

Spiking Neural Network Architecture Search: A Survey

Kama Svoboda and Tosiron Adegbija, *Senior Member, IEEE*

Abstract—This survey paper presents a comprehensive examination of Spiking Neural Network (SNN) architecture search (SNNaS) from a unique hardware/software co-design perspective. SNNs, inspired by biological neurons, have emerged as a promising approach to neuromorphic computing. They offer significant advantages in terms of power efficiency and real-time resource-constrained processing, making them ideal for edge computing and IoT applications. However, designing optimal SNN architectures poses significant challenges, due to their inherent complexity (e.g., with respect to training) and the interplay between hardware constraints and SNN models. We begin by providing an overview of SNNs, emphasizing their operational principles and key distinctions from traditional artificial neural networks (ANNs). We then provide a brief overview of the state of the art in NAS for ANNs, highlighting the challenges of directly applying these approaches to SNNs. We then survey the state-of-the-art in SNN-specific NAS approaches. Finally, we conclude with insights into future research directions for SNN research, emphasizing the potential of hardware/software co-design in unlocking the full capabilities of SNNs. This survey aims to serve as a valuable resource for researchers and practitioners in the field, offering a holistic view of SNNaS and underscoring the importance of a co-design approach to harness the true potential of neuromorphic computing.

Index Terms—Spiking neural networks, neural architecture search.

I. INTRODUCTION

As the demand for artificial intelligence (AI) in resource-constrained environments grows, the need for energy-efficient neural architectures has become increasingly critical. Spiking neural networks (SNNs), inspired by the brain’s biological processes, have emerged as a promising solution for low-power AI applications, particularly in edge computing and real-time processing scenarios. These networks, which communicate through discrete temporal spikes rather than continuous values, represent a fundamental departure from conventional artificial neural networks (ANNs) and align closely with the principles of neuromorphic computing, an area that has attracted significant research attention for its potential to enable ultra-low-power AI systems.

SNNs offer several potential advantages over ANNs. Spike timing in SNNs captures temporal patterns better than static ANNs, mirroring real-world signals that are inherently temporal, such as audio, visual, and sensory data [1], [2]. Their sparsity ensures energy efficiency, making them ideal for constrained environments [3]. SNNs’ ability to mimick biological neurons enables the potential to derive insight into

The authors are with the Department of Electrical and Computer Engineering, The University of Arizona, USA, email: {ksvoboda, tosiron}@arizona.edu.

brain function, advancing neuroscience and AI while boosting network fault tolerance [4], [5]. The temporal nature and unique structure of SNNs make them efficient, versatile, and closer to natural neural processes than traditional artificial networks [6].

The architecture of an SNN plays a critical role in both its performance and efficiency [7]. However, designing optimal architectures poses unique challenges due to the intricate interplay between network topology, neuronal dynamics, and hardware constraints. Unlike traditional ANNs, SNNs employ discrete spike-based communication, necessitating specialized design and training strategies [8], [9]. Most SNN architectures have followed one of two suboptimal paths: either adapting existing ANN architectures (such as VGG-Net [10] or ResNet [11]) that were not designed for spike-based communication, or relying on expert-driven manual design through time-consuming trial-and-error processes [7], [8], [12]. Neither approach fully capitalizes on the unique characteristics of spiking neurons.

Given these challenges, *Spiking Neural Network Architecture Search (SNNaS)* has emerged as a critical research direction [7], [8], [12]–[28]. SNNaS aims to automate the discovery of optimal network architectures by simultaneously optimizing multiple objectives, including inference accuracy, processing speed, and energy consumption. Importantly, this optimization process must consider both software and hardware constraints, particularly when targeting deployment on neuromorphic hardware [29].

This survey paper provides a comprehensive analysis of SNNaS research through the lens of hardware-software co-design. We begin by establishing fundamental SNN concepts, including spike encoding/decoding mechanisms, neuron models, and training methodologies, while highlighting key distinctions from traditional ANNs. Building on this foundation, we examine how Neural Architecture Search (NAS) approaches from the ANN domain must be adapted for the unique characteristics of SNNs. The survey’s core contribution lies in its systematic exploration of SNNaS methodologies, encompassing search space design, optimization algorithms, acceleration strategies, and performance estimation techniques, with particular emphasis on hardware constraint integration. Finally, looking forward, we identify promising research directions, including surrogate gradient optimization, Bayesian approaches, hybrid methodologies, and emerging neuromorphic hardware platforms. Through this analysis, we aim to illuminate both the challenges and opportunities in developing efficient SNN architectures for next-generation AI systems.

