# ON ACCELERATING EDGE AI: OPTIMIZING RESOURCE-CONSTRAINED ENVIRONMENTS

A PREPRINT

**Jacob Sander**
Department of Computer Science
University of West Florida
Pensacola, FL, USA
jhs39@students.uwf.edu

**Achraf Cohen**
Department of Mathematics and Statistics
University of West Florida
Pensacola, FL, USA
acohen@uwf.edu

**Venkat R. Dasari**
DEVCOM Army Research Laboratory
Aberdeen Proving Ground, MD, USA
venkateswara.r.dasari.civ@army.mil

**Brent Venable**
Department of Intelligent Systems and Robotics
University of West Florida
Institute for Human & Machine Cognition
Pensacola, FL, USA
bvenable@uwf.edu

**Brian Jalaian***
Department of Computer Science
Department of Intelligent Systems and Robotics
University of West Florida
Institute for Human & Machine Cognition
Pensacola, FL, USA
bjalaian@uwf.edu

January 30, 2025

## ABSTRACT

Resource-constrained edge deployments demand AI solutions that balance high performance with stringent compute, memory, and energy limitations. In this survey, we present a comprehensive overview of the primary strategies for accelerating deep learning models under such constraints. First, we examine *model compression* techniques-pruning, quantization, tensor decomposition, and knowledge distillation-that streamline large models into smaller, faster, and more efficient variants. Next, we explore *Neural Architecture Search (NAS)*, a class of automated methods that discover architectures inherently optimized for particular tasks and hardware budgets. We then discuss *compiler and deployment frameworks*, such as TVM, TensorRT, and OpenVINO, which provide hardware-tailored optimizations at inference time. By integrating these three pillars into unified pipelines, practitioners can achieve multi-objective goals, including latency reduction, memory savings, and energy efficiency-all while maintaining competitive accuracy. We also highlight emerging frontiers in *hierarchical NAS*, *neurosymbolic approaches*, and *advanced distillation* tailored to large language models, underscoring open challenges like pre-training pruning for massive networks. Our survey offers practical insights, identifies current research gaps, and outlines promising directions for building scalable, platform-independent frameworks to accelerate deep learning models at the edge.

*Keywords* Edge AI · Optimization · Resource-Constrained Environments · Neural Architecture Search · Model Compression

# 1 Introduction

Integrating Artificial Intelligence (AI) models into tactical edge computing systems is fundamentally transforming military operations and evolving battlefield technology over time. In practical applications of tactical edge computing, AI models are deployed on constrained platforms that have significant limitations in processing power, memory, and communication capabilities. These constraints can hinder the use of complex AI models for real-time tasks such as surveillance, enemy identification, predictive simulations, and navigation.

The goal is to develop intelligent and adaptive optimization strategies that can overcome these resource limitations while enhancing the performance of AI models without significantly compromising their accuracy. This is particularly critical in military contexts, where timely and precise decision-making can be pivotal to the success of a mission. Therefore, the ability to effectively utilize AI models in tactical edge computing environments is crucial for military applications. This capability allows for the creation of intelligent systems and machines that can operate in harsh conditions with limited network access, ultimately providing invaluable support to troops in the field.
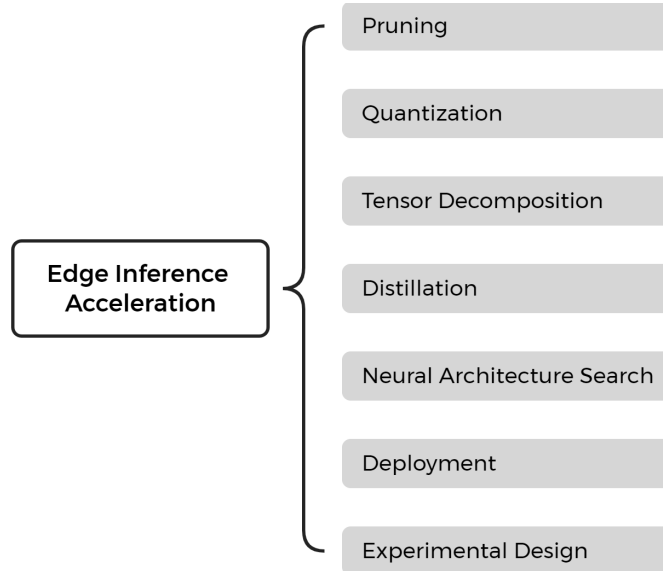


Figure 1: Strategies and Tools for Enhancing AI Edge Computing

Our survey aims to identify emerging approaches for optimizing Artificial Intelligence (AI) models to accelerate model inference. We focus on strategies that meet both model-specific and platform-specific constraints while maintaining accuracy. The survey will examine techniques such as model pruning, quantization, knowledge distillation, neural architecture search, and hardware-aware design. Fig. 1 presents a map of strategies and tools for edge computing acceleration. By exploring the latest advancements in AI model optimization, we hope to uncover promising methods that address the challenges of deploying AI models in resource-constrained environments. Ultimately, our goal is to contribute to developing more efficient, effective, and reliable AI-powered systems for critical applications.

We categorize optimization strategies into two main families, although in practice, scientists and engineers often combine approaches to enhance performance. First, we will examine *Model Compression* techniques, which focus on transforming a trained model into a smaller version that still maintains high performance. However, Model Compression techniques lack performance guarantees, leading researchers to seek higher-performing and more innovative architectures. Consequently, we will also address *Neural Architecture Search* (NAS) and Hyper-Parameter Optimization (HPO) methods that explore the configuration space to identify more efficient architectures based on defined objectives. Although these methods are formally separate from the aforementioned optimizations, we will discuss compile-time techniques that reduce model inference time during deployment. Finally, we will provide a detailed overview of experiments involving composite methods, which integrate Model Compression techniques with customized objectives and constraints to explore the NAS and HPO space.

This paper is organized as follows: Section 2 introduces model compression techniques and their characteristics, such as quantization, pruning, and knowledge distillation; Section 3 presents neural architecture search; Section 4 discusses the deployment aspects related to compilers and hardware. Section 5 discusses case studies and integrated approaches. Finally, in Section 6, conclusions and some possible research directions are presented.

## 2    Model Compression Techniques

AI models, especially deep learning ones, tend to be over-parameterized. This over-parameterization is crucial, making model compression feasible and highly effective. Model compression refers to reducing the size and computational complexity while preserving performance (e.g., accuracy).

In tactical edge computing environments, the size and computational demands of AI models can hinder their deployment on resource-constrained platforms. Model compression techniques aim to reduce the size of larger AI models and lower computational requirements during inference, all without significantly compromising performance. Compressing a model enables the use of advanced AI capabilities in tactical edge computing environments with limited computational power, memory, and storage.

### 2.1    Quantization

Quantization is the process of reducing the precision of weights and activations, often stored as 32-bit floating-point values, to lower bit representations such as 8-bit or binary, to optimize memory and computational efficiency Jacob et al. [2018]. Recent research on quantization in large language models (LLMs) highlights its importance for reducing model size and improving efficiency. Quantization techniques include post-training quantization (PTQ) and quantization-aware training (QAT) Chen et al. [2024a], with algorithms like LLM-QAT and SmoothQuant addressing challenges such as outliers and activation quantization Wang et al. [2024]. SmoothQuant addresses outliers in activation values by redistributing them to reduce their impact during quantization, enabling stable and efficient low-bit representations. Post-training quantization (PTQ) modifies the model's parameters after training without additional data or computational overhead. In contrast, quantization-aware training (QAT) incorporates quantization during training, allowing the model to adapt to the reduced precision, often yielding better performance. LLM-QAT is specifically tailored for large language models, adapting quantization-aware training techniques to handle the unique challenges posed by their massive parameter scales and diverse activation ranges. CVXQ leverages convex optimization techniques to enable highly flexible and efficient post-training quantization, particularly for extremely large models with billions of parametersYoung [2024]. Activation-Aware Weight Quantization (AWQ) Lin et al. [2024a], EfficientQAT Chen et al. [2024a], and QLoRA Dettmers et al. [2023] all present alternatives to training a full-precision model from scratch to implement quantization-aware models. For example, AWQLin et al. [2024a] introduces a novel mechanism to optimize weight quantization based on activation patterns, ensuring high accuracy even under aggressive quantization regimes. Comparative studies reveal that the optimal quantization format (integer or floating-point) varies across layers, leading to the proposal of Mixture of Formats Quantization (MoFQ) for improved performance in both weight-only and weight-activation scenarios Zhang et al. [2024]. The RPTQ employs a channel-wise rearrangement and clustering strategy to manage activation range variations, enabling robust 3-bit activation quantization for large language modelsYuan et al. [2023]. Outlier Suppression+ innovates by introducing channel-wise shifting and scaling, effectively addressing asymmetric outliers and achieving near-floating-point performance for both small and large modelsWei et al. [2023]. These techniques can significantly reduce memory consumption, with OPT-175B (Meta's Open Pre-trained Transformer) quantization potentially leading to an 80% reduction Yuan et al. [2023].

Mixed precision training has emerged as an effective technique for improving the efficiency of large-scale AI models. Researchers introduced a method using half-precision floating-point numbers for weights, activations, and gradients, while maintaining a single-precision copy of weights to accumulate gradients Micikevicius et al. [2017]. Li et al. Li et al. [2023a] proposed a layered mixed-precision approach, adjusting training precisions for each layer based on its contribution to the training effect. The Channel-Wise Mixed-Precision Quantization (CMPQ) was developed Chen et al. [2024b], allocating quantization precision in a channel-wise pattern for large language models (OPT-2.7B and OPT-6.7B). Other research explored mixed precision low-bit quantization for neural network language models in speech recognition Xu et al. [2021], using techniques such as KL-divergence, Hessian trace weighted quantization perturbation, and mixed precision neural architecture search. These methods have shown notable enhancements in training speed, memory efficiency, and model compression, all while preserving performance in various deep learning applications.

Applying weight and activation quantization together improves inference speed and model memory metrics, especially for quantization-aware training Menghani [2023], AI [2021]. See Table 1.

In certain situations, hardware memory constraints can affect inference performance. Using reduced-size weight representations can help lower latency. However, this process introduces a quantization error, which can be significant. Various algorithms for quantization and de-quantization have been developed to minimize the impact of this error Menghani [2023].

Weight and activation quantization can be applied before or after training; for example, authors have taken large, pre-trained models and used a quantization policy to their weights and activations AI [2021]. This is an example of Post-

| Model | Accuracy (Original) | Accuracy (PTQ) | Accuracy (QAT) | Latency (Original) | Latency (PTQ) | Latency (QAT) | Original Size (MB) | Optimized Size (MB) |
|---|---|---|---|---|---|---|---|---|
| Mobilenet-v1-1-224 | .709 | .657 | .70 | 124 | 112 | 64 | 16.9 | 4.3 |
| Mobilenet-v2-1-224 | .719 | .637 | .709 | 89 | 98 | 54 | 14 | 3.6 |
| Inception v3 | .78 | .772 | .775 | 1130 | 845 | 543 | 95.7 | 23.9 |
| Resnet v2 101 | .770 | .768 | N/A | 3973 | 2868 | N/A | 178.3 | 44.9 |

Table 1: Comparison of original, post-training quantized (PTQ), and quantization-aware training (QAT) in terms of accuracy, latency (ms), and size as reported in AI [2021].

Training Quantization. This leads to some amount of quantization error and can occasionally even increase inference time, especially if weights and activations are dynamically quantized and de-quantized (see Mobilenet-v2-1-224 in Table 1 for an example of this effect).

These advancements in quantization techniques are crucial for deploying LLMs on resource-constrained devices, reducing computational costs, and mitigating the environmental impact of large-scale AI systems Lang et al. [2024]. Quantization can be applied before (QAT) or after training (PTQ). The strength of PTQ methods is that they do not require further data or training; a model can be trained, quantized, and deployed in a streamlined pipeline Menghani [2023]. This is of the most significant benefit when considering large models that are impractical to re-train. On the other hand, QAT methods Jacob et al. [2018] generally result in better performance, as the model can adjust to the introduced quantization error but at the cost of additional training. This may be tractable for small models, but it is impractical for large transformer models with prohibitive training costs.

## 2.2 Pruning

Pruning is a model compression technique that reduces the size of a neural network by removing unnecessary connections, aiming to improve computational efficiency while maintaining performance. It is well-understood that neural nets are over-parameterized to guarantee a path to a reasonable local loss minimum during training Choromanska et al. [2014]. By dropping low-importance neurons or channels from the model, we can preserve high performance while minimizing the model's size and, thus, the resulting compute and energy costs.

Pruning techniques can be categorized by their timing (pre-training, during training, post-training), granularity (structured, unstructured), and decision criteria (rule-based, learning-based) as shown in Fig. 2. Pruning methods can be
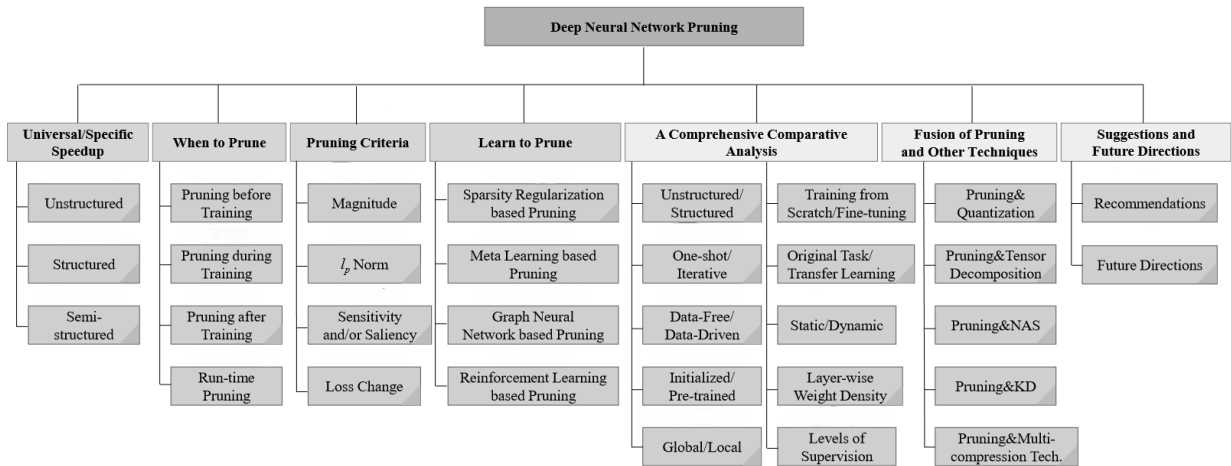


Figure 2: Pruning Taxonomy as formulated by Cheng et al. [2023]

classified based on their structural impact on deep neural networks (DNNs). Unstructured, or weight-wise pruning, is a fine-grained approach where individual weights in the network are pruned Lee et al. [2019]. This method is efficient for small networks, as it reduces the overall number of weights by applying binary masks to the least essential parameters Wang et al. [2020a]. In larger networks like LLMs, authors may merely set pruned weights to 0 instead of maintaining a separate binary mask for all parameters Frantar and Alistarh [2023]. Regardless, to realize the inference speed gains of a pruned network, specialized hardware or deployment-time compilers must be used to realize the gains from this fine-grained sparsification Cheng et al. [2023]. Structured pruning operates at a coarser level by removing groups of

objects together: channels He et al. [2017], filters You et al. [2019], transformer attention heads Shim et al. [2021], neurons Ashkboos et al. [2024], layers Men et al. [2024], and blocks Ma et al. [2023]. These structured methods offer inference acceleration regardless of hardware because the sparsification is coarse. If an entire channel is pruned, for example, then at inference, it is ignored; if instead the neurons in that channel are pruned in an unstructured way, the channel is still called at inference time, and the compiler and hardware must make sure to skip 0 weight operations efficiently inside the retained channel. While structured pruning offers consistent inference acceleration across hardware, unstructured pruning often achieves higher compression ratios but requires specialized hardware or compiler optimizations to realize speedups.

Structured pruning maximizes speed improvements by carefully selecting architectural elements to prune in a way that minimizes performance degradation. Semi-structured pruning is a catch-all category for many methods that blend elements of structured and unstructured pruning Xu et al. [2024a], Ma et al. [2020], Meng et al. [2020]. By sequentially performing coarse pruning and then fine-grained pruning of the remaining structures, authors achieve greater compression levels while retaining performance.

As we mentioned, pruning seeks to decrease neural networks' size and computational expenses by eliminating weights or structural elements that minimally affect overall performance. Following this, we will introduce various pruning strategies—which differ mainly by the timing of their application—and emphasize significant research gaps.

**Pre-Training Pruning:** One forward-looking method is to prune networks based on randomly initialized weights Wang et al. [2020a]. Though appealing for saving training time (because pruned weights require no subsequent computation), this approach risks issues like layer collapse Tanaka et al. [2020]. It has been applied to convolutional neural nets (CNNs) Lee et al. [2019, 2020], but scaling it up to large models remains challenging, given the expense of training even pruned networks.

**Pruning During Training:** Another strategy embeds pruning into the training loop, iteratively updating which weights remain active. For instance, the RigL algorithm Evci et al. [2021] periodically removes and regrows weights to maintain model capacity. Structured sparsity learning (SSL) Wen et al. [2016], network slimming Liu et al. [2017], and differentiable methods like Differential Sparsity Allocation (DSA) Ning et al. [2020] similarly integrate pruning decisions with gradient-based updates. In neural architecture search, a related concept called Progressive Shrinking alternates between pruning and training to explore potential architectures Wang et al. [2020b]. However, due to high computational costs, these techniques see limited exploration for large-scale models.

**Pruning After Training:** Pruning after training is widely used thanks to the abundance of large pre-trained models, which can be pruned and then fine-tuned on downstream tasks. Several algorithms focus on mitigating accuracy loss without retraining Frantar and Alistarh [2023], Ashkboos et al. [2024], Kwon et al. [2022], while others adopt a pipeline of pruning plus fine-tuning on a smaller dataset Liu et al. [2021], Ma et al. [2023]. The Lottery Ticket Hypothesis Frankle and Carbin [2019] is particularly influential in this space: it suggests dense, trained networks harbor sub-networks ("winning tickets") that can match original performance at a fraction of the size. Research has confirmed these ideas in CNNs You et al. [2022] and transformer-based language models Chen et al. [2021], including variations that prune large networks post-training (e.g., APQ Wang et al. [2020b]) without full retraining. When retraining is needed but cost-prohibitive, lightweight fine-tuning schemes (e.g., LoRA) have shown promise for recovering performance Ma et al. [2023].

**Inference Time Pruning:** Inference time pruning (or dynamic pruning) adjusts the network *per input* Rao et al. [2019], bypassing full computation for simpler samples. By pruning based on input complexity Tang et al. [2021], resource usage can be reduced considerably without heavily compromising performance.

**Criteria for Pruning:** Different pruning criteria guide which weights or modules to remove. Magnitude-based methods prune weights below a certain threshold Han et al. [2016], while norm-based methods discard entire channels (e.g., using $L1$ norm Li et al. [2017], He et al. [2017]). Alternative criteria include channel saliency Zhao et al. [2019], neuron sensitivity Santacroce et al. [2023], and module relevance Dery et al. [2024]. Likewise, some work leverages structural graphs to locate highly connected components for more efficient pruning Ma et al. [2023], Zhang et al. [2021]. Researchers have also introduced reinforcement learning to learn pruning policies He et al. [2018], Bencsik and Szemenyei [2022] automatically.

Selecting a pruning strategy depends on model size, the target deployment scenario, and the available computational budget. Notably, while pruning yields smaller networks, it does not always guarantee proportional speedups on certain hardware backends. To address these gaps, a more quantitative outlook-comparing pruning to alternatives like quantization or tensor decomposition across metrics such as inference speed, memory footprint, and accuracy trade-offs-is needed. Such evaluations would be especially insightful for large-scale models that remain prohibitively expensive to retrain or fine-tune.

## 2.3   Tensor Decomposition

Tensor decomposition is a recognized model compression method that approximates a high-rank weight tensor using products of lower-rank factors. This technique decreases the parameter count and computation costs associated with a neural network. As shown in Fig. 3, substantial weight matrices and convolutional kernels can be broken down into smaller parts that effectively perform the same function, significantly reducing memory usage.
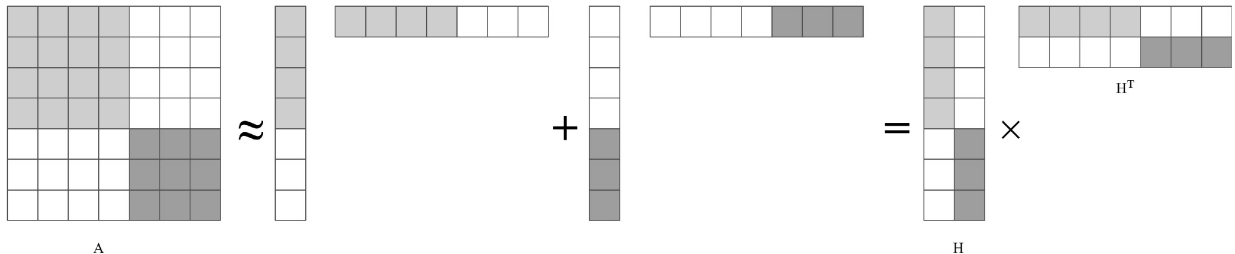


Figure 3: Kernel approximation of low-rank decomposition matrix as depicted in Li et al. [2023b]

Historically, tensor decomposition was first applied to fully connected layers, factorizing large weight matrices into the product of two much smaller matrices. Subsequent research extended this idea to convolutional neural networks by decomposing large convolutional kernels into equivalent operations that increase the number of channels while reducing overall computation Li et al. [2023b]. Various decomposition strategies, such as CP, Tucker, and Tensor-Train Decomposition (TTD), provide different ways to balance representation capacity and resource efficiency.

In the context of foundation models, tensor decomposition shows particular promise. Because these models typically contain billions of parameters, decomposition can dramatically reduce memory usage and computational overhead, especially when paired with a short period of domain-specific fine-tuning Hajimolahoseini et al. [2021], Saha et al. [2024]. Recent research has focused on compressing large language models (LLMs) using tensor decomposition techniques. TensorGPT applies Tensor-Train Decomposition (TTD) to GPT-2 to compress token embeddings, achieving up to 65.64x compression ratios while maintaining comparable performance to the original model Xu et al. [2023a]. MoDeGPT introduces a modular decomposition framework that partitions Transformer blocks into matrix pairs. It applies various matrix decomposition algorithms, achieving 25-30% compression rates while maintaining 90-95% zero-shot performance on LLAMA-2/3 and OPT models Lin et al. [2024b]. These methods significantly reduce LLM size and computational requirements without substantial performance loss.

However, implementing decomposition-based compression on real-world hardware still poses challenges, including needing specialized kernels to maintain throughput gains. Future directions in tensor decomposition include automated rank selection, hybrid approaches that integrate decomposition with other compression methods (e.g., quantization and pruning), and rigorous exploration of decomposition within advanced architectures like Transformers.

## 2.4   Knowledge Distillation

Distillation is a model compression technique that transfers knowledge from a large, complex model (teacher) to a smaller, efficient model (student) with minimal loss in performance, reducing resource costs and enabling deployment in resource-constrained environments Hinton et al. [2015]. Fig. 4 referees to distillation components and approaches.

The two essential steps in knowledge distillation (KD) are knowledge elicitation and student distillation Xu et al. [2024b]. The initial step involves extracting knowledge from the teacher model. For instance, the teacher can be given input data to produce the corresponding output logits Hinton et al. [2015]. White-box methods, like feature extraction, provide direct access to the intermediate activations of the teacher model. In contrast, black-box techniques only utilize the output logits from both the teacher and student models Xu et al. [2024b]. Furthermore, in the context of foundation models, other knowledge elicitation strategies can be employed; the teacher might label a dataset for training the student, generate synthetic data through input expansion (creating new inputs based on initial seed data), or use data curation, where teacher feedback is applied to refine the dataset. After eliciting the knowledge, the student model is trained to mimic or approximate the teacher's knowledge. This training typically involves a loss function that minimizes the divergence between student and teacher logits but can also incorporate supervised fine-tuning (which maximizes the likelihood of sequences generated by the teacher) and reinforcement learning, where teacher feedback enhances the student's performance. These strategies aim to narrow the gap between the outputs of the student and teacher, facilitating effective knowledge transfer Xu et al. [2024b].
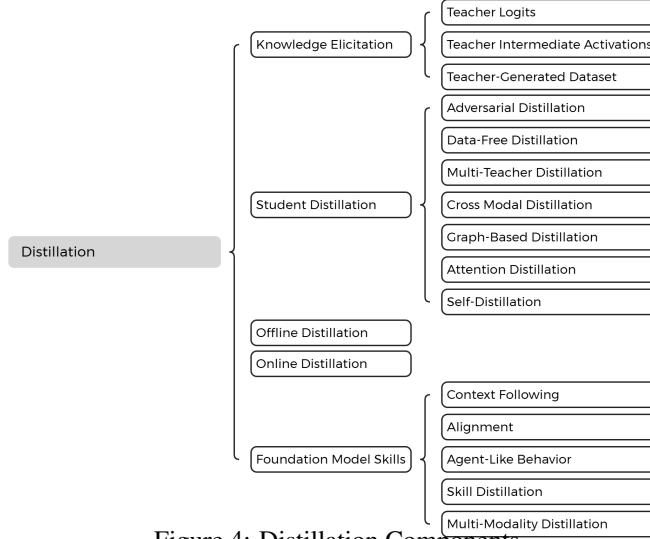
Figure 4: Distillation Components

Numerous knowledge distillation methods have been established, each with design considerations, techniques, and uses. **Adversarial Distillation** highlights that the compact model produced through distillation may be susceptible to adversarial threats Goldblum et al. [2020] By employing the Generative Adversarial Network (GAN) framework, Goodfellow et al. [2014] researchers have created several training pathways that include Adversary, Teacher, and Student networks to enhance the final student's robustness Jung et al. [2024]. Additionally, GANs can be utilized for **Data-Free Distillation**, where the adversarial network produces synthetic examples, allowing the student to be trained without requiring extra data Chen et al. [2019a]. Although adversarial distillation bolsters robustness, it frequently raises computational demands and training complexity, which may hinder scalability.

In **Multi-Teacher Distillation**, multiple teacher models distill their knowledge into a single student model. The knowledge is usually represented as logits or features, providing a comprehensive and balanced distillation from diverse sources Hinton et al. [2015], Zhang et al. [2017]. Authors have developed alternative approaches by combining multi-teacher methods with Bayesian Neural Nets. For example, Vadera et al. show that a Bayesian Neural Net can generate an ensemble of teacher networks; this ensemble can then be used to train a student network as a multi-teacher distillation, with measurable benefits to the student's uncertainty calibration Vadera et al. [2020]. **Cross-Modal Distillation** involves transferring knowledge from a teacher model trained on one modality to a student model designed to work with a different modality. Authors propose that students deployed on small or expensive-to-label datasets may benefit from this kind of transfer learning, as a large teacher model in one modality may incorporate the knowledge that improves the performance of the student in another modality Gupta et al. [2015]. Recent advancements include combining iterations of self-distillation with cross-modal data for regularization Andonian et al. [2022], and filtering out cross-modal samples during distillation that do not add information to the student modality Huo et al. [2024]. In **Graph-Based Knowledge Distillation**, the data modality is restricted to graphs, and models to Graph Neural Networks (GNNs) Liu et al. [2023a]. In systems that utilize attention mechanisms, like transformers, Vaswani et al. [2023] a distillation loss can be defined using the difference in attention mechanism between student and teacher; this is defined as **Attention-Based Distillation**. It has found use across modalities, including text Wang et al. [2020c] and images Wang et al. [2022a]. The attention-based distillation loss is defined in Eq. (1).

$$\mathcal{L}_{AT} = \frac{1}{A_h|x|} \sum_{a=1}^{A_h} \sum_{t=1}^{|x|} D_{KL}(\mathbf{A}_{\mathbf{L,a,t}}^{\mathbf{T}}||\mathbf{A}_{\mathbf{M,a,t}}^{\mathbf{S}}), \qquad (1)$$

where $|x|$ and $A_h$ represent the sequence length and the number of attention heads; more details can be found in this paper Wang et al. [2020c]. The number of teacher and student network layers are denoted with $L$ and $M$; their last transformer layers are written $A_L^T$ and $A_M^S$.

This is not to be confused with other methods that define an activation-based distillation loss, also sometimes called attention Crowley et al. [2019], See Eq. (2).

$$\mathcal{L}_{AT} = \mathcal{L}_{CE}(y, \sigma(s)) + \beta \sum_{i=1}^{N_L} \left\| \frac{f(A_i^t)}{\|f(A_i^t)\|_2} - \frac{f(A_i^s)}{\|f(A_i^s)\|_2} \right\|_2, \tag{2}$$

where $s$ is the output logits of the student network. The cross-entropy loss, $\mathcal{L}_{CE}$, is added to the attention transfer component, which ensures that the difference between the spatial distributions of the student and teacher activations at selected layers in the network, $f(A_i^t)$, and $f(A_i^s)$, is minimized. This can be viewed as ensuring that the student network pays attention to the same things as the teacher network at those layers.

In addition to base model type and data modality, the distillation training task can also be defined in different ways. **Offline Distillation** is the conventional approach where the teacher model is pre-trained, and its knowledge is then used to guide the training of a smaller student model. Hinton et al. [2015] By contrast, **Online Distillation** occurs simultaneously, where both teacher and student models are trained together. This approach is appropriate when a high-capacity teacher model is unavailable, and knowledge needs to be distilled on the fly Gou et al. [2021]. This idea can even be extended to teacher-free frameworks, where ensembles of student models distill knowledge into each other simultaneously Zhang et al. [2017], Chen et al. [2019b]. **Self-Distillation** involves the same model acting as both teacher and student, iteratively refining its own outputs. Some authors argue that self-distillation serves as a form of regularization, preventing overfitting by softening the output logits Mobahi et al. [2020]. Other researchers utilize self-distillation to reduce the model size by transferring knowledge from deeper layers to shallower layers through an iterative process Zhang et al. [2019].

Given the size and scope of tasks for Foundation Models (FMs), there are numerous additional requirements defined by authors and various formulations of distillation to meet these requirements. In our survey, we adopt the framework proposed by Xu et al. Gou et al. [2021], which encompasses a broad definition of distillation. For instance, if a teacher model labels a dataset used to train a student model, we consider this a form of distillation since the underlying knowledge is still transferred, even if specific distillation losses are not explicitly defined.

**Context Following** refers to the student model's ability to effectively interpret and respond to complex user inputs or contexts. Authors may view self-instruction as a form of self-distillation, where models generate their own input-output examples to further fine-tune themselves for specific tasks Wang et al. [2022b]. This approach can enhance the reasoning capabilities of smaller FMs by leveraging the strengths of larger models. Similarly, student models may select or filter data used by the teacher for fine-tuning, representing another form of self-distillation Li et al. [2024]. Additionally, schemes like ORCA Mukherjee et al. [2023] utilize step-by-step explanations from large foundation teacher models to train smaller models in reasoning through tasks. For conversational models, self-distillation with feedback involves the large teacher model ranking the utility of multiple proposed outputs from the small model, using this feedback to refine the small model's responses Xu et al. [2023b]. Recently, Knowledge-Augmented Reasoning Distillation (KARD) Kang et al. [2023] has been introduced to address the challenge of smaller FMs memorizing vast amounts of data required for expert-level question answering. KARD involves retrieving data from external sources and using it to fine-tune the small model's responses, enhancing both distillation for small FMs and the performance of general-purpose models Asai et al. [2023].

**Alignment** pertains to the qualitative ability of a machine learning algorithm to understand and respond in accordance with human-intuitive values. Recent benchmarks for this quality in FMs include MT-Bench Zheng et al. [2023] and AlpacaEval Dubois et al. [2024]. In this context, distillation into a smaller model has been observed to decrease alignment with user intent Tunstall et al. [2023]. Ongoing efforts aim to develop methods for distilling teacher-model alignment into student models and enhancing the alignment of the resulting students. For example, Anthropic's Constitutional AI et. al. [2022] evaluates responses based on a concise set of human-generated rules to ensure alignment.

**Agent-Like Behavior** involves enhancing the autonomous capabilities of the student model, enabling it to plan and execute actions in response to inputs by utilizing external "tools" (e.g., other models or APIs) to assist in task completion. One approach involves distilling two student models: the first, a 'sub-goal generator', takes the last ten actions of the agent and an overarching task as input to output the current sub-goal. The second student, an 'action generator', takes the current sub-goal, overarching task, and last ten actions to predict the next action Hashemzadeh et al. [2024]. Similar methodologies are found in the Lumos framework Yin et al. [2024a] and the FireAct framework Chen et al. [2023], which utilize a limited number of agent-action trajectories for fine-tuning. Other researchers have focused on enabling FMs to accurately call APIs, minimize hallucinations, and remain robust to distribution shifts over time in API calls by incorporating distillation techniques specific to API-call tasks Patil et al. [2023].

**Skill Distillation** enables the student model to specialize in a particular domain, such as Natural Language Processing (NLP) or computer vision. Training a general-purpose student model using a large teacher across extensive datasets is often infeasible due to time and cost constraints. However, by carefully selecting a smaller, domain-specific data distribution, relevant skills can be effectively transferred to the student. Efforts to streamline the distillation process

include MiniLLM Gu et al. [2024] and DistillLLM Ko et al. [2024]. Alternative methods, such as Impossible Distillation Jung et al. [2023], use paraphrastic proximity to retrieve high-quality paraphrase datasets and distill student models that outperform their GPT-2 teachers, providing both datasets and trained student models. In NLP search and data augmentation tasks, QUILL employs a two-stage distillation process where a 'Professor' model generates long, retrieval-augmented responses. These responses are distilled into a 'Teacher' model without retrieval capabilities, which then annotates an unlabeled dataset for distillation into the final 'Student' model Srinivasan et al. [2022]. Skill distillation can also be performed across different data modalities, similar to non-FM approaches.

**Multi-Modality Distillation** involves transferring knowledge from a teacher model operating in one modality to a student model in a different modality (e.g., visual to textual). This is typically achieved through the teacher model's labeling of a dataset, which the student model then uses for training. Notable examples of such methods include LLaVa Liu et al. [2023b] and Macaw-LLM Lyu et al. [2023].

Model compression addresses researchers' needs for streamlined search processes and rapid deployment. However, there are scenarios where achieving higher performance necessitates more advanced architectures beyond compression. In such cases, it becomes essential to employ methods that can search within an appropriate architectural space to optimize specified objectives while adhering to imposed constraints, which are often dictated by hardware limitations and mission requirements at the edge. This necessity leads to the exploration of **Neural Architecture Search** (NAS), a pivotal area in model optimization that systematically designs and identifies optimal neural network architectures tailored to specific tasks and constraints.

**Pruning vs. Quantization:** Recent studies argue that quantization may, in some cases, be superior to pruning Yin et al. [2024b]. Weights deemed unimportant in one domain might still be crucial for challenging downstream tasks, and quantizing them (rather than removing them) can better preserve performance under resource constraints.

# 3 Neural Architecture Search

Neural Architecture Search (NAS) automates the design of neural network architectures by systematically exploring a search space of possible configurations. As highlighted by He et al. [2021], NAS is a sub-field of Automated Machine Learning (AutoML) specializing in neural architectures. Fig. 5 outlines concepts and tools related to NAS. It involves:

1. **Defining the search space** (e.g., convolutional, recurrent, or fully connected layers),
2. **Selecting a search strategy** (or Architecture Optimization method),
3. **Choosing an evaluation method** to gauge candidate architectures' performance.

Additionally, researchers often distinguish **architecture optimization (AO)** from **hyperparameter optimization (HPO)**: AO concerns layer configurations and connectivity, while HPO adjusts non-architectural factors such as batch size or learning rate.

## 3.1 Defining the Search Space

**Neural Architecture Search** begins with specifying which architectures are valid candidates. For instance, the architecture might be a sequence of convolutional and pooling layers, or a more complex graph of connected modules. This definition greatly impacts both the quality of results and the computational cost of NAS.

**Primitive vs. Composite Components**   Researchers can allow only a set of primitive operations (e.g., "3×3 convolution, ReLU, max pooling"), creating a large, expressive search space Zoph and Le [2017], Cai et al. [2017]. Such expressiveness can discover novel architectures but may demand intensive computation. Conversely, a higher-level "cell-based" or "motif-based" approach restricts the search space to composite building blocks (cells), greatly reducing complexity but potentially missing innovative designs.

**Leveraging Encodings**   An effective way to limit or navigate the search space is by encoding architectures into concise representations:

- **String-Based Encoding**: The network is represented as a sequence of tokens denoting operations and hyperparameters (e.g., `[Conv 32 3x3, ReLU, MaxPool, Dense 128]`) Zoph and Le [2017], Cai et al. [2017].
- **Graph-Based Encoding**: A directed acyclic graph (DAG) describes data flow, with edges as operations and nodes as data tensors (or vice versa). This DAG can be flattened as an adjacency matrix or represented path-by-path (see Fig. 6) White et al. [2020].
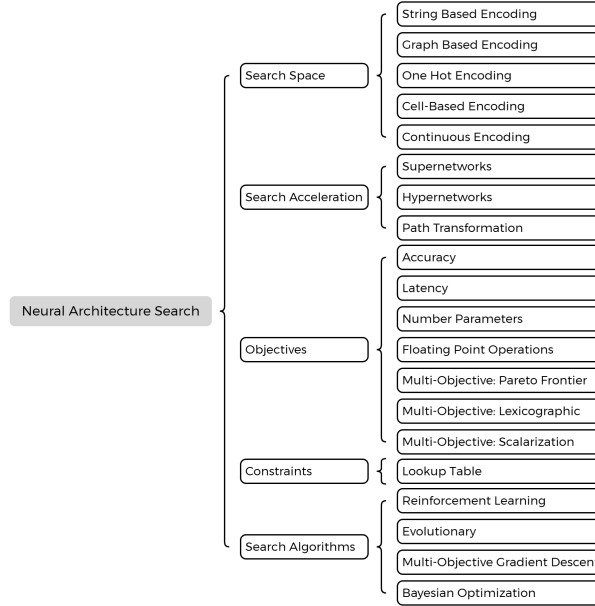
Figure 5: Concepts and Tools in Neural Architecture Search

- **Cell-Based Encoding**: Instead of searching large architectures at once, one can *first* search for a relatively small "cell" (a DAG of operations) and then stack or concatenate multiple copies of this cell to form the full network White et al. [2023]. This hierarchical approach is especially popular for large models, reducing a potentially massive search space to a smaller cell space plus macroscale configuration.
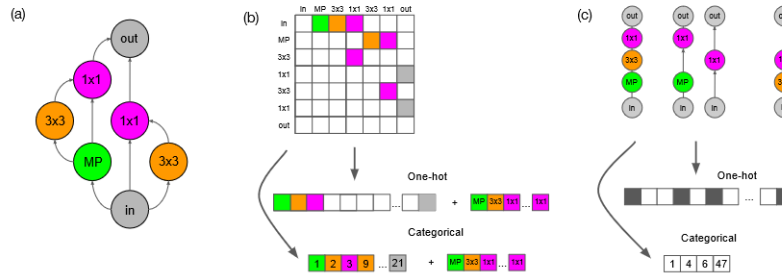


Figure 6: A Study on Encodings for NAS White et al. [2020]. (a) A CNN visualized as a DAG; (b) Adjacency matrix of the DAG with one-hot and categorical encodings; (c) A path-based graph representation.

**Differentiable Architectures**    While many search spaces are discrete, **differentiable NAS** (e.g., DARTS Liu et al. [2018]) continuously parameterizes architectural choices. For a DAG with multiple candidate operations at each edge, differentiable weights determine which operations are active. Training these weights via gradient descent dramatically accelerates the search but can produce suboptimal or unstable results if the relaxed approximation is biased.

**Abstract Architectures**    Recent work notes that many architectures with similar global properties (e.g., total depth or width) often exhibit similar performance Jin et al. [2022], Wan et al. [2021]. Hence, some authors reduce the search space to high-level abstractions, evaluating only one representative network per abstract configuration.

## 3.2   Search Strategies

**Supernetworks    Once-For-All (OFA) training** Cai et al. [2020] is a widely-known supernetwork approach. One trains a large model that conceptually contains multiple sub-networks (created by pruning layers or channels). Each sub-network is then evaluated with minimal additional fine-tuning, drastically cutting the repeated training cost. This

idea is related to **Progressive Shrinking**, where the model is iteratively pruned and fine-tuned to preserve performance Wang et al. [2020b]. However, supernetworks can be impractical for extremely large architectures (e.g., LLMs) because the combined supernetwork might exceed feasible memory or training time Xie et al. [2020]. Moreover, if the supernetwork's weight-sharing scheme is biased, the sub-network performance can be noisy or misleading.

**Hypernetworks**    A **hypernetwork** predicts the weights for candidate architectures, bypassing the need to train each from scratch Li et al. [2020], Liu et al. [2019]. *Graph Hypernetworks (GHNs)* Zhang et al. [2018] (see Fig. 7) extend this to a graph neural network (GNN) that reads an input DAG (the candidate architecture) and outputs all free parameters for that network. Evaluating many architectures thus involves:

1. Feeding each architecture's graph into the GHN;
2. Generating the corresponding weights;
3. Quickly measuring its performance (e.g., on a small validation set).

Since the GHN itself is trained only once, large-scale searches become more computationally tractable.
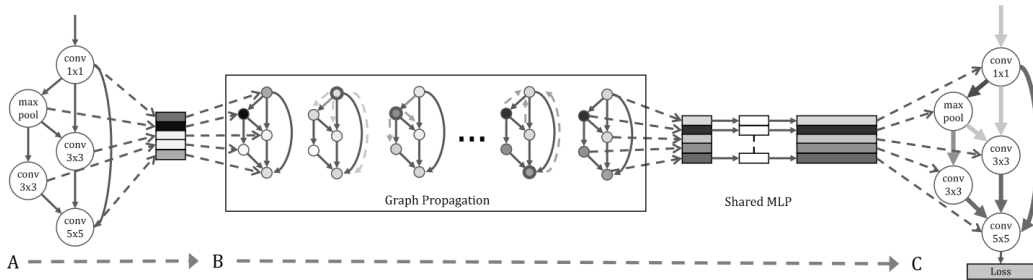


Figure 7: Graph Hypernetwork framework Zhang et al. [2018]. **(A)** A neural network architecture is sampled and input to a GHN. **(B)** The GHN performs graph propagation to generate weights. **(C)** The GHN is trained to minimize the training loss of the generated network, ranking architectures by performance.

**Path-Transformation**    In **path-transformation** Cai et al. [2018], one starts with a large model and gradually transforms it by adding or pruning layers, guided by reinforcement learning (RL) or heuristic rules. Weight sharing between unchanged components avoids retraining from scratch each time. After each transformation, the model is quickly fine-tuned; the search strategy monitors validation metrics to decide the next step.

### 3.3   Objectives in NAS

NAS typically seeks to maximize performance (*e.g.*, accuracy) and/or minimize resource use (*e.g.*, latency, memory). However, fully training *each* candidate to measure accuracy is often impractical with a large search space. Hence, authors have proposed multiple *surrogates* or *predictors* to approximate accuracy.

**Accuracy Predictors and Early Stopping**    Instead of full training, an **accuracy predictor** Wang et al. [2020b], Cai et al. [2020] can take an architecture encoding as input and output an estimated accuracy. The predictor itself is learned from a small subset of architectures that have been fully trained. Alternatively, **early stopping** techniquesZhang et al. [2022] use partial training to guess final accuracy, which cuts search time. However, these predictions can be highly variable if the underlying loss surface is rugged or if the training set for predictor modeling is too small Choromanska et al. [2014], Xie et al. [2020], White et al. [2021].

**Latency and Parameters**    **Latency** (inference time) is crucial for real-time tasks. Yet, parameter count (or FLOPs) and latency do not necessarily correlate on modern hardware Li et al. [2021], Zhang et al. [2022]. In some cases, a model with more parameters can run faster if its architecture aligns better with a device's parallelization. Consequently, many researchers measure latency on target hardware, use a **lookup table**, or train a **latency predictor** model for quick estimates Tan et al. [2018], Wang et al. [2020b].

**Energy and Robustness**    In battery-limited edge environments, **energy consumption** can be a direct objective Zhou et al. [2024]. Other specialized objectives include **adversarial robustness** Wu et al. [2024] and **uncertainty calibration**, depending on domain requirements.

**Neural Scaling Laws**   Large-scale studies Kaplan et al. [2020], Bahri et al. [2024] show that performance often follows empirical *scaling laws* w.r.t. data size ($D$), compute ($C$), and parameter count ($N$). While these laws do not prescribe a specific search method, they reveal diminishing returns beyond certain scales.

## 3.4   Multi-Objective Optimization and the Pareto Frontier

Real-world applications frequently require **multi-objective optimization (MOO)**. For instance, accuracy, latency, and memory can all matter simultaneously. A configuration is *Pareto-optimal* if no objective can be improved without worsening another Freitas [2024].

**Pareto Frontier**   Collectively, all Pareto-optimal solutions form the **Pareto frontier**. Practitioners often select from these trade-off points according to specific hardware or mission needs.

**Lexicographic Optimization**   Some approaches prioritize objectives in a strict hierarchy Abernethy et al. [2024]. For example, one might first maximize accuracy, then among top candidates minimize latency, and only then consider FLOPs. This method yields a single final solution guaranteed to lie on the Pareto frontier but might ignore equally valid trade-offs.

**Scalarization**   Alternatively, multiple objectives can be merged into a single metric via additive or multiplicative factors. MnasNet Tan et al. [2018], for instance, uses a function that blends accuracy with latency constraints. One drawback is that scalarization can bias solutions toward certain regions of the Pareto frontier Zhang et al. [2022].

**Hybrid Approaches**   Some authors combine lexicographic methods with scalarization or even dynamically shift objective weights during search Freitas [2024]. This can systematically explore multiple zones of the frontier.

## 3.5   Constraints

Apart from explicit objectives, constraints (e.g., hardware memory limits, real-time latency caps) often further restrict NAS. A common technique is the **lookup table (LUT)** approach Wang et al. [2020b], which stores per-layer or per-block cost (latency, energy, etc.) on the target hardware. The total cost of a candidate architecture is approximated by summing the LUT entries. While computationally cheap, LUT-based sums can be imprecise when concurrency or caching effects become significant Tan et al. [2018]. Hence, LUT methods typically function as a quick screening step to discard obviously infeasible designs.

## 3.6   NAS Optimization Algorithms

Given the (often large) search space and potential multi-objective setting, authors have proposed various optimization algorithms:

**Reinforcement Learning (RL)**   Early NAS work by Zoph and Le Zoph and Le [2017] used a policy-gradient RL approach to iteratively propose architectures, receiving rewards based on validation accuracy. This framework naturally extends to **pruning** or **quantization** decisions Bencsik and Szemenyei [2022], He et al. [2018] by defining actions (e.g., "remove 10% of weights") and rewards (e.g., trade-off in accuracy vs. model size). Parameter sharing Pham et al. [2018] can reduce redundant computation among candidates.

**Evolutionary Algorithms**   Evolutionary strategies view each architecture as an organism in a population Real et al. [2018, 2017]. After evaluating each architecture's metrics, top performers "survive" to the next generation, while new ones are created by mutation (altering some layers/operations) or crossover (combining parts of two architectures).

**Multi-Objective Gradient Descent**   When architectures can be continuously relaxed (as in DARTS Liu et al. [2018]), **multi-objective gradient descent** Désidéri [2012], Sukthanker et al. [2024] can simultaneously optimize multiple objectives. This requires carefully balancing gradients for each objective at each iteration. However, pure differentiable methods may struggle with the inherently discrete nature of many architectural choices.

**Bayesian Optimization (BO)**   In the context of NAS + HPO, a **surrogate model** predicts performance from an architecture encoding; an *acquisition function* (e.g., expected improvement, upper confidence bound) balances exploration of uncertain regions against exploitation of promising areas. BANANAS White et al. [2019], AutoDistill Zhang et al. [2022], and Transfer NAS Shala et al. [2023] exemplify BO for NAS. Notably, BO can maintain uncertainty estimates

over the performance predictor to adaptively guide the search toward regions most likely to yield better Pareto-optimal architectures.

In summary, **Neural Architecture Search** offers a systematic way to discover efficient models under diverse objectives (accuracy, latency, energy) and constraints (memory, real-time performance). By carefully defining the search space, leveraging surrogates or hypernetworks to reduce training costs, and choosing suitable algorithms, NAS can automate the design of neural networks for resource-constrained edge computing scenarios.

## 4    Compiler and Deployment Frameworks

While model compression and neural architecture search aim to yield lightweight networks, additional *deployment-specific* optimizations can be achieved via specialized compilers and hardware platforms. These compilers transform trained models into optimized code that exploits hardware-specific features (e.g., tensor cores, vector units, sparse processing) for faster inference and lower power consumption. Below, we summarize several notable solutions.

### 4.1    Industry-Developed Compilers

**TVM    Tensor Virtual Machine (TVM)** Chen et al. [2018] is an open-source compiler that parses a trained neural network model into a hardware-agnostic computational graph. It then applies graph-level optimizations (e.g., operator fusion) and generates low-level code tuned to the target hardware. Experiments show speedups of 1.2–3.8× on CPU, GPU, and FPGA backends.

**XLA    Accelerated Linear Algebra (XLA)**[1] is another open-source compiler (primarily developed by Google) that similarly decomposes a model into a high-level graph. It fuses operations into hardware-aware kernels, often achieving lower latency than unoptimized TensorFlow or JAX execution.

**GroqFlow    Groq** provides custom hardware based on the Tensor-Streaming Processor (TSP) architecture Gwennap [2020]. Its compiler, GroqFlow, converts trained models into programs for Groq devices. Internal benchmarks report a 2× increase in images-per-second on ResNet-50 compared to an Nvidia V100 GPU.

**OpenVINO    OpenVINO** Ria Chruvu is Intel's open-source toolkit that includes hardware-aware optimizations for CPUs, GPUs, and other Intel accelerators. It also ships with a *Neural Network Compression Framework (NNCF)* supporting pruning and quantization (both post-training and during training). This integration can further reduce inference time and memory footprint on Intel platforms.

**TensorRT    TensorRT** Zhou and Yang [2022] is Nvidia's open-source compiler for its GPUs and specialized accelerators. It offers optimizations for quantization and structured sparsity, and includes **TensorRT-LLM** for large language models. Empirical results suggest at least a 50% reduction in latency and energy usage for LLM inference Zhou et al. [2024].

### 4.2    Specialized Hardware for Edge AI: IBM NorthPole

Beyond general-purpose compilers, certain vendors design hardware specifically for *edge* inference. **IBM NorthPole** Cassidy et al. [2024] is a neuromorphic chip that uses event-driven processing to achieve high energy efficiency. By activating computation only on relevant data and pruning unnecessary paths, it significantly reduces power draw. On-chip memory further cuts down latency and off-chip transfers. NorthPole also supports weight pruning and quantization at the hardware level, compressing model layers and thereby accelerating inference. This combination of sparse data processing, event-driven operation, and local memory makes NorthPole particularly suitable for low-power, real-time tasks in edge environments.

### 4.3    Relation to the Overall Pipeline

These compilers and specialized hardware are often the final step in a deployment pipeline. Even after *model compression* (Section 2) or *NAS* (Section 3), the resulting model can benefit further from hardware-tailored graph optimizations. Consequently, the choice of compiler or hardware backend directly impacts whether or not the theoretical speedups from compression or NAS fully manifest in real-world inference. Practitioners aiming for efficient edge deployment typically combine all three:

---

[1] https://github.com/openxla/xla

1. Model compression or NAS to produce a compact, high-performing architecture,

2. Possibly re-quantizing or pruning for the specific device,

3. Compiling with a toolchain (TVM, TensorRT, etc.) specialized for the target hardware.

In summary, compiler frameworks and specialized hardware backends constitute a critical link in the optimization chain, translating theoretical gains from compression and architecture search into tangible, real-time performance improvements.

# 5 Integrated Approaches and Case Studies

In this section, we first discuss common *experimental methods* and publicly available *NAS benchmarks* that facilitate reproducible research. We then highlight several *case studies* and *integrated approaches* that combine the design elements introduced earlier: compression (pruning, quantization, distillation), NAS + HPO, and multi-objective optimization. The examples illustrate how these techniques work together to meet the unique demands of edge computing.

## 5.1 Common NAS Benchmarks

**NAS Benchmarks** provide standardized environments for evaluating and comparing different neural architecture search algorithms. Such benchmarks pair a fixed *search space* with tasks and often include pre-computed results or surrogates to minimize the computational overhead of repeatedly training candidate architectures.

**NAS-Bench-101**    Proposed in Ying et al. [2019], **NAS-Bench-101** targets the CIFAR-10 dataset, representing each candidate architecture as a Directed Acyclic Graph (DAG) with up to 7 vertices and 9 edges. Allowed operations are $3 \times 3$ convolution, $1 \times 1$ convolution, and $3 \times 3$ max pooling. Each valid DAG is then repeated (stacked) three times to form the final network. This search space encompasses $\sim 423{,}000$ distinct architectures after excluding duplicates or invalid topologies. For each architecture, the authors trained 3 random initializations to convergence and averaged the final accuracy. This yields a *lookup table* mapping {architecture $\rightarrow$ accuracy, training time}, enabling researchers to compare NAS algorithms without the heavy cost of full model training each time. However, NAS-Bench-101 only provides *single-objective* data (accuracy and runtime), limiting its direct utility for multi-objective (e.g. accuracy-latency) research.

**NAS-Bench-201**    Dong and Yang [2020] extend the idea to three datasets: CIFAR-10, CIFAR-100, and ImageNet-16-20 (see Fig. 8). This benchmark defines a more compact search space of about 15,600 architectures. Each is formed by stacking five identical cells, where each cell has four internal nodes and five possible operations. Performance metrics include model accuracy, latency, floating-point operations (FLOPs), parameter count, and training time. Notably, NAS-Bench-201 also provides full learning curves over 200 epochs, facilitating deeper analyses of training dynamics. Its main limitation is that the architectures themselves are still relatively small compared to modern large-scale applications.
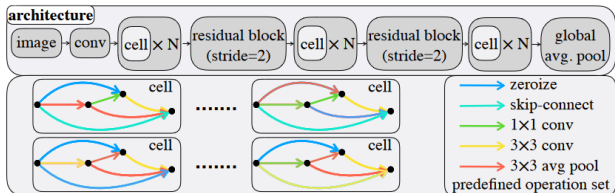


Figure 8: NAS-Bench-201: Each candidate cell has four nodes, and the overall network consists of five stacked cells Dong and Yang [2020].

**NAS-Bench-301**    Recognizing that NAS-Bench-101 and 201 cover relatively small design spaces, **NAS-Bench-301** Siems et al. [2020] takes a different approach. Instead of enumerating all possible architectures and storing their trained results in a table, it trains a *surrogate model* to predict performance in the *DARTS-based* search space on CIFAR-10. This surrogate leverages an ensemble of Graph Isomorphism Networks (GINs) to predict accuracy and an additional LGBoost-based predictor for latency. The dataset includes $\sim 60{,}000$ fully trained architectures. By sampling new designs and querying the surrogate, researchers can approximate performance quickly, enabling multi-objective or hardware-aware NAS studies at larger scales than NAS-Bench-101/201.

**NAS-Bench-Suite**    Finally, Mehta et al. [2022] aggregate multiple NAS benchmarks (including the ones above) into a single open-source repository (NasLib[2]). The authors observe that approaches or hyperparameters that excel in one benchmark may not generalize to others. Consequently, NAS-Bench-Suite helps evaluate NAS algorithms more robustly across different tasks and search spaces.

## 5.2    Case Studies Integrating Compression and NAS

With the foundational benchmarks in mind, we now examine several integrated approaches. These methods combine elements of **model compression** (pruning, quantization, distillation), **NAS + HPO**, and **hardware constraints** to achieve efficient inference on resource-limited devices.

### 5.2.1    APQ: Pruning, Quantization, and Supernetworks Wang et al. [2020b]

**APQ** (*Accuracy Predictor for Quantization*) unifies pruning, quantization, supernetwork-based NAS, multi-objective optimization, and a hardware-driven *lookup table* approach (see Fig. 9):

1. **Supernetwork Training**: First, a large convolutional neural network (CNN) supernetwork is trained on an image classification task. This supernetwork implicitly contains many sub-architectures (via pruning).

2. **Progressive Shrinking and Validation**: By pruning different subsets of channels or layers, one obtains sub-architectures, each with its own validation accuracy. Encoding each sub-architecture as a bit-vector plus its measured accuracy builds a dataset of {architecture, accuracy} pairs.

3. **Quantization Policies**: The authors apply different quantization policies (e.g., 8-bit, mixed precision) to those sub-architectures and record resulting accuracies. These tuples train an *accuracy predictor* that simultaneously captures the impact of pruning *and* quantization.

4. **Evolutionary Search**: With this accuracy predictor, APQ searches over the joint space of architecture (pruning) and quantization policies.

5. **Lookup Table Constraints**: Finally, a hardware LUT (layer-level latency and energy estimates) is summed for each candidate. Subnetworks exceeding allowable constraints are discarded. This yields a Pareto set of feasible (architecture, quantization) solutions.

Experiments show that APQ can effectively find high-accuracy, low-latency CNNs by combining these modular compression and search steps in a single pipeline.
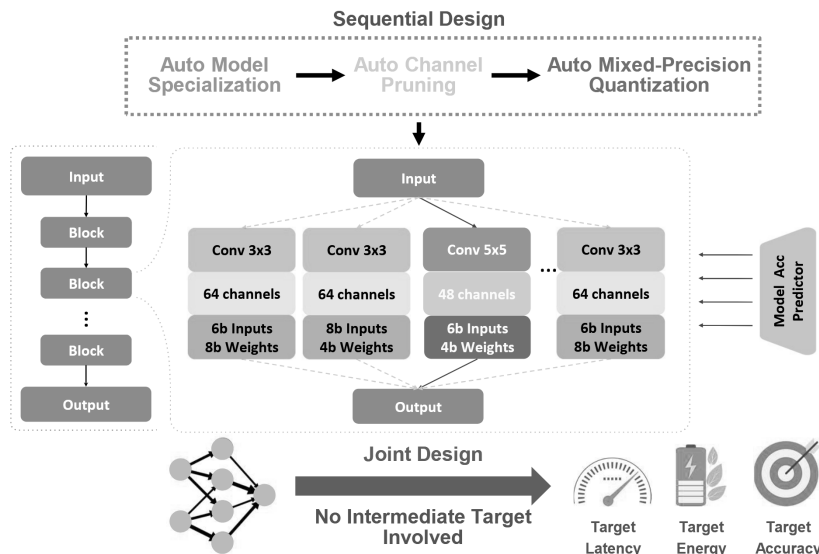


Figure 9: APQ Framework Wang et al. [2020b]. The supernetwork (left) supports multiple sub-architectures via pruning, and each sub-architecture can have different quantization policies. An accuracy predictor plus hardware LUT guide an evolutionary search.

---

[2]https://github.com/automl/naslib

### 5.2.2  DARTS: Differentiable Architecture Search Liu et al. [2018]

**Differentiable ARchiTecture Search (DARTS)** tackles NAS by relaxing discrete architecture choices into continuous, learnable parameters. Concretely:

- A DAG is defined where each edge can be one of several operations (e.g., $3 \times 3$ conv, $5 \times 5$ conv, skip connection).
- Real-valued coefficients (architecture parameters) weigh these operations.
- A *bilevel* optimization scheme alternates updating *model weights* (on the training set) and *architecture parameters* (on the validation set).

After convergence, the operation with the highest coefficient on each edge is selected, yielding a discrete architecture (See Fig. 10). The authors demonstrate that DARTS can reduce search time for CIFAR-10 from thousands of GPU-days (NASNet-A) to about 4 GPU-days, drastically lowering the computational barrier.
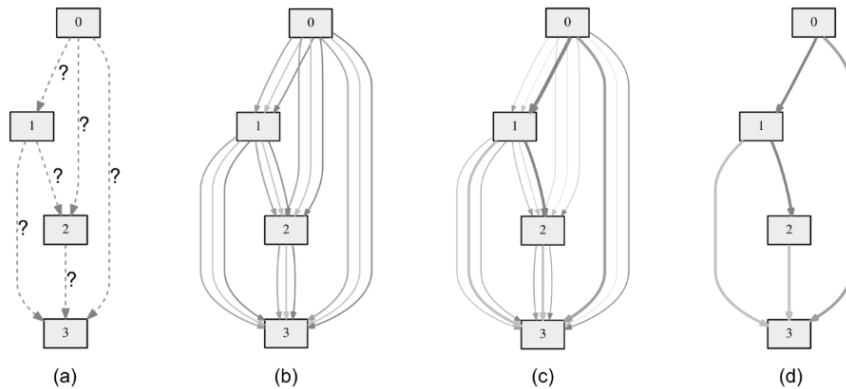


Figure 10: DARTS Framework Liu et al. [2018]. A DAG with multiple candidate operations per edge is learned via continuous weights, then discretized after training.

Notably, DARTS itself does not inherently integrate pruning or quantization, but variants have extended the differentiable approach to incorporate these techniques.

### 5.2.3  AWQ: Activation-Aware Weight Quantization Lin et al. [2024a]

**AWQ** focuses on compressing large transformer-based models via *channel-level quantization*. Observing that most weights can be quantized aggressively with minimal loss, the authors design a scheme to treat only the top $\sim 1\%$ of high-activation channels with higher precision:

1. Measure *average activation* per channel on a validation set.
2. Retain full precision for the channels with the largest activations (e.g., top 1%), quantize the rest to 4-bit floats.
3. Optionally allow fine-tuning if training data is available.

This selective strategy yields speedups of 3–4$\times$ for large LLMs (e.g., LLaMA, OPT) while preserving perplexity or accuracy on benchmarks (captioning tasks, visual-language tasks, etc.). Latency improvements are also reported. Although AWQ focuses on quantization, it complements the notion of NAS or pruning by providing a targeted approach to per-channel precision adjustments.

### 5.2.4  AutoDistill: Bayesian NAS + Distillation Zhang et al. [2022]

**AutoDistill** unifies NAS, teacher-student distillation, and hardware-aware objectives. The pipeline, see Fig. 11:

1. **Input Model and Constraints**: Start with a large, pre-trained model (the teacher), target hardware constraints (e.g., max memory, latency), and a discrete search space (24-layer stack of cells Sun et al. [2020]).
2. **Bayesian Optimization (BO)**: Use Google Vizier Golovin et al. [2017] to select candidate architectures for the student. A *flash distillation* (early stopping) approach provides a fast proxy for the final accuracy after knowledge distillation, mitigating the cost of full retraining.

3. **Distillation Loss**: Combine multi-head attention (MHA) loss, feature map (FM) loss, and logit loss Sun et al. [2020] to train the student. Mismatched teacher-student layers are handled via an intermediate fully-connected layer.

4. **Iterative Search**: Candidate architectures are evaluated on hardware for latency, FLOPs, memory usage, etc. These measurements inform the BO's acquisition function for the next iteration.

5. **Full Training**: Promising architectures undergo complete distillation (200 epochs) for final accuracy. The end result is a smaller model that meets hardware constraints and approximates the teacher's performance.


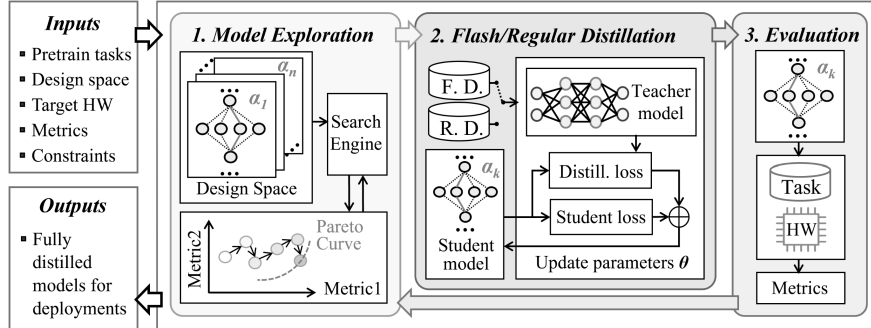
Figure 11: **AutoDistill** pipeline Zhang et al. [2022]. The system repeatedly samples architectures (driven by Bayesian Optimization), does a quick "flash" distillation, and measures hardware metrics.

By combining Bayesian NAS and knowledge distillation in a single loop, AutoDistill exemplifies how compression and architecture design can be simultaneously optimized for edge deployment.

**Flash Distillation Correlation** An intriguing component of AutoDistill is **Flash Distillation**, which uses only $\sim 5\%$ of the full training steps to predict final performance. Fig. 12 shows a high correlation between flash-distilled accuracy and fully-distilled accuracy, validating the method as a surrogate.
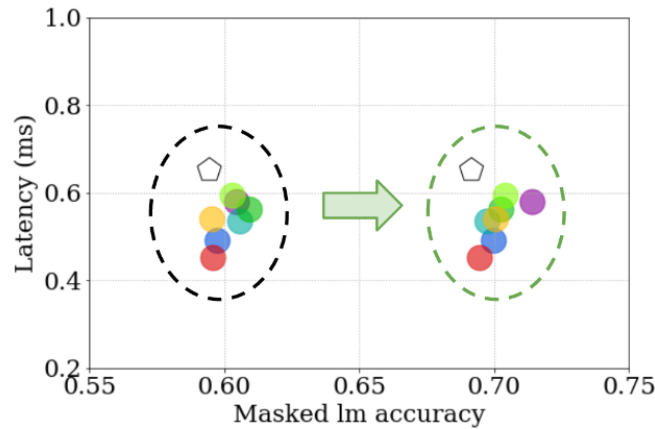


Figure 12: Flash Distillation correlates strongly with Full Distillation Zhang et al. [2022]. Each point represents a candidate architecture.

## 5.3 Summary of Integrated Methods

Overall, these case studies demonstrate how combining **model compression**, **NAS + HPO**, and **hardware constraints** can lead to more specialized, efficient, and robust models for edge devices. The synergy arises because each approach addresses different stages of the pipeline:

- *Compression* (pruning, quantization, distillation) reduces model size and complexity.

- *NAS + HPO* explores architectural variations to find topologies inherently suited to resource constraints.

- *Hardware-Aware Constraints* (LUT-based or direct measurement) ensure feasibility and speed on real devices.

By leveraging established benchmarks like NAS-Bench-101/201/301 and integrated frameworks such as APQ, DARTS, AWQ, and AutoDistill, researchers can systematically evaluate and refine AI models for edge deployment.

# 6  Conclusion and Discussion

Our comprehensive survey of the AI model optimization landscape aimed to identify a set of effective approaches that could inform the development of a scalable, generalized framework for optimizing AI models - one that is both model-agnostic and platform-independent. This survey not only sheds light on the complexities of a generalized model optimization but also paves the way for formulating a mathematically tractable problem formulation that can facilitate inference acceleration across diverse AI models deployed over resource-constrained edge computing platforms. We identified three broad methods for inference acceleration. The first method is model compression, composed of pruning, quantization, tensor decomposition, and distillation. These methods start with a large, high-performing model and find ways to achieve comparable performance with less resources. The second method is Neural Architecture Search, which searches for the optimal architecture for a given task and sometimes hardware. The third method of inference acceleration is choice of compiler; at deployment time, how the model packaged and and sent for computation. This is rapidly changing given the proliferation of new hardware architectures.

While the work in distillation has had notable victories, we note room for more fundamental work comparing loss topologies for different distillation formulations. As a practical matter, a researcher needs to estimate how much data and training time are needed to fairly evaluate a model. A comparison of this kind on a fixed task, dataset, and architecture but varying the distillation loss formulation may give the researcher an intuition for how different formulations may affect the topology. What do different distillation formulations do to the loss topology? Are there more minima? Are they lower overall? Are there more paths to the same minima? Researchers may need to know how different distillations may be suited for specific tasks. If a topology is jagged, current evidence suggests it could benefit from self-distillation which regularizes and softens the landscape.Mobahi et al. [2020] What other experiments and conclusions can we draw like this?

Neural Architecture Search is well studied by the AutoML community, but research frontiers still abound. Of particular interest are the Program Synthesis approaches that search for hierarchical motifs and compose new architectures in terms of those motifs, fighting the combinatorial explosion in the search space of large, complex architectures. While still an unsolved problem, advances in Neurosymbolic AI may advance the field dramatically in the near future, building custom model architectures while finding hardware-aware motifs.

## 6.1  Discussion and Future Directions

In this survey, we examined a range of strategies for optimizing AI models under the stringent resource constraints often found in edge computing environments. Three major categories of inference acceleration methods emerged:

1. **Model Compression**: Techniques such as pruning, quantization, tensor decomposition, and distillation compress large, over-parameterized models to more manageable forms. These methods aim to preserve near state-of-the-art performance while significantly reducing memory, compute, and energy overhead.

2. **Neural Architecture Search (NAS)**: NAS systematically explores architectural variations to identify designs inherently suited to the performance, latency, and memory requirements of a given task or hardware. When combined with hyperparameter optimization (HPO), it provides a comprehensive strategy for discovering configurations that exploit model redundancy and meet specific deployment constraints.

3. **Compiler and Deployment Frameworks**: Specialized toolchains (e.g., TVM, TensorRT, OpenVINO) optimize model graphs at the implementation level, leveraging hardware-specific features (such as operator fusion or sparse kernels) to further accelerate inference.

**Challenges and Open Questions in Distillation.** Although knowledge distillation has yielded important real-world benefits-particularly for large language models-its theoretical underpinnings remain partially unexplored. For instance, self-distillation appears to regularize the training process Mobahi et al. [2020], but there is limited understanding of why certain distillation objectives might induce flatter minima or better generalization. Furthermore, mismatches between teacher and student architectures raise questions about how to measure "distance" in latent feature space, how to handle domain gaps (e.g., cross-modal transfers), and how much data or training time is required to achieve a robust student

model. Deeper investigations of the *loss landscapes* shaped by distinct distillation formulations could clarify which methods are most appropriate for particular tasks or modalities.

**Frontiers in NAS.** While NAS has seen considerable progress in the AutoML community, several frontiers remain relatively unexplored:

- *Hierarchical / Program-Synthesis Approaches*: Searching for high-level motifs or symbolic expressions within vast design spaces can mitigate combinatorial blow-ups, especially for very deep networks. Advances in neurosymbolic AI may one day automate hardware-specific or domain-specific architecture discovery with minimal manual intervention.
- *Scalability to Very Large Models*: Applying NAS to large language models (LLMs) remains computationally daunting; supernetwork or hypernetwork-based approaches face scalability hurdles with billions of parameters. Hybrid techniques (e.g., partial weight sharing) or approximate evaluation methods may be necessary to push NAS further into the LLM realm.

**Sparse Literature on Pre-Training Pruning.** Despite the potential of pruning *during* or *before* large-scale training—an approach that could save considerable compute costs—existing literature is sparse for foundation-scale models. The expense of even partially training a massive model often deters researchers. More efficient or incremental pre-training pruning algorithms could enable scaling down LLMs at earlier stages, unlocking substantial resource savings.

**Toward a Unified, Model-Agnostic Framework.** Collectively, these insights suggest a path toward a *unified optimization pipeline* for edge AI, one that integrates:

- *Adaptive Compression* (e.g., distillation, quantization, pruning) guided by data- or task-specific requirements;
- *Neural Architecture Search* that systematically navigates both macro- and micro-architecture choices (potentially via hierarchical or symbolic encodings);
- *Hardware-Aware Compilation* ensuring that any discovered design truly yields speedups and memory savings in deployment.

A central challenge remains formalizing this pipeline as a tractable multi-objective optimization problem, balancing metrics like accuracy, latency, energy, and model size under varying device constraints.

**Concluding Remarks.** Although our survey did not explicitly focus on automated multi-modal architecture design, many of the discussed techniques-particularly cross-modal distillation and hardware-aware NAS-could be extended to multi-modal tasks. With deep learning models continuing to grow in complexity and scale, the convergence of model compression, NAS, and advanced compiler optimizations promises significant gains in efficiency and deployability. Bridging theoretical understanding (e.g., loss-topology analyses in distillation or hierarchical motif synthesis in NAS) with the practical realities of modern hardware stands as an exciting frontier. We hope this survey spurs further research and collaboration across these disciplines, ultimately enabling AI systems that are both powerful and feasible in real-world, resource-constrained deployments.

## Acknowledgements

## References

Benoit Jacob, Skirmantas Kligys, Bo Chen, Menglong Zhu, Matthew Tang, Andrew Howard, Hartwig Adam, and Dmitry Kalenichenko. Quantization and training of neural networks for efficient integer-arithmetic-only inference. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2704–2713, 2018.

Mengzhao Chen, Wenqi Shao, Peng Xu, Jiahao Wang, Peng Gao, Kaipeng Zhang, Yu Qiao, and Ping Luo. Efficientqat: Efficient quantization-aware training for large language models, 2024a. URL `https://arxiv.org/abs/2407.11062`.

Yanshu Wang, Tong Yang, Xiyan Liang, Guoan Wang, Hanning Lu, Xu Zhe, Yaoming Li, and Li Weitao. Art and science of quantizing large-scale models: A comprehensive overview. *arXiv preprint arXiv:2409.11650*, 2024.

Sean I Young. Foundations of large language model compression–part 1: Weight quantization. *arXiv preprint arXiv:2409.02026*, 2024.

Ji Lin, Jiaming Tang, Haotian Tang, Shang Yang, Wei-Ming Chen, Wei-Chen Wang, Guangxuan Xiao, Xingyu Dang, Chuang Gan, and Song Han. AWQ: Activation-aware Weight Quantization for LLM Compression and Acceleration, July 2024a. URL http://arxiv.org/abs/2306.00978. arXiv:2306.00978 [cs].

Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. Qlora: Efficient finetuning of quantized llms, 2023. URL https://arxiv.org/abs/2305.14314.

Yijia Zhang, Lingran Zhao, Shijie Cao, Sicheng Zhang, Wenqiang Wang, Ting Cao, Fan Yang, Mao Yang, Shanghang Zhang, and Ningyi Xu. Integer or floating point? new outlooks for low-bit quantization on large language models. In *2024 IEEE International Conference on Multimedia and Expo (ICME)*, page 1–6. IEEE, July 2024. doi:10.1109/icme57554.2024.10688089. URL http://dx.doi.org/10.1109/ICME57554.2024.10688089.

Zhihang Yuan, Lin Niu, Jiawei Liu, Wenyu Liu, Xinggang Wang, Yuzhang Shang, Guangyu Sun, Qiang Wu, Jiaxiang Wu, and Bingzhe Wu. Rptq: Reorder-based post-training quantization for large language models, 2023. URL https://arxiv.org/abs/2304.01089.

Xiuying Wei, Yunchen Zhang, Yuhang Li, Xiangguo Zhang, Ruihao Gong, Jinyang Guo, and Xianglong Liu. Outlier suppression+: Accurate quantization of large language models by equivalent and optimal shifting and scaling, 2023. URL https://arxiv.org/abs/2304.09145.

Paulius Micikevicius, Sharan Narang, Jonah Alben, Gregory Frederick Diamos, Erich Elsen, David García, Boris Ginsburg, Michael Houston, Oleksii Kuchaiev, Ganesh Venkatesh, and Hao Wu. Mixed precision training. *ArXiv*, abs/1710.03740, 2017. URL https://api.semanticscholar.org/CorpusID:3297437.

Hao Li, Yuzhu Wang, Yan Hong, Fei Li, and Xiaohui Ji. Layered mixed-precision training: A new training method for large-scale ai models. *Journal of King Saud University - Computer and Information Sciences*, 35(8):101656, September 2023a. ISSN 1319-1578. doi:10.1016/j.jksuci.2023.101656. URL http://dx.doi.org/10.1016/j.jksuci.2023.101656.

Zihan Chen, Bike Xie, Jundong Li, and Cong Shen. Channel-wise mixed-precision quantization for large language models, 2024b. URL https://arxiv.org/abs/2410.13056.

Junhao Xu, Jianwei Yu, Shoukang Hu, Xunying Liu, and Helen Meng. Mixed precision low-bit quantization of neural network language models for speech recognition. *IEEE/ACM Transactions on Audio, Speech, and Language Processing*, 29:3679–3693, 2021. ISSN 2329-9304. doi:10.1109/taslp.2021.3129357. URL http://dx.doi.org/10.1109/TASLP.2021.3129357.

Gaurav Menghani. Efficient deep learning: A survey on making deep learning models smaller, faster, and better. *ACM Comput. Surv.*, 55(12), March 2023. ISSN 0360-0300. doi:10.1145/3578938. URL https://doi.org/10.1145/3578938.

Google AI. Model optimization - google ai. https://ai.google.dev/edge/litert/models/model_optimization, 2021.

Jiedong Lang, Zhehao Guo, and Shuyu Huang. A comprehensive study on quantization techniques for large language models, 2024. URL https://arxiv.org/abs/2411.02530.

Anna Choromanska, Mikael Henaff, Michaël Mathieu, Gérard Ben Arous, and Yann LeCun. The loss surface of multilayer networks. *CoRR*, abs/1412.0233, 2014. URL http://arxiv.org/abs/1412.0233.

Hongrong Cheng, Miao Zhang, and Javen Qinfeng Shi. A survey on deep neural network pruning-taxonomy, comparison. *Analysis, and Recommendations*, 2023.

Namhoon Lee, Thalaiyasingam Ajanthan, and Philip H. S. Torr. Snip: Single-shot network pruning based on connection sensitivity, 2019. URL https://arxiv.org/abs/1810.02340.

Chaoqi Wang, Guodong Zhang, and Roger Grosse. Picking winning tickets before training by preserving gradient flow, 2020a. URL https://arxiv.org/abs/2002.07376.

Elias Frantar and Dan Alistarh. Sparsegpt: Massive language models can be accurately pruned in one-shot, 2023. URL https://arxiv.org/abs/2301.00774.

Yihui He, Xiangyu Zhang, and Jian Sun. Channel pruning for accelerating very deep neural networks, 2017. URL https://arxiv.org/abs/1707.06168.

Zhonghui You, Kun Yan, Jinmian Ye, Meng Ma, and Ping Wang. Gate Decorator: Global Filter Pruning Method for Accelerating Deep Convolutional Neural Networks. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc., 2019. URL https://proceedings.neurips.cc/paper_files/paper/2019/hash/b51a15f382ac914391a58850ab343b00-Abstract.html.

Kyuhong Shim, Iksoo Choi, Wonyong Sung, and Jungwook Choi. Layer-wise pruning of transformer attention heads for efficient language modeling. In *2021 18th International SoC Design Conference (ISOCC)*, pages 357–358, 2021. doi:10.1109/ISOCC53507.2021.9613933.

Saleh Ashkboos, Maximilian L. Croci, Marcelo Gennari do Nascimento, Torsten Hoefler, and James Hensman. Slicegpt: Compress large language models by deleting rows and columns, 2024. URL `https://arxiv.org/abs/2401.15024`.

Xin Men, Mingyu Xu, Qingyu Zhang, Bingning Wang, Hongyu Lin, Yaojie Lu, Xianpei Han, and Weipeng Chen. Shortgpt: Layers in large language models are more redundant than you expect, 2024. URL `https://arxiv.org/abs/2403.03853`.

Xinyin Ma, Gongfan Fang, and Xinchao Wang. Llm-pruner: On the structural pruning of large language models, 2023. URL `https://arxiv.org/abs/2305.11627`.

Kaixin Xu, Zhe Wang, Chunyun Chen, Xue Geng, Jie Lin, Xulei Yang, Min Wu, Xiaoli Li, and Weisi Lin. Lpvit: Low-power semi-structured pruning for vision transformers, 2024a. URL `https://arxiv.org/abs/2407.02068`.

Xiaolong Ma, Wei Niu, Tianyun Zhang, Sijia Liu, Sheng Lin, Hongjia Li, Xiang Chen, Jian Tang, Kaisheng Ma, Bin Ren, and Yanzhi Wang. An image enhancing pattern-based sparsity for real-time inference on mobile devices, 2020. URL `https://arxiv.org/abs/2001.07710`.

Fanxu Meng, Hao Cheng, Ke Li, Huixiang Luo, Xiaowei Guo, Guangming Lu, and Xing Sun. Pruning filter in filter, 2020. URL `https://arxiv.org/abs/2009.14410`.

Hidenori Tanaka, Daniel Kunin, Daniel L. K. Yamins, and Surya Ganguli. Pruning neural networks without any data by iteratively conserving synaptic flow, 2020. URL `https://arxiv.org/abs/2006.05467`.

Namhoon Lee, Thalaiyasingam Ajanthan, Stephen Gould, and Philip H. S. Torr. A signal propagation perspective for pruning neural networks at initialization, 2020. URL `https://arxiv.org/abs/1906.06307`.

Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners, 2021. URL `https://arxiv.org/abs/1911.11134`.

Wei Wen, Chunpeng Wu, Yandan Wang, Yiran Chen, and Hai Li. Learning structured sparsity in deep neural networks, 2016. URL `https://arxiv.org/abs/1608.03665`.

Zhuang Liu, Jianguo Li, Zhiqiang Shen, Gao Huang, Shoumeng Yan, and Changshui Zhang. Learning efficient convolutional networks through network slimming, 2017. URL `https://arxiv.org/abs/1708.06519`.

Xuefei Ning, Tianchen Zhao, Wenshuo Li, Peng Lei, Yu Wang, and Huazhong Yang. Dsa: More efficient budgeted pruning via differentiable sparsity allocation, 2020. URL `https://arxiv.org/abs/2004.02164`.

Tianzhe Wang, Kuan Wang, Han Cai, Ji Lin, Zhijian Liu, and Song Han. APQ: Joint Search for Network Architecture, Pruning and Quantization Policy, June 2020b. URL `http://arxiv.org/abs/2006.08509`. arXiv:2006.08509 [cs, stat].

Woosuk Kwon, Sehoon Kim, Michael W. Mahoney, Joseph Hassoun, Kurt Keutzer, and Amir Gholami. A fast post-training pruning framework for transformers, 2022. URL `https://arxiv.org/abs/2204.09656`.

Liyang Liu, Shilong Zhang, Zhanghui Kuang, Aojun Zhou, Jing-Hao Xue, Xinjiang Wang, Yimin Chen, Wenming Yang, Qingmin Liao, and Wayne Zhang. Group fisher pruning for practical network compression, 2021. URL `https://arxiv.org/abs/2108.00708`.

Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks, 2019. URL `https://arxiv.org/abs/1803.03635`.

Haoran You, Chaojian Li, Pengfei Xu, Yonggan Fu, Yue Wang, Xiaohan Chen, Richard G. Baraniuk, Zhangyang Wang, and Yingyan Lin. Drawing early-bird tickets: Towards more efficient training of deep networks, 2022. URL `https://arxiv.org/abs/1909.11957`.

Xiaohan Chen, Yu Cheng, Shuohang Wang, Zhe Gan, Zhangyang Wang, and Jingjing Liu. Earlybert: Efficient bert training via early-bird lottery tickets, 2021. URL `https://arxiv.org/abs/2101.00063`.

Yongming Rao, Jiwen Lu, Ji Lin, and Jie Zhou. Runtime network routing for efficient image classification. *IEEE Trans. Pattern Anal. Mach. Intell.*, 41(10):2291–2304, oct 2019. ISSN 0162-8828. doi:10.1109/TPAMI.2018.2878258. URL `https://doi.org/10.1109/TPAMI.2018.2878258`.

Yehui Tang, Yunhe Wang, Yixing Xu, Yiping Deng, Chao Xu, Dacheng Tao, and Chang Xu. Manifold regularized dynamic network pruning, 2021. URL `https://arxiv.org/abs/2103.05861`.

Song Han, Huizi Mao, and William J. Dally. Deep compression: Compressing deep neural networks with pruning, trained quantization and huffman coding, 2016. URL `https://arxiv.org/abs/1510.00149`.

Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient convnets, 2017. URL https://arxiv.org/abs/1608.08710.

Chenglong Zhao, Bingbing Ni, Jian Zhang, Qiwei Zhao, Wenjun Zhang, and Qi Tian. Variational convolutional neural network pruning. In *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2775–2784, 2019. doi:10.1109/CVPR.2019.00289.

Michael Santacroce, Zixin Wen, Yelong Shen, and Yuanzhi Li. What matters in the structured pruning of generative language models?, 2023. URL https://arxiv.org/abs/2302.03773.

Lucio Dery, Steven Kolawole, Jean-François Kagy, Virginia Smith, Graham Neubig, and Ameet Talwalkar. Everybody prune now: Structured pruning of llms with only forward passes, 2024. URL https://arxiv.org/abs/2402.05406.

Mingyang Zhang, Xinyi Yu, Jingtao Rong, and Linlin Ou. Graph pruning for model compression, 2021. URL https://arxiv.org/abs/1911.09817.

Yihui He, Ji Lin, Zhijian Liu, Hanrui Wang, Li-Jia Li, and Song Han. *AMC: AutoML for Model Compression and Acceleration on Mobile Devices*, page 815–832. Springer International Publishing, 2018. ISBN 9783030012342. doi:10.1007/978-3-030-01234-2_48. URL http://dx.doi.org/10.1007/978-3-030-01234-2_48.

Blanka Bencsik and Márton Szemenyei. Efficient neural network pruning using model-based reinforcement learning. In *2022 International Symposium on Measurement and Control in Robotics (ISMCR)*, pages 1–8, 2022. doi:10.1109/ISMCR56534.2022.9950598.

Zhuo Li, Hengyi Li, and Lin Meng. Model compression for deep neural networks: A survey. *Computers*, 12(3), 2023b. ISSN 2073-431X. doi:10.3390/computers12030060. URL https://www.mdpi.com/2073-431X/12/3/60.

Habib Hajimolahoseini, Mehdi Rezagholizadeh, Vahid Partovinia, Marzieh Tahaei, Omar Mohamed Awad, and Yang Liu. Compressing pre-trained language models using progressive low rank decomposition. *Advances in Neural Information Processing Systems*, 2021.

Rajarshi Saha, Naomi Sagan, Varun Srivastava, Andrea J. Goldsmith, and Mert Pilanci. Compressing large language models using low rank and low precision decomposition, 2024. URL https://arxiv.org/abs/2405.18886.

Mingxue Xu, Yao Lei Xu, and Danilo P. Mandic. Tensorgpt: Efficient compression of large language models based on tensor-train decomposition, 2023a. URL https://arxiv.org/abs/2307.00526.

Chi-Heng Lin, Shangqian Gao, James Seale Smith, Abhishek Patel, Shikhar Tuli, Yilin Shen, Hongxia Jin, and Yen-Chang Hsu. Modegpt: Modular decomposition for large language model compression, 2024b. URL https://arxiv.org/abs/2408.09632.

Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network, 2015. URL https://arxiv.org/abs/1503.02531.

Xiaohan Xu, Ming Li, Chongyang Tao, Tao Shen, Reynold Cheng, Jinyang Li, Can Xu, Dacheng Tao, and Tianyi Zhou. A Survey on Knowledge Distillation of Large Language Models, March 2024b. URL http://arxiv.org/abs/2402.13116. arXiv:2402.13116 [cs].

Micah Goldblum, Liam Fowl, Soheil Feizi, and Tom Goldstein. Adversarially robust distillation. *Proceedings of the AAAI Conference on Artificial Intelligence*, 34(04):3996–4003, April 2020. ISSN 2159-5399. doi:10.1609/aaai.v34i04.5816. URL http://dx.doi.org/10.1609/aaai.v34i04.5816.

Ian J. Goodfellow, Jean Pouget-Abadie, Mehdi Mirza, Bing Xu, David Warde-Farley, Sherjil Ozair, Aaron Courville, and Yoshua Bengio. Generative adversarial networks, 2014. URL https://arxiv.org/abs/1406.2661.

Jaewon Jung, Hongsun Jang, Jaeyong Song, and Jinho Lee. Peeraid: Improving adversarial distillation from a specialized peer tutor, 2024. URL https://arxiv.org/abs/2403.06668.

Hanting Chen, Yunhe Wang, Chang Xu, Zhaohui Yang, Chuanjian Liu, Boxin Shi, Chunjing Xu, Chao Xu, and Qi Tian. Data-free learning of student networks, 2019a. URL https://arxiv.org/abs/1904.01186.

Ying Zhang, Tao Xiang, Timothy M. Hospedales, and Huchuan Lu. Deep mutual learning, 2017. URL https://arxiv.org/abs/1706.00384.

Meet P. Vadera, Brian Jalaian, and Benjamin M. Marlin. Generalized bayesian posterior expectation distillation for deep neural networks. *CoRR*, abs/2005.08110, 2020. URL https://arxiv.org/abs/2005.08110.

Saurabh Gupta, Judy Hoffman, and Jitendra Malik. Cross modal distillation for supervision transfer, 2015. URL https://arxiv.org/abs/1507.00448.

Alex Andonian, Shixing Chen, and Raffay Hamid. Robust cross-modal representation learning with progressive self-distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16430–16441, 2022.

Fushuo Huo, Wenchao Xu, Jingcai Guo, Haozhao Wang, and Song Guo. C2kd: Bridging the modality gap for cross-modal knowledge distillation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16006–16015, 2024.

Jing Liu, Tongya Zheng, Guanzheng Zhang, and Qinfen Hao. Graph-based knowledge distillation: A survey and experimental evaluation, 2023a. URL `https://arxiv.org/abs/2302.14643`.

Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. Attention is all you need, 2023. URL `https://arxiv.org/abs/1706.03762`.

Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan Yang, and Ming Zhou. Minilm: Deep self-attention distillation for task-agnostic compression of pre-trained transformers. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 5776–5788. Curran Associates, Inc., 2020c. URL `https://proceedings.neurips.cc/paper_files/paper/2020/file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf`.

Kai Wang, Fei Yang, and Joost van de Weijer. Attention distillation: self-supervised vision transformer students need more guidance, 2022a. URL `https://arxiv.org/abs/2210.00944`.

Elliot J. Crowley, Gavin Gray, and Amos Storkey. Moonshine: Distilling with cheap convolutions, 2019. URL `https://arxiv.org/abs/1711.02613`.

Jianping Gou, Baosheng Yu, Stephen J. Maybank, and Dacheng Tao. Knowledge distillation: A survey. *International Journal of Computer Vision*, 129(6):1789–1819, March 2021. ISSN 1573-1405. doi:10.1007/s11263-021-01453-z. URL `http://dx.doi.org/10.1007/s11263-021-01453-z`.

Defang Chen, Jian-Ping Mei, Can Wang, Yan Feng, and Chun Chen. Online knowledge distillation with diverse peers. *CoRR*, abs/1912.00350, 2019b. URL `http://arxiv.org/abs/1912.00350`.

Hossein Mobahi, Mehrdad Farajtabar, and Peter Bartlett. Self-distillation amplifies regularization in hilbert space. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and H. Lin, editors, *Advances in Neural Information Processing Systems*, volume 33, pages 3351–3361. Curran Associates, Inc., 2020. URL `https://proceedings.neurips.cc/paper_files/paper/2020/file/2288f691b58edecadcc9a8691762b4fd-Paper.pdf`.

Linfeng Zhang, Jiebo Song, Anni Gao, Jingwei Chen, Chenglong Bao, and Kaisheng Ma. Be your own teacher: Improve the performance of convolutional neural networks via self distillation. *CoRR*, abs/1905.08094, 2019. URL `http://arxiv.org/abs/1905.08094`.

Yizhong Wang, Yeganeh Kordi, Swaroop Mishra, Alisa Liu, Noah Smith, Daniel Khashabi, and Hannaneh Hajishirzi. Self-instruct: Aligning language model with self generated instructions, 12 2022b.

Ming Li, Lichang Chen, Jiuhai Chen, Shwai He, Jiuxiang Gu, and Tianyi Zhou. Selective reflection-tuning: Student-selected data recycling for llm instruction-tuning. *ArXiv*, abs/2402.10110, 2024. URL `https://api.semanticscholar.org/CorpusID:267682220`.

Subhabrata Mukherjee, Arindam Mitra, Ganesh Jawahar, Sahaj Agarwal, Hamid Palangi, and Ahmed Awadallah. Orca: Progressive learning from complex explanation traces of gpt-4, 2023. URL `https://arxiv.org/abs/2306.02707`.

Canwen Xu, Daya Guo, Nan Duan, and Julian McAuley. Baize: An open-source chat model with parameter-efficient tuning on self-chat data, 2023b. URL `https://arxiv.org/abs/2304.01196`.

Minki Kang, Seanie Lee, Jinheon Baek, Kenji Kawaguchi, and Sung Ju Hwang. Knowledge-augmented reasoning distillation for small language models in knowledge-intensive tasks, 2023. URL `https://arxiv.org/abs/2305.18395`.

Akari Asai, Zeqiu Wu, Yizhong Wang, Avirup Sil, and Hannaneh Hajishirzi. Self-rag: Learning to retrieve, generate, and critique through self-reflection, 2023. URL `https://arxiv.org/abs/2310.11511`.

Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang, Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica. Judging llm-as-a-judge with mt-bench and chatbot arena, 2023. URL `https://arxiv.org/abs/2306.05685`.

Yann Dubois, Balázs Galambosi, Percy Liang, and Tatsunori B. Hashimoto. Length-controlled alpacaeval: A simple way to debias automatic evaluators, 2024. URL `https://arxiv.org/abs/2404.04475`.

Lewis Tunstall, Edward Beeching, Nathan Lambert, Nazneen Rajani, Kashif Rasul, Younes Belkada, Shengyi Huang, Leandro von Werra, Clémentine Fourrier, Nathan Habib, Nathan Sarrazin, Omar Sanseviero, Alexander M. Rush, and Thomas Wolf. Zephyr: Direct distillation of lm alignment, 2023. URL `https://arxiv.org/abs/2310.16944`.

Yuntao Bai et. al. Constitutional ai: Harmlessness from ai feedback, 2022. URL `https://arxiv.org/abs/2212.08073`.

Maryam Hashemzadeh, Elias Stengel-Eskin, Sarath Chandar, and Marc-Alexandre Cote. Sub-goal distillation: A method to improve small language agents, 2024. URL `https://arxiv.org/abs/2405.02749`.

Da Yin, Faeze Brahman, Abhilasha Ravichander, Khyathi Chandu, Kai-Wei Chang, Yejin Choi, and Bill Yuchen Lin. Agent lumos: Unified and modular training for open-source language agents, 2024a. URL `https://arxiv.org/abs/2311.05657`.

Baian Chen, Chang Shu, Ehsan Shareghi, Nigel Collier, Karthik Narasimhan, and Shunyu Yao. Fireact: Toward language agent fine-tuning, 2023. URL `https://arxiv.org/abs/2310.05915`.

Shishir G. Patil, Tianjun Zhang, Xin Wang, and Joseph E. Gonzalez. Gorilla: Large language model connected with massive apis, 2023. URL `https://arxiv.org/abs/2305.15334`.

Yuxian Gu, Li Dong, Furu Wei, and Minlie Huang. Minillm: Knowledge distillation of large language models, 2024. URL `https://arxiv.org/abs/2306.08543`.

Jongwoo Ko, Sungnyun Kim, Tianyi Chen, and Se-Young Yun. Distillm: Towards streamlined distillation for large language models, 2024. URL `https://arxiv.org/abs/2402.03898`.

Jaehun Jung, Peter West, Liwei Jiang, Faeze Brahman, Ximing Lu, Jillian Fisher, Taylor Sorensen, and Choi Yejin. Impossible distillation: from low-quality model to high-quality dataset and model for summarization and paraphrasing, 05 2023.

Krishna Srinivasan, Karthik Raman, Anupam Samanta, Lingrui Liao, Luca Bertelli, and Mike Bendersky. Quill: Query intent with large language models using retrieval augmentation and multi-stage distillation, 2022. URL `https://arxiv.org/abs/2210.15718`.

Haotian Liu, Chunyuan Li, Qingyang Wu, and Yong Jae Lee. Visual instruction tuning, 2023b. URL `https://arxiv.org/abs/2304.08485`.

Chenyang Lyu, Minghao Wu, Longyue Wang, Xinting Huang, Bingshuai Liu, Zefeng Du, Shuming Shi, and Zhaopeng Tu. Macaw-llm: Multi-modal language modeling with image, audio, video, and text integration, 2023. URL `https://arxiv.org/abs/2306.09093`.

Lu Yin, AJAY KUMAR JAISWAL, Shiwei Liu, Souvik Kundu, and Zhangyang Wang. Junk DNA hypothesis: Pruning small pre-trained weights $\textit{Irreversibly}$ and $\textit{Monotonically}$ impairs "difficult" downstream tasks in LLMs. In *Forty-first International Conference on Machine Learning*, 2024b. URL `https://openreview.net/forum?id=EfUrTeuUfy`.

Xin He, Kaiyong Zhao, and Xiaowen Chu. AutoML: A survey of the state-of-the-art. *Knowledge-Based Systems*, 212: 106622, January 2021. ISSN 0950-7051. doi:10.1016/j.knosys.2020.106622. URL `https://www.sciencedirect.com/science/article/pii/S0950705120307516`.

Barret Zoph and Quoc V. Le. Neural Architecture Search with Reinforcement Learning, February 2017. URL `http://arxiv.org/abs/1611.01578`. arXiv:1611.01578 [cs].

Han Cai, Tianyao Chen, Weinan Zhang, Yong Yu, and Jun Wang. Efficient Architecture Search by Network Transformation, November 2017. URL `http://arxiv.org/abs/1707.04873`. arXiv:1707.04873 [cs].

Colin White, Willie Neiswanger, Sam Nolen, and Yash Savani. A study on encodings for neural architecture search. *CoRR*, abs/2007.04965, 2020. URL `https://arxiv.org/abs/2007.04965`.

Colin White, Mahmoud Safari, Rhea Sukthanker, Binxin Ru, Thomas Elsken, Arber Zela, Debadeepta Dey, and Frank Hutter. Neural architecture search: Insights from 1000 papers, 2023. URL `https://arxiv.org/abs/2301.08727`.

Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: differentiable architecture search. *CoRR*, abs/1806.09055, 2018. URL `http://arxiv.org/abs/1806.09055`.

Charles Jin, Phitchaya Mangpo Phothilimthana, and Sudip Roy. Neural architecture search using property guided synthesis. *Proceedings of the ACM on Programming Languages*, 6(OOPSLA2):1150–1179, October 2022. ISSN 2475-1421. doi:10.1145/3563329. URL `http://dx.doi.org/10.1145/3563329`.

Xingchen Wan, Binxin Ru, Pedro M. Esperança, and Fabio Maria Carlucci. Approximate neural architecture search via operation distribution learning. *CoRR*, abs/2111.04670, 2021. URL `https://arxiv.org/abs/2111.04670`.

Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once-for-All: Train One Network and Specialize it for Efficient Deployment, April 2020. URL `http://arxiv.org/abs/1908.09791`. arXiv:1908.09791 [cs, stat].

Lingxi Xie, Xin Chen, Kaifeng Bi, Longhui Wei, Yuhui Xu, Zhengsu Chen, Lanfei Wang, An Xiao, Jianlong Chang, Xiaopeng Zhang, and Qi Tian. Weight-sharing neural architecture search: A battle to shrink the optimization gap. *CoRR*, abs/2008.01475, 2020. URL `https://arxiv.org/abs/2008.01475`.

Yawei Li, Shuhang Gu, Kai Zhang, Luc Van Gool, and Radu Timofte. Dhp: Differentiable meta pruning via hypernetworks, 2020. URL `https://arxiv.org/abs/2003.13683`.

Zechun Liu, Haoyuan Mu, Xiangyu Zhang, Zichao Guo, Xin Yang, Tim Kwang-Ting Cheng, and Jian Sun. Metapruning: Meta learning for automatic neural network channel pruning, 2019. URL `https://arxiv.org/abs/1903.10258`.

Chris Zhang, Mengye Ren, and Raquel Urtasun. Graph hypernetworks for neural architecture search. *CoRR*, abs/1810.05749, 2018. URL `http://arxiv.org/abs/1810.05749`.

Han Cai, Jiacheng Yang, Weinan Zhang, Song Han, and Yong Yu. Path-Level Network Transformation for Efficient Architecture Search, June 2018. URL `http://arxiv.org/abs/1806.02639`. arXiv:1806.02639 [cs, stat].

Xiaofan Zhang, Zongwei Zhou, Deming Chen, and Yu Emma Wang. Autodistill: an end-to-end framework to explore and distill hardware-efficient language models. *CoRR*, abs/2201.08539, 2022. URL `https://arxiv.org/abs/2201.08539`.

Colin White, Arber Zela, Robin Ru, Yang Liu, and Frank Hutter. How powerful are performance predictors in neural architecture search? In M. Ranzato, A. Beygelzimer, Y. Dauphin, P.S. Liang, and J. Wortman Vaughan, editors, *Advances in Neural Information Processing Systems*, volume 34, pages 28454–28469. Curran Associates, Inc., 2021. URL `https://proceedings.neurips.cc/paper_files/paper/2021/file/ef575e8837d065a1683c022d2077d342-Paper.pdf`.

Sheng Li, Mingxing Tan, Ruoming Pang, Andrew Li, Liqun Cheng, Quoc V. Le, and Norman P. Jouppi. Searching for fast model families on datacenter accelerators. *CoRR*, abs/2102.05610, 2021. URL `https://arxiv.org/abs/2102.05610`.

Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, and Quoc V. Le. Mnasnet: Platform-aware neural architecture search for mobile. *CoRR*, abs/1807.11626, 2018. URL `http://arxiv.org/abs/1807.11626`.

Yuxiao Zhou, Zhishan Guo, Zheng Dong, and Kecheng Yang. Multi-accelerator neural network inference via tensorrt in heterogeneous embedded systems. In *2024 IEEE 48th Annual Computers, Software, and Applications Conference (COMPSAC)*, pages 463–472, 2024. doi:10.1109/COMPSAC61105.2024.00070.

Yongtao Wu, Fanghui Liu, Carl-Johann Simon-Gabriel, Grigorios G Chrysos, and Volkan Cevher. Robust nas under adversarial training: benchmark, theory, and beyond, 2024. URL `https://arxiv.org/abs/2403.13134`.

Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.

Yasaman Bahri, Ethan Dyer, Jared Kaplan, Jaehoon Lee, and Utkarsh Sharma. Explaining neural scaling laws. *Proceedings of the National Academy of Sciences*, 121(27):e2311878121, 2024.

Alex A. Freitas. The case for hybrid multi-objective optimisation in high-stakes machine learning applications. *SIGKDD Explor. Newsl.*, 26(1):24–33, jul 2024. ISSN 1931-0145. doi:10.1145/3682112.3682116. URL `https://doi.org/10.1145/3682112.3682116`.

Jacob Abernethy, Robert E. Schapire, and Umar Syed. Lexicographic optimization: Algorithms and stability, 2024. URL `https://arxiv.org/abs/2405.01387`.

Hieu Pham, Melody Y. Guan, Barret Zoph, Quoc V. Le, and Jeff Dean. Efficient neural architecture search via parameter sharing, 2018. URL `https://arxiv.org/abs/1802.03268`.

Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V. Le. Regularized evolution for image classifier architecture search. *CoRR*, abs/1802.01548, 2018. URL `http://arxiv.org/abs/1802.01548`.

Esteban Real, Sherry Moore, Andrew Selle, Saurabh Saxena, Yutaka I. Leon-Suematsu, Quoc V. Le, and Alex Kurakin. Large-scale evolution of image classifiers. *CoRR*, abs/1703.01041, 2017. URL `http://arxiv.org/abs/1703.01041`.

Jean-Antoine Désidéri. Multiple-gradient descent algorithm (MGDA) for multiobjective optimization. *Comptes Rendus Mathematique*, 350(5):313–318, 2012. ISSN 1631-073X. doi:https://doi.org/10.1016/j.crma.2012.03.014. URL `https://www.sciencedirect.com/science/article/pii/S1631073X12000738`.

Rhea Sanjay Sukthanker, Arber Zela, Benedikt Staffler, Samuel Dooley, Josif Grabocka, and Frank Hutter. Multi-objective differentiable neural architecture search, 2024. URL `https://arxiv.org/abs/2402.18213`.

Colin White, Willie Neiswanger, and Yash Savani. BANANAS: bayesian optimization with neural architectures for neural architecture search. *CoRR*, abs/1910.11858, 2019. URL `http://arxiv.org/abs/1910.11858`.

Gresa Shala, Thomas Elsken, Frank Hutter, and Josif Grabocka. Transfer NAS with meta-learned bayesian surrogates. In *The Eleventh International Conference on Learning Representations*, 2023. URL `https://openreview.net/forum?id=paGvsrl4Ntr`.

Tianqi Chen, Thierry Moreau, Ziheng Jiang, Lianmin Zheng, Eddie Yan, Meghan Cowan, Haichen Shen, Leyuan Wang, Yuwei Hu, Luis Ceze, et al. Tvm: An automated end-to-end optimizing compiler for deep learning. *arXiv preprint arXiv:1802.04799*, 2018.

Linley Gwennap. GROQ ROCKS NEURAL NETWORKS. 2020.

Ryan Loney Ria Chruvu. Optimizing large language models with the OpenVINO™ toolkit. URL `https://www.intel.com/content/www/us/en/content-details/817010/optimizing-large-language-models-with-the-openvino-toolkit.html`.

Yuxiao Zhou and Kecheng Yang. Exploring tensorrt to improve real-time inference for deep learning. In *2022 IEEE 24th Int Conf on High Performance Computing & Communications; 8th Int Conf on Data Science & Systems; 20th Int Conf on Smart City; 8th Int Conf on Dependability in Sensor, Cloud & Big Data Systems & Application (HPCC/DSS/SmartCity/DependSys)*, pages 2011–2018, 2022. doi:10.1109/HPCC-DSS-SmartCity-DependSys57074.2022.00299.

Andrew S. Cassidy, John V. Arthur, Filipp Akopyan, Alexander Andreopoulos, Rathinakumar Appuswamy, Pallab Datta, Michael V. Debole, Steven K. Esser, Carlos Ortega Otero, Jun Sawada, Brian Taba, Arnon Amir, Deepika Bablani, Peter J. Carlson, Myron D. Flickner, Rajamohan Gandhasri, Guillaume J. Garreau, Megumi Ito, Jennifer L. Klamo, Jeffrey A. Kusnitz, Nathaniel J. McClatchey, Jeffrey L. McKinstry, Yutaka Nakamura, Tapan K. Nayak, William P. Risk, Kai Schleupen, Ben Shaw, Jay Sivagnaname, Daniel F. Smith, Ignacio Terrizzano, Takanori Ueda, and Dharmendra Modha. 11.4 ibm northpole: An architecture for neural network inference with a 12nm chip. In *2024 IEEE International Solid-State Circuits Conference (ISSCC)*, volume 67, pages 214–215, 2024. doi:10.1109/ISSCC49657.2024.10454451.

Chris Ying, Aaron Klein, Esteban Real, Eric Christiansen, Kevin Murphy, and Frank Hutter. Nas-bench-101: Towards reproducible neural architecture search. *CoRR*, abs/1902.09635, 2019. URL `http://arxiv.org/abs/1902.09635`.

Xuanyi Dong and Yi Yang. Nas-bench-201: Extending the scope of reproducible neural architecture search. *CoRR*, abs/2001.00326, 2020. URL `http://arxiv.org/abs/2001.00326`.

Julien Siems, Lucas Zimmer, Arber Zela, Jovita Lukasik, Margret Keuper, and Frank Hutter. Nas-bench-301 and the case for surrogate benchmarks for neural architecture search. *arXiv preprint arXiv:2008.09777*, 4:14, 2020.

Yash Mehta, Colin White, Arber Zela, Arjun Krishnakumar, Guri Zabergja, Shakiba Moradian, Mahmoud Safari, Kaicheng Yu, and Frank Hutter. Nas-bench-suite: NAS evaluation is (now) surprisingly easy. *CoRR*, abs/2201.13396, 2022. URL `https://arxiv.org/abs/2201.13396`.

Zhiqing Sun, Hongkun Yu, Xiaodan Song, Renjie Liu, Yiming Yang, and Denny Zhou. Mobilebert: a compact task-agnostic bert for resource-limited devices, 2020. URL `https://arxiv.org/abs/2004.02984`.

Daniel Golovin, Benjamin Solnik, Subhodeep Moitra, Greg Kochanski, John Karro, and D. Sculley. Google vizier: A service for black-box optimization. In *Proceedings of the 23rd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, KDD '17, page 1487–1495, New York, NY, USA, 2017. Association for Computing Machinery. ISBN 9781450348874. doi:10.1145/3097983.3098043. URL `https://doi.org/10.1145/3097983.3098043`.