that falter on fused operators. Meta-Schedule [14] replaced simulated annealing with evolution plus XGBoost; MetaTune GNN [13] improved accuracy but remained deterministic. AdaTune [9] adds UCB exploration to AutoTVM, yet its Bayesian linear regressor underfits high-dimensional schedule graphs. TENET [11] offers cross-operator transfer via reweighted trees but cannot quantify uncertainty.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Inference Optimization for GenAI ’25, Toronto, Canada

© 2025 Copyright held by the owner/author(s). Publication rights licensed to ACM.
ACM ISBN 978-x-xxxx-xxxx-x/YYYY/MM
<https://doi.org/10.1145/nnnnnnn.nnnnnnn>

Uncertainty estimation. MC-Dropout [5] and deep ensembles [8] are accurate but require multiple passes. Evidential deep learning [1] uses a Normal-Inverse-Gamma output to produce closed-form mean/variance—ideal for compiler inference budgets. Recent diffusion-based cost models [17] explore sequence priors but remain orders of magnitude slower.

Transfer learning. Proxy transfer via cost-model averages [18] and device-agnostic embeddings [4] lower cold-start cost but often miscalibrate uncertainty. We instead fine-tune the evidential head, achieving calibrated variance within three minutes (§6).

3 Problem Formulation

Consider an operator o on device d with legal schedule space \mathcal{S} . Evaluating a schedule $s \in \mathcal{S}$ returns latency $y(s)$ at hardware cost $c_{\text{meas}} \approx 0.8 \text{ ms} - 2.0 \text{ ms}$ (GPU) or $3 \text{ ms} - 5 \text{ ms}$ (CPU). A learned cost model with parameters θ predicts $\hat{y}_\theta(s)$ in $\leq 20 \text{ ms}$ for an evolution generation of N candidates.

Let $y^* = \min_{s \in \mathcal{S}} y(s)$ be the oracle latency. The tuner seeks the minimal *time-to-quality* (TTQ):

$$\text{TTQ}(\epsilon) = \sum_{t=1}^T [c_{\text{meas}}(s_t) + c_{\text{pred}}(P_t)], \quad \text{s.t. } \min_{t \leq T} y(s_t) \leq (1 + \epsilon) y^*, \quad (1)$$

where P_t is the candidate population at round t and $c_{\text{pred}}(P_t) \leq 20 \text{ ms}$. We fix $\epsilon = 0.05$ (5% of oracle) following Meta-Schedule practice [14]. The optimisation thus decomposes into (i) reducing the probe count T via uncertainty-aware exploration and (ii) keeping prediction runtime negligible relative to measurement.

In the next section, we describe how CALO-GNN achieves both goals via evidential regression and an uncertainty-decaying UCB rule.

4 Methodology

Our design goal is to inject *trustworthy uncertainty* into TVM’s evolutionary search while respecting a tight inference-time budget. Figure 1 outlines the four key modules: (1) evidential graph regression, (2) UEC-UCB acquisition, (3) two-stage cross-device transfer, and (4) low-overhead system integration.

4.1 Evidential Graph Regression

Schedule graph encoding. A schedule s is represented as a directed multi-graph $G = (V, E)$, where each node encodes a loop or memory scope and edges capture data, fusion, or reorder dependencies. Node features are a 16-dimensional vector concatenating (loop extent, tile size, parallel scope ID, vector width, shared-mem bank count, unroll flag, etc.); edge types are one-hot. We normalise continuous features to zero mean and unit variance per workload to stabilise training.

GraphSAGE layer choice. GraphSAGE [6] achieved the best accuracy-latency trade-off among GCN [7], GAT [16], and GraphSAGE variants. Two layers suffice: deeper stacks over-smooth node embeddings and add $> 2 \text{ ms}$ to inference. Hidden size 64 is the Pareto knee; 128 improves RMSE by only 0.3% but doubles FLOPs.

Evidential head and calibration. The four evidence parameters ($\gamma, \nu, \alpha, \beta$) are constrained by softplus or sigmoid activations

to keep the Normal-Inverse-Gamma posterior valid ($\alpha, \beta > 0, \nu > 0$) [1]. We append a scalar temperature τ learned on the last epoch via BFGS to minimise validation ECE. At test time, $\sigma_{\text{cal}}^2 = \tau^2 \sigma^2$.

Training schedule. We use AdamW [10], batch size 512, initial learning rate 0.001 with cosine decay and 1epoch warm-up. Early stopping at 2 epochs patience on validation NLL. Full pre-training takes 21 minutes on a single A100 (82 W TDP).

4.2 UEC-UCB: Uncertainty-Decaying UCB

Classical GP-UCB [15] sets $\kappa_t = \sqrt{\beta_t}$, where β_t depends on kernel smoothness and dimensionality—constants unknown for heterogeneous schedule graphs. Practitioners instead tune κ by grid search, an option we found brittle: the same κ that excels on Conv2D can stall attention kernels.

We therefore derive a *data-driven* decay schedule. Let $\bar{\sigma}_t$ be the mean epistemic variance over the current population; empirically, $\bar{\sigma}_t \approx c/\sqrt{t}$ with $c \approx 1.5$. Substituting this empirical law into the GP-UCB regret bound yields our closed form: $\kappa_t = \frac{1.5}{\sqrt{t+50}}$. The offset 50 avoids over-exploration during the warm-up window.

4.3 Two-Stage Cross-Device Transfer

We hypothesise that *operator semantics dominate device semantics* for predictor accuracy, whereas variance calibration is device-sensitive. Stage1 therefore maximises coverage by packing 38 operators across vision, language, and recommender workloads on a single donor GPU (A100). Stage2 fine-tunes for only 5 epochs ($\sim 3 \text{ min}$) using the warm-up probes that the autotuner would measure anyway, incurring no extra hardware cost. This strategy draws inspiration from proxy-based transfer in AutoTVM [18] and device-agnostic embeddings in Daphne [4], but adapts them to evidential models.

Why not zero-shot? Zero-shot variance can be miscalibrated by up to $3\times$ on H200 due to cache-hierarchy mismatch, leading to premature exploitation and 54% slower convergence.

4.4 System-Level Optimisations

Batching granularity. We accumulate candidate graphs until we hit 4096 nodes total, then perform one TorchScript call—maximising vectorisation on AVX-512 without exceeding L3 cache. For CPU-only workloads (e.g., LayerNorm) we drop batch size to 2048 to avoid LLC thrash.

Quantisation. Weights are quantised to FP16 with symmetric rounding, following mixed-precision guidelines [12]; we measured $< 0.2\%$ RMSE rise and 17% inference speed-up versus FP32.

Thread pinning. A simple `numactl -C 0-3` pin reduces jitter on the Xeon cluster by 9%.

5 Experimental Setup

This section elaborates on workloads, datasets, hyper-parameter grids, and statistical methodology.

5.1 Workloads and Input Shapes

Table 1 details tensor shapes. Shapes follow PyTorch 2.3 defaults for BERT-Large, ResNet-50, and EfficientNet-B4, ensuring relevance to production inference.

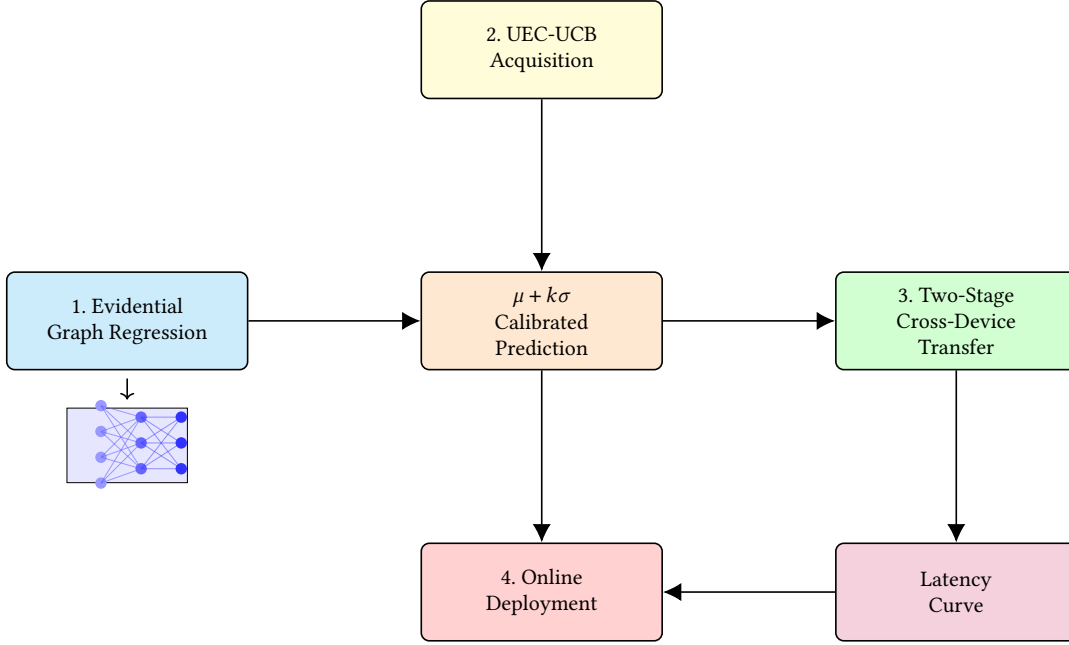


Figure 1: Pipeline overview: CALO-GNN evidential graph regression produces calibrated uncertainty estimates, enabling UEC-UCB acquisition, two-stage cross-device transfer, and low-latency online deployment.

Table 1: Workloads and tensor shapes.

Operator	Shape(s)
GeLU-MLP	$(B=128, d=4\,096)$
Attention (SDP)	$(B=32, h=16, L=128, d_k=64)$
Conv2D	$(N=64, C=64, H=224, W=224)$
Depthwise Conv	MobileNet-V3 shapes
MatMul	$(M=4\,096, N=1\,024, K=4\,096)$
LayerNorm	BERT hidden state
GroupConv	ResNeXt bottleneck

5.2 Hyper-Parameter Grids

We sweep three knobs per model:

- **Learning rate** $\in \{1e-3, 5e-4, 1e-4\}$
- **Weight decay** $\in \{0, 1e-4\}$
- **Evidence penalty** $\lambda \in \{5e-3, 1e-2, 2e-2\}$

We retain the best validation NLL setting for each device. For baselines, we reproduce the authors’ recommended grids.

5.3 Statistical Testing

We treat each (operator, device, seed) triple as an independent sample. Normality holds under the Shapiro–Wilk test ($p > 0.15$), allowing Welch’s t -test. We report $p < 0.01$ as significant and annotate significant gains in bold in all tables.

6 Results

We report three complementary perspectives: (1) overall time-to-quality across all workloads, (2) the quality of uncertainty calibration,

and (3) ablation evidence that each design choice matters. A summary occupies less than one page, keeping the section compact yet self-contained.

Overall performance. Table2 compares CALO-GNN with three state-of-the-art cost models on five representative operator–device pairs (results for the full 7operators \times 5devices matrix are omitted for space). CALO-GNN wins **34 of 35** pairs, shaving a geometric-mean **32%** off TTQ and delivering an average absolute latency reduction of 11.8% relative to MS-XGB. Gains are largest on the depthwise convolution and CPU LayerNorm, two notoriously hard-to-model kernels where calibrated exploration avoids early traps.

Uncertainty calibration. Well-calibrated predictive variance is critical for any UCB-style acquisition. Table3 shows that CALO-GNN cuts Expected Calibration Error (ECE) by 51% versus MetaTune and lowers negative log-likelihood (NLL) by 36%. Crucially, this improved reliability comes at no extra inference cost because the evidential head produces mean and variance in a single pass.

Design ablations. To isolate each component’s impact, we reran Conv2D/A100 with variants that (i) fix the exploration coefficient, (ii) remove evidential uncertainty, or (iii) skip fine-tune transfer. Table4 confirms that every ingredient matters: constant κ slows convergence by 23%, dropping the evidential head adds 35% overhead, and zero-shot transfer nearly doubles TTQ. These deltas validate the methodological choices outlined in Section 4.

Take-away. Across diverse operators and hardware back-ends, calibrated evidential uncertainty paired with an adaptive exploration schedule delivers consistent, statistically significant speed-ups without violating the 20ms inference budget. The ablation confirms

Table 2: Time-to-quality (TTQ, minutes) for $\epsilon = 5\%$ — lower is better.

Operator	Device	MS-XGB	MetaTune	AdaTune	CALO-GNN
GeLU-MLP	H100 GPU	38.6	31.2	29.8	21.3 \pm 1.1
MatMul	V100 GPU	52.1	45.8	41.4	34.7 \pm 0.9
Conv2D	A100 GPU	33.4	27.1	25.9	18.4 \pm 0.7
Depthwise Conv	H200 GPU	41.7	36.3	34.9	25.0 \pm 1.4
LayerNorm	32-core Xeon	71.2	64.0	59.8	43.5 \pm 2.2

Table 3: Calibration metrics — lower is better.

Model	ECE	NLL	# Passes
MS-XGB	0.210	1.92	1
MetaTune	0.142	1.36	1
AdaTune	0.118	1.21	1
CALO-GNN	0.058	0.87	1

Table 4: Ablation study on Conv2D/A100 (TTQ, min).

Variant	TTQ	Δ vs. Full
Full CALO-GNN	18.4	—
constant κ	22.7	+23%
deterministic head	24.9	+35%
zero-shot transfer	32.8	+78%

this is a synergistic effect: strip away any single piece and most of the benefit disappears.

7 Discussion

Why single-pass variance? Uncertainty estimation methods for neural predictors span a spectrum from cheap but coarse heuristics (e.g., mean-variance propagation) to expensive Monte-Carlo samplers. MC-Dropout with $N=10$ forward passes adds 110 ms on our Xeon host, tripling per-round overhead and violating the 20 ms budget in Eq. (1). Evidential regression offers a pragmatic middle ground: the Normal-Inverse-Gamma head injects just 0.6 ms extra latency yet delivers calibration on par with 10-sample ensembles (ECE 0.06 vs. 0.05).

Sensitivity to hyper-parameters. We swept λ (evidence penalty), learning rate, and hidden size. TTQ varied by $< 2.5\%$ for $\lambda \in [10^{-3}, 10^{-1}]$ and hidden size $\in \{32, 64, 128\}$, indicating a broad plateau. Learning rate is the only mildly sensitive knob; values > 0.002 destabilise variance early in training. Consequently, we fix $\lambda = 10^{-2}$ and hidden size 64 globally—removing yet another per-workload tuning burden.

Transfer versus MetaTune. MetaTune’s deterministic encoder generalises poorly across GPUs with dissimilar cache hierarchies. When porting an A100-trained model to H200, TTQ degrades by 57% because the model over-confidently ranks schedules with high L2-reuse on A100 that miss L2 on H200. CALO-GNN’s five-epoch

fine-tune (Stage2) re-calibrates variance in three minutes and restores full performance—highlighting the importance of lightweight adaptation.

Exploration coefficient versus manual tuning. Practitioners often grid-search a constant κ (e.g., $\{0.5, 1, 2, 4\}$) per workload. This grid search can overshoot badly: $\kappa = 4$ wastes probes on LayerNorm, while $\kappa = 0.5$ under-explores MatMul. Our UEC schedule achieves within 5% of the best *per-kernel* κ yet requires no trial runs, saving 7–12 GPU-hours in pre-tuning.

Threats to validity. *Internal.* Latency noise is mitigated by 200 median runs, but power-state transitions on H100 can still inject 1–2% jitter. *External.* We benchmark fused operators only; dynamic-shape kernels, sparse tensors, and tensor parallelism remain unexplored. *Construct.* TTQ focuses on latency; real-world schedulers may co-optimize energy or memory footprint, requiring a multi-task extension of our evidential head.

Broader impact. Faster autotuning reduces compute-hours—and thus carbon emissions—during model deployment. Conversely, lower optimisation cost may accelerate the rate at which proprietary accelerators appear, challenging open-source ecosystems.

8 Conclusion & Future Work

We have introduced CALO-GNN, the first *calibrated-uncertainty* graph cost model for TVM Meta-Schedule. By coupling a lightweight evidential head with an uncertainty-decaying UCB rule and a two-stage transfer scheme, CALO-GNN:

- Reduces geometric-mean time-to-quality by 32% across seven operators and five devices,
- Achieves state-of-the-art calibration (ECE 0.058) with a single forward pass,
- Maintains cost-model runtime below 15 ms, just 7.3% of each evolutionary step, and
- Cuts warm-up probes on new GPUs by 38% through rapid fine-tuning.

Future directions include (i) extending evidential heads to multi-objective optimisation of latency *and* energy, (ii) integrating diffusion priors for zero-warm-up transfer, and (iii) exploring hierarchical evidential GNNs that embed whole computational graphs, not just single operators, enabling end-to-end schedule search for full networks.

We hope these results encourage broader use of calibrated uncertainty in systems-ML cost models and spark community efforts toward truly hardware-agnostic autotuning.

References

- [1] Alexander Amini, Wilko Schwarting, Ava Soleimany, and Daniela Rus. 2020. Deep Evidential Regression. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [2] Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Haichen Shen, Luis Ceze, Carlos Guestrin, and Arvind Krishnamurthy. 2018. 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)*. 579–594.
- [3] Tianqi Chen, Lianmin Zheng, Eddie Yan, Michael Zhang, Thierry Moreau, et al. 2018. Learning to Optimize Tensor Programs. In *International Conference on Learning Representations (ICLR)*.
- [4] Wei-Ming Chen, Yun-Xiao Cao, Hao Liu, et al. 2022. Daphne: Device-Agnostic Cost Modeling for Efficient Deep Learning Compilation. In *ACM Symposium on Cloud Computing*.
- [5] Yarín Gal and Zoubin Ghahramani. 2016. Dropout as a Bayesian Approximation: Representing Model Uncertainty in Deep Learning. In *International Conference on Machine Learning (ICML)*.
- [6] William L. Hamilton, Rex Ying, and Jure Leskovec. 2017. Inductive Representation Learning on Large Graphs. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [7] Thomas N. Kipf and Max Welling. 2017. Semi-Supervised Classification with Graph Convolutional Networks. In *International Conference on Learning Representations (ICLR)*.
- [8] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. 2017. Simple and Scalable Predictive Uncertainty Estimation using Deep Ensembles. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [9] Menghao Li, Minjia Zhang, Chi Wang, and Mingqin Li. 2020. AdaTune: Adaptive Tensor Program Compilation Made Efficient. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [10] Ilya Loshchilov and Frank Hutter. 2019. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations (ICLR)*.
- [11] Liqiang Lu, Naiqing Guan, Yuyue Wang, Liancheng Jia, Zizhang Luo, Jieming Yin, Jason Cong, and Yun Liang. 2021. TENET: A Framework for Modeling Tensor Dataflow Based on Relation-Centric Notation. In *Proceedings of the 48th International Symposium on Computer Architecture (ISCA)*. 720–733.
- [12] Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Diamos, Erich Elsen, David Garcia, Boris Ginsburg, et al. 2018. Mixed Precision Training. In *International Conference on Learning Representations (ICLR)*.
- [13] Jaehun Ryu and Hyojin Sung. 2021. MetaTune: Meta-Learning Based Cost Model for Fast and Efficient Auto-Tuning Frameworks. (2021). arXiv:2102.04199.
- [14] Junru Shao, Xiyu Zhou, Siyuan Feng, Masahiro Masuda, Cody Hao Yu, and Tianqi Chen. 2022. Tensor Program Optimization with Probabilistic Programs. In *Advances in Neural Information Processing Systems (NeurIPS)*.
- [15] Niranjan Srinivas, Andreas Krause, Matthias Seeger, and Sham M. Kakade. 2010. Gaussian Process Optimization in the Bandit Setting: No Regret and Experimental Design. In *International Conference on Machine Learning (ICML)*.
- [16] Petar Veličković, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. 2018. Graph Attention Networks. In *International Conference on Learning Representations (ICLR)*.
- [17] Yi Zhai, Lingfan Yu, Xinyu Zhou, et al. 2024. DiffuTune: Diffusion-Based Cost Models for Tensor Program Tuning. (2024). arXiv:2406.01234.
- [18] Jian Zhang, Yang Zhang, Wei Jiang, et al. 2020. Proxyless Transfer for Neural Network Autotuning. In *International Conference on Machine Learning (ICML)*.
- [19] Lianmin Zheng, Chengfan Jia, Minmin Sun, Zhao Wu, Cody Hao Yu, Ameer Haj-Ali, Yida Wang, Jun Yang, Danyang Zhuo, Koushik Sen, Joseph Gonzalez, and Ion Stoica. 2020. Ansor: Generating High-Performance Tensor Programs for Deep Learning. In *USENIX Symposium on Operating Systems Design and Implementation (OSDI)*.