The remainder of this paper is organized as follows: Section

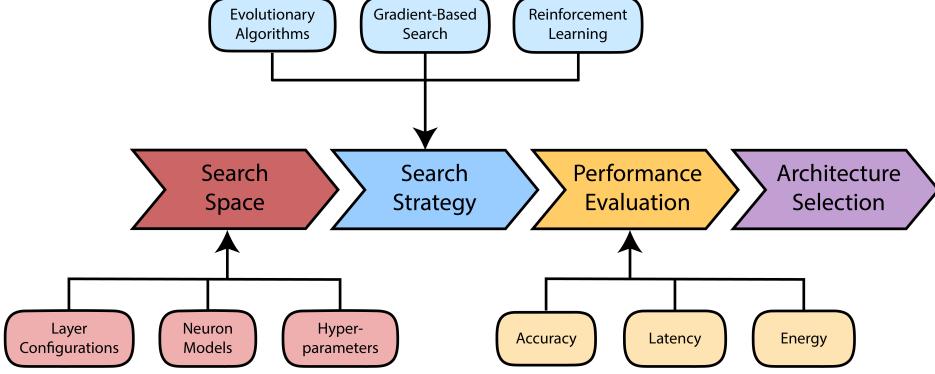


Fig. 1. Overview of the Neural Architecture Search (NAS) workflow for Spiking Neural Networks: From defining the search space (e.g., layer configurations, neuron models, and hyperparameters) to employing search strategies (e.g., evolutionary algorithms, gradient-based search, and reinforcement learning, etc.) for performance evaluation (e.g., accuracy, latency, and energy) and, finally, optimal architecture selection.

II introduces the fundamental concepts of SNNs, including their operational principles, neuron models, spike encoding mechanisms, and training methodologies. Section III examines Neural Architecture Search, beginning with the unique challenges of applying NAS to SNNs, followed by an extensive analysis of search spaces, search strategies, and evaluation methods. We also discuss various acceleration techniques and hardware considerations in this section. Section IV explores current challenges and future research directions, including scalability issues, energy efficiency optimization, and the development of more sophisticated training algorithms. Finally, Section V concludes by summarizing the key insights and emphasizing the importance of hardware/software co-design in advancing SNNaS research.

II. BACKGROUND

Spiking Neural Networks (SNNs) mimic biological nervous systems, using spiking neurons to transmit data through precisely timed electrical pulses [30]. This temporal element provides fresh perspectives on brain dynamics and neural network efficiency. SNNs excel at processing time-dependent patterns, making them ideal for complex recognition tasks [31]. By emulating nature's design, these networks offer a powerful approach to artificial intelligence that leverages the intricacies of biological information processing.

A. Key Components Affecting Architecture Search

1) *Neuron Models*: SNNs employ diverse neuron models to emulate biological neurons, presenting different trade-offs for architecture search. For example, the Leaky Integrate-and-Fire (LIF) model [30], [32] is a widely used model that describes neuron behavior through membrane potential dynamics and spike generation, balancing computational efficiency with biological plausibility. The Hodgkin-Huxley model [33], on the other hand, offers high biological accuracy but demands significant computational resources due to its granular modeling of ion channel kinetics and membrane potentials. The Izhikevich

model [34] strikes a balance between the biological plausibility of the Hodgkin-Huxley model and the computational efficiency of the LIF model [35]. It abstracts the low-level details present in the Hodgkin-Huxley models while maintaining key dynamic properties.

Each model presents unique advantages, allowing researchers to select the most appropriate approach for their specific SNN applications, often balancing the trade-off between biological accuracy and computational efficiency. For example, while the precision of the Hodgkin-Huxley model makes it ideal for studies requiring biophysical detail, its computational intensity can limit scalability. In contrast, the Izhikevich model can replicate a wide range of spike behaviors with significantly lower computational demands. However, this efficiency comes at the expense of detailed biophysical insights, which may be critical for applications requiring fine-grained neuronal dynamics.

2) *Spike Encoding*: Information encoding significantly impacts architecture performance. SNNs encode information using discrete neuronal events or spike trains, in which precise spike patterns and timing are used to represent data.

Various encoding algorithms transform input data into spike trains, each suited to specific applications and hardware constraints [36], [37]. Rate-based encoding, for example, uses the firing frequency to represent signal magnitude [38]. Temporal encoding methods focus on the timing of spikes, offering different strategies: time-to-first-spike emits earlier spikes for larger signals [39]; phase coding periodically modulates spike weights [40]; and burst encoding utilizes inter-spike intervals and bursts to convey information [41].

Spatio-temporal encoding captures both spike timing and spatial distribution across neurons or layers, enhancing the representation of neural dynamics and supporting complex pattern learning [42]. Population encoding further extends this by involving neuron groups to represent single inputs, with each neuron responding to distinct signal features to enable robust and distributed representations.

The choice of encoding scheme significantly influences the efficiency of information transfer and the performance of both hardware and software. Spatio-temporal encoding offers high efficiency for tasks requiring dynamic interaction between neurons [42], while temporal pattern coding can reduce power consumption in computationally demanding tasks [43]. For instance, in ImageNet classification, temporal pattern coding achieved up to $35\times$ reduction compared to rate-coded SNNs and $42\times$ versus traditional ANNs, with only a minimal increase in error [43]. These trade-offs highlight the importance of aligning encoding strategies with the specific requirements of both the application and the underlying hardware.

3) *Training Approaches*: Synaptic plasticity underpins learning in biological brains, with connection strengths changing over time [44]. SNNs aim to mimic this but face unique training challenges. The non-differentiable spike function complicates traditional backpropagation methods [45]. Furthermore, training deep SNNs with multi-layer learning is particularly challenging, often resulting in learning being confined to a single layer in many SNNs [46]–[48]. While progress has been made with two-layer [49]–[51] and recurrent SNNs [52], [53], training deeper networks remains difficult. Hidden layers are crucial for tackling complex real-world problems, making this a key hurdle to overcome [54].

Three primary methods are used for training SNNs and influence architecture design: *surrogate gradients*, *ANN-to-SNN conversion*, and *spike-timing-dependent plasticity (STDP)*.

Surrogate gradients enable SNN training by approximating non-differentiable spike functions, allowing backpropagation while preserving temporal dynamics [54]. Building on initial work [55], [56], various backpropagation through time (BPTT) methods have demonstrated effectiveness in image classification [54], [57], [58]. While these approaches achieve high accuracy with low latency, their resource intensity limits application to simpler architectures [59]. The optimal choice of surrogate function must balance accuracy and implementation efficiency. Although software implementations offer flexibility [54], hardware considerations impact spiking activity sparsity and implementation efficiency [60].

ANN-to-SNN Conversion (or Shadow Training) transforms ANNs into SNNs by training a ‘shadow’ ANN, then replacing analog neurons with spiking neurons, while retaining weights and reconfiguring neurons to mimic original network. This bypasses the issue of nondifferentiable spiking activity and reduces training overhead [61]. The spike rate or latency codes approximate the ANN’s behavior, but this conversion poses challenges such as requiring lengthy simulations that can negate SNNs’ power and speed advantages [62]. Moreover, the approximation process may lead to performance degradation, especially as latency is decreased [63], [64]. Additionally, since ANNs are often trained on static datasets like MNIST and CIFAR-10, the converted SNNs may underutilize temporal dynamics, limiting their potential advantages [3].

Spike-Timing Dependent Plasticity (STDP) [65] adjusts synapse strength based on spike timing, acting as an unsupervised learning rule for SNNs. Hebbian STDP strengthens synapses when pre-synaptic neurons fire before post-synaptic neurons and weakens them otherwise, while anti-Hebbian

STDP inverts this pattern. The proximity of spike pairs determines weight adjustments, making STDP sensitive to order and timing. While STDP naturally adapts to temporal pattern learning, its unsupervised nature complicates network control and does not guarantee optimal performance, especially in complex tasks [66]. Its effectiveness diminishes in deeper networks, limiting its scalability in multilayered SNNs [67], and STDP-trained networks often underperform other learning algorithms [68].

B. Challenges of SNNs

Table I illustrates the fundamental paradox of SNNs: while offering compelling advantages like energy efficiency and temporal processing, each benefit introduces corresponding challenges that must be carefully balanced in the design process. First, as previously alluded to, their non-differentiable spiking mechanism complicates the training process, making standard techniques like backpropagation unsuitable without adaptations [3]. While techniques like surrogate gradients have become a workhorse for training SNNs, they remain less mature and efficient compared to those used in traditional ANNs.

SNNs can be computationally costly due to their temporal dynamics. Unlike ANNs, which compute activations in a single pass, SNNs require the simulation of multiple time steps to capture spike-based neuronal dynamics, increasing both computational cost and latency [69]. This complexity becomes more pronounced when scaling SNNs to large networks or datasets.

Additionally, SNN research suffers from a lack of standardized tools, models, and best practices. Unlike ANNs, which benefit from mature ecosystems such as TensorFlow and PyTorch, SNNs rely on emerging tools like NEST, BindsNET, Brian, and SNNtorch, which are still evolving [6], [70]. This disparity hinders widespread adoption and slows down the development of optimized solutions. Similarly, SNNs often fail to match ANN-level accuracies on traditional benchmarks like ImageNet and CIFAR, due to the nature of these datasets, which are acquired using frame-based sensors rather than event-driven cameras [5], [6]. This discrepancy highlights the need for benchmarks that better reflect the strengths of SNNs, such as temporal and event-driven tasks.

Key issues such as identifying optimal SNN designs, understanding the impact of spike encoding and decoding, and scaling SNNs to handle larger and more complex datasets remain open challenges [71]. Addressing these challenges requires systematic exploration of the design space to develop efficient and scalable architectures. Given the computational constraints and temporal dynamics of SNNs, hardware-software co-design plays a critical role in ensuring that these architectures perform effectively on neuromorphic hardware [37]. This necessitates the studies on Neural Architecture Search (NAS), which has emerged as a promising approach for automating the discovery of optimal SNN designs. By integrating NAS with hardware-aware techniques, researchers can identify architectures that balance accuracy, efficiency, and resource constraints, unlocking the full potential of SNNs for real-world applications.

TABLE I
KEY ADVANTAGES AND CHALLENGES OF SNNs

Feature	Advantages	Challenges
Energy Efficiency	Highly energy-efficient due to sparse and event-driven computation.	Simulation of multiple timesteps increases computational complexity.
Temporal Dynamics	Excellent for time-dependent tasks (e.g., audio or sensory input).	Lack of standardized tools and methods for representing temporal data.
Biological Plausibility	Mimics biological neurons, offering insights into neuroscience and fault-tolerant systems.	Training non-differentiable spike-based mechanisms remains challenging.
Real-Time Processing	Ideal for edge computing and IoT use cases requiring low-latency response.	Event-driven hardware platforms are still emerging and lack widespread adoption.
Scalability	Potential to scale with neuromorphic hardware for efficient parallel processing.	Large datasets and deep architectures drive up computational costs, limiting scalability.
Training	Surrogate gradient and conversion techniques enable learning despite spiking non-differentiability.	Deep multi-layer SNNs remain difficult to train effectively, needing further research.
Hardware Integration	Efficiently integrates with neuromorphic chips (e.g., Loihi) for improved hardware utilization.	Neuromorphic hardware often supports only limited neuron and synapse model diversity.
Accuracy	Competitive results on event-driven or sensor-based benchmarks.	Underperforms ANNs on traditional tasks (e.g., MNIST, CIFAR-10) due to sensor mismatches.

III. NEURAL ARCHITECTURE SEARCH

NAS automates neural network design by treating it as an optimization problem, seeking architectures that maximize performance—typically measured as validation accuracy—while potentially considering other objective functions like computational efficiency. Instead of manually designing networks, NAS employs optimization techniques to identify optimal structures for a given task. An early NAS approach by Stanley et al. [72] used evolutionary methods, inspiring numerous refinements. Although these advances have benefited ANNs, SNNs pose distinct challenges due to their event-driven processing and temporal dynamics, necessitating specialized architectural considerations. To fully exploit SNNs’ potential, NAS methods must be adapted accordingly. This section explores why traditional NAS techniques designed for ANNs may be ineffective for SNNs and examines existing NAS approaches for SNNs, focusing on the three core components: search space, search strategy, and evaluation [73] (Figure 2).

A. Unique Challenges of Applying NAS to SNNs

SNNs’ temporal dynamics and event-driven processing demand specialized NAS approaches distinct from conventional ANN methods. Several key challenges complicate this adaptation:

First, while standard NAS relies on gradient-based optimization, SNNs’ discrete spike functions are inherently non-differentiable. Though surrogate gradient methods exist, they introduce approximation errors and computational overhead that impact large-scale architecture searches [74]. The prevalence of spike-based learning rules like STDP, distinct from gradient-based approaches, further complicates traditional NAS application.

Second, SNNs require specialized performance metrics beyond standard ANN measures. These include spike timing precision, firing rates, and temporal dynamics [3], necessitating evaluation functions that specifically address spiking behavior.

Third, neuromorphic hardware considerations introduce unique constraints such as limited precision and asynchronous computation [75]. Traditional NAS methods, which do not account for these constraints, may produce architectures incompatible with neuromorphic platforms.

Addressing these challenges requires specialized techniques that incorporate spiking neuron models and suitable network topologies [14], [24]. Solutions must balance temporal coding performance, energy efficiency, and hardware compatibility [12], [24], while employing training mechanisms adapted for spike-based computation [8], [14].

B. Search Space

The search space for a NAS method is the collection of feasible network configurations, typically defined by parameters such as the number and type of layers, and the layer hyperparameters (e.g., kernel size, filters, or number of neurons). Because these dimensions combine to yield billions of potential architectures, the search space can be prohibitively large and often becomes the primary bottleneck for architecture search performance [76]. Moreover, different architectural choices can significantly impact hardware accelerators, affecting power efficiency, memory usage, and throughput [37], [60]. Consequently, it is important to design the search space to account for practical hardware constraints. Balancing performance and hardware friendliness ensures that the discovered architectures are accurate and can be efficiently implemented in hardware.

The search can be simplified, and the search space reduced by incorporating prior knowledge about well-performing architectural patterns. However, doing so may introduce human bias that undermines the discovery of novel solutions. On the other hand, overly constraining the search space might exclude high-performing architectures, resulting in suboptimal performance. Hence, selecting an appropriate search space remains a critical challenge in SNNaS, requiring consideration of both algorithmic and hardware factors.

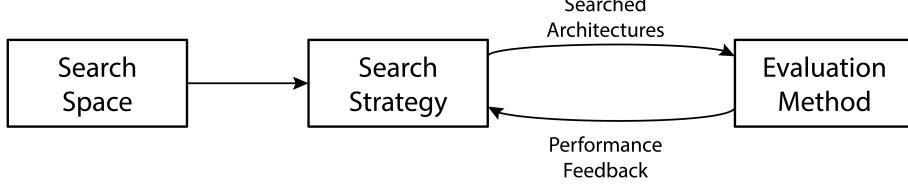


Fig. 2. Neural Architecture Search has three important components: search space, search strategy, and evaluation [73].

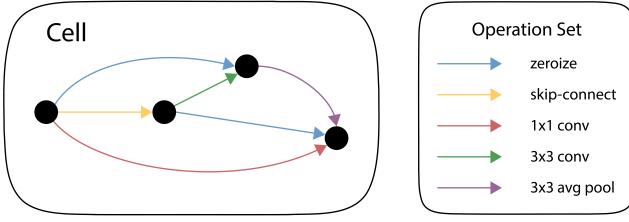


Fig. 3. Example of a cell with four nodes and five possible operations. The cell is a directed acyclic graph, with each edge representing an operation [80].

Current NAS methods use a variety of search spaces.

a) Global Search Space: Global search spaces consider the entire network configuration at once. They allow for tremendous diversity but can be extremely large. For instance, MixedSNN [19] searches among different spiking neuron models and thresholds, and Auto-Spikformer [17] combines Transformer and SNN parameters. Methods such as Zhou et al. [13], Liu and Yi [20], Yan et al. [24], AutoST [28], GAQ-SNN [16], and ANAS [77] illustrate the complexity of global search—spanning from trillions of spiking configurations to hardware-related parameters like processing elements and buffer sizes.

b) Sequential Search Space: Inspired by early NAS works [78], [79], sequential designs treat a network as a layer-by-layer sequence of operations (e.g., convolution parameters). While powerful, these methods can require enormous compute resources. Sequential spaces remain largely unexplored in SNNaS..

c) Cell-Based Search Space: Cell-based approaches define networks at macro and micro levels. Macro-architectures arrange repeated cells, while micro-architectures define operations within each cell. SNASNet [8], based on NAS-Bench-201 [80], pioneered cell-based SNN search using directed acyclic graphs (DAGs) of operations (Zeroize, Skip-Con, 1×1Conv, 3×3Conv, 3×3AvgPool). Extensions include SSTNAS [81], SpikeNAS [12], larger-kernel exploration [21], hardware-friendly constraints [22], and variations in cell design [14]. SpikeDHS [15], [82] also employs DAG-based cells. An example cell-based diagram is shown in Figure 3.

d) Hierarchical Structure Search Space: Hierarchical NAS arranges architectures in multi-level motifs, mirroring modular designs like VGGNet [83] or ResNet [11]. MSE-NAS [23] organizes SNNs into microscopic, mesoscopic, and macroscopic levels; AutoSNN [7] and ANCoEF [84] combine macro-level backbones with micro-level spiking

blocks. SpikeDHS [15], [82] extends from cell-level to layer-level search, and LitE-SNN [25] adds pruning and mixed-precision quantization. SpikeExplorer [26] employs Bayesian optimization for multi-objective search. DESpine [27] and NeuEvo [18] further emphasize biologically inspired motifs, excitatory/inhibitory neurons, and feedback connections. Figure 4 illustrates a typical three-level hierarchy.

e) Memory Bank Representation Search Space: Instead of conventional graphs, SMASH [86] encodes a network as a memory bank system where each operation reads/writes from distinct memory blocks. Figure 5 contrasts single-branch and branching designs. Although SMASH’s representation reduces training overhead, it remains unexplored in SNNaS.

f) Knowledge Distillation: Distillation [87] transfers “teacher” model knowledge into “student” networks (Figure 6). While not yet widely used for SNNaS, SAKD (self-architectural knowledge distillation) [88] shows promise by matching spiking and ANN features under ultra-low latency. Distillation can reduce training costs and improve student generalization, though it adds complexity by requiring balance between teacher–student dynamics.

g) Hardware Search Space: As network architectures directly affect implementation efficiency (e.g., power, memory), hardware aspects must be considered early in SNNaS. Parameter-based frameworks [84], [90], [91] adjust PEs, buffer sizes, and interfaces, whereas template-based methods [92] reconfigure known hardware design templates. ANCoEF [84] highlights that incorrect neuron-to-PE ratios can waste bits in lookup tables, weight SRAM, and network-on-chip flits. Benmeziane et al. [93] classify hardware into server-grade CPUs/GPUs/FPGAs/ASICs, mobile devices, and tiny devices—each requiring distinct design optimizations.

C. Hardware-Aware NAS

Hardware-aware NAS (HW-NAS) seeks to optimize neural network architectures for specific hardware platforms, balancing traditional performance metrics (e.g., accuracy) with constraints like energy, latency, memory usage, and computational requirements [93], [94]. Importantly, temporal and event-driven SNN computations behave differently across hardware implementations [6], so even theoretically optimal SNNs may be impractical if they fail to align with platform-level design choices. As summarized in Table II, HW-NAS ensures that the final architectures exploit hardware features—like event-driven processing or sparse computation [95]—to achieve substantial efficiency gains.

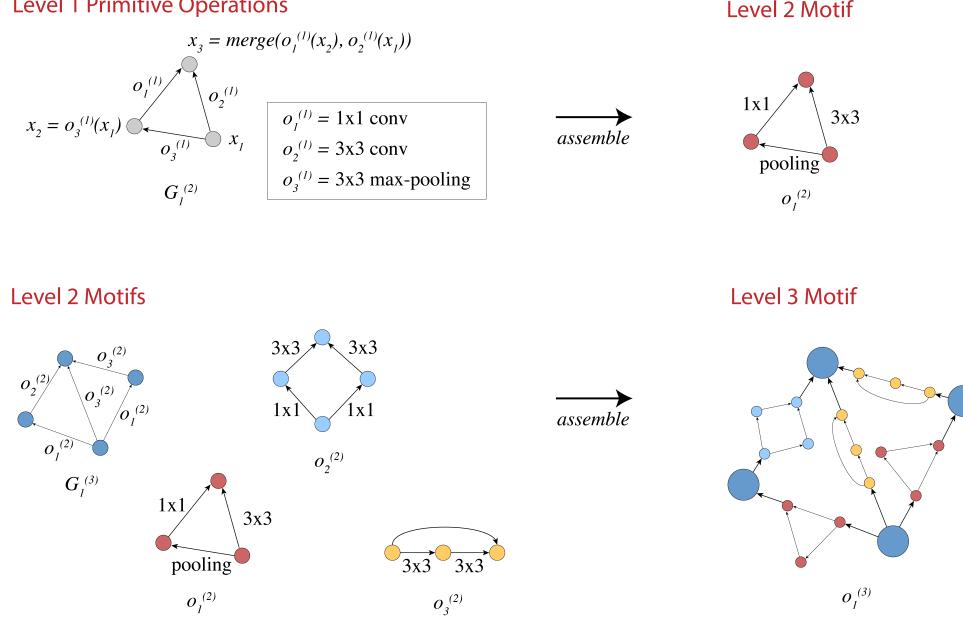


Fig. 4. Example of a three-level hierarchy taken from Liu et al. [85]. Level-1 primitive operations $o_1^{(1)}$, $o_2^{(1)}$, and $o_3^{(1)}$, representing 1×1 convolution, 3×3 convolution, and 3×3 max-pooling respectively, are assembled into a level-2 motif, $o_1^{(2)}$ (TOP). Level-2 motifs $o_1^{(2)}$, $o_2^{(2)}$, and $o_3^{(2)}$ are then assembled into a level-3 motif $o_1^{(3)}$ (BOTTOM).

TABLE II
HARDWARE CONSTRAINTS AND OPTIMIZATION TECHNIQUES

Hardware Constraint	Description	Examples of Optimization Techniques
Memory Usage	Refers to the amount of memory required to store model parameters and operations.	Mixed-precision quantization (e.g., LitE-SNN [25]); Memory-aware kernel size exploration (e.g., SpikeNAS [12]); Weight pruning and compression (e.g., GAQ-SNN [16]).
Latency	Time taken to complete inference or training, crucial for real-time applications.	Hardware-specific optimization (e.g., HASNAS [22], ANCoEF [84]; Timestep compression in SNNs (e.g., LitE-SNN [25]))
Energy Consumption	The power required to perform computations, important for low-power devices.	Minimizing spike counts (e.g., NeuEvo [18], AutoSNN [7]); Synaptic operation (SynOps)-based evaluation (e.g., Auto-Spikformer [17]); Sparse computation using event-driven architectures.
Hardware Area	The physical size of the hardware components required for implementation.	Compute-in-memory architectures (e.g., HASNAS [22]); Optimization of neuron and synapse allocations.
Hardware Compatibility	Ensures the designed SNN can efficiently operate on specific hardware.	Co-exploration of architecture and hardware (e.g., ANAS [77], ANCoEF [84]); Resource-balancing algorithms for neuromorphic platforms (e.g., ANCoEF [84]).
Bandwidth	Data transfer capacity between components, critical in parallel architectures.	Data partitioning and mapping (e.g., ANCoEF [84]); Efficient routing and scheduling strategies.

A variety of accelerators support SNNs, ranging from GPGPU systems [96]–[98] and FPGA-based designs [99]–[104] to neuromorphic platforms like Loihi [105] and TrueNorth [106], mixed-signal VLSI [107], [108], and nanotechnology [109]. Although neuromorphic hardware can be highly energy-efficient, it often supports only a subset of neuron and synaptic models; more flexible general-purpose computing systems can simulate diverse designs but may consume more resources [110]. Thus, HW-NAS must consider whether the targeted hardware is specialized (e.g., tailored to a certain neuron model) or general-purpose (e.g., a GPU).

In general, there are two main strategies for HW-NAS:

single-objective and multi-objective optimization [93]. Single-objective methods either use a two-stage process—first searching for accuracy, then optimizing for hardware—or incorporate hardware limits as constraints from the onset. Multi-objective frameworks, such as SpikeExplorer [26] and LitE-SNN [25], balance accuracy with hardware metrics (e.g., power, area, and memory). Many SNNAs methods rely on Synaptic Operations (SynOps) to model computational and energy costs [16], [17], [24], [25], while spike-count minimization (e.g., [7], [18], [23]) reduces power. Quantization [16] and kernel-size optimization [21] further shrink memory footprints. HW-NAS methods like HASNAS [22], NeuEvo [18], GAQ-SNN [16],

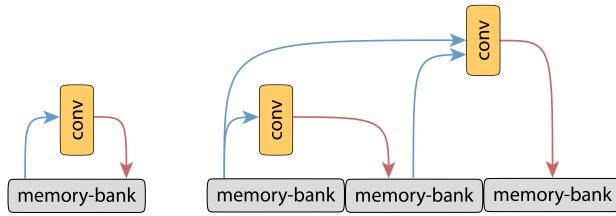


Fig. 5. Example of a memory-bank representation taken from Brock et al. [86]. In a single-branch architecture (LEFT), the network utilizes a singular, large memory bank that it both reads and overwrites at each operation. Branching architectures (RIGHT) read from all previously written banks and write to an empty bank.

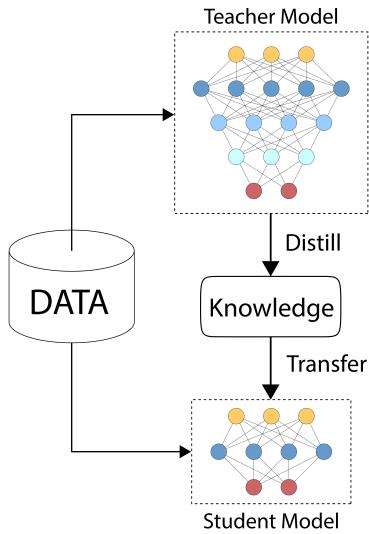


Fig. 6. Example of knowledge distillation [89].

and Auto-Spikformer [17] show that carefully coupling architecture search with hardware-specific constraints yields high accuracy alongside notable improvements in energy and resource efficiency.

D. Co-Exploring Neural Architecture Search Space and Hardware Design Space

Co-exploration simultaneously searches both the neural architecture and the hardware design, going beyond standard HW-NAS frameworks by optimizing every dimension of the design jointly [111]. Unlike HW-NAS—which assumes a fixed hardware platform—co-exploration may consider different or new hardware architectures to ensure designs are feasible, efficient, and well-aligned with real-world deployment constraints [94]. This approach typically delivers better performance and lower costs than separately optimizing hardware or SNNs, but at the expense of greater search complexity, substantial GPU hours, and a large carbon footprint [112].

ANAS [77] exemplifies co-exploration by using an evolutionary algorithm to search a comprehensive design space, evaluating hardware performance with the CanMore simulator. Compared to random or grid search, ANAS achieves up to an

order-of-magnitude improvement in energy-delay and is four orders of magnitude faster to search. Similarly, ANCoEF [84] combines an RL-based multi-objective hardware search with a supernet-based SNN search [7], partially training candidate architectures and discarding those that fail hardware constraints. This iterative process yields a $1.81\times$ reduction in energy-delay product (EDP) and a $2.73\times$ speedup over ANAS [77] on the N-MNIST dataset, highlighting the promise of fully joint NAS and hardware design optimization.

E. Neural Network Architecture Search Strategies

Neural architecture search involves exploring a vast design space to identify high-performing configurations for a specific task while minimizing computational resources and time. A successful strategy must evaluate accuracy (often by training on a reduced dataset subset), estimate hardware cost (via real-world measurement or modeling), and select an appropriate search algorithm (e.g., grid search, random search, evolutionary algorithms, reinforcement learning, Bayesian optimization, or gradient-based methods) [93]. Here we briefly describe different popular NAS strategies, with an emphasis on their applicability to SNNAS, as summarized in Table III.

1) *Grid Search*: Grid search systematically explores every combination of layer counts, types, and hyperparameters to identify an optimal architecture. Although simple and effective for small or constrained search spaces [115], it quickly becomes impractical when the number of configurations grows exponentially. Constrained grid search narrows the range of hyperparameters to control computational cost, but advanced methods (e.g., evolutionary algorithms, reinforcement learning) are typically used when the search space is large. In SNNAS, HASNAS [22] demonstrates how exhaustive methods can be adapted by removing low-impact SNN operations and introducing hardware constraints (e.g., memory, latency, area, energy) prior to the search. This strategy yields an $11\times$ speedup compared to SNASNet [8] on CIFAR10 and CIFAR100 while meeting strict hardware limits and maintaining accuracy, underlining grid search's viability when the space can be reduced effectively.

2) *Random Search*: Random search samples configurations from the architecture space at random, forgoing exhaustive exploration and often discovering novel designs missed by more structured methods [116]–[119]. Although it does not guarantee quick convergence on optimal solutions, it can yield strong results for SNNs when combined with additional constraints or metrics. For example, Kim et al. [8] introduced a training-free measure of spike activation diversity to identify promising random SNN topologies, resulting in SNASNet. Putra and Shafique [21] adapted this approach to construct memory-efficient networks, showing that larger kernels boost accuracy at the expense of increased memory usage. In another study, Liu et al. [20] randomly sampled 1,000 architectures for the Loihi chip, selected 20 based on gradient-based metrics, then fully trained and evaluated these candidates. Their final architectures achieved competitive accuracy, lower energy consumption, and compact designs compared to state-of-the-art SNNs, highlighting the versatility of random search when guided by well-chosen criteria or hardware-specific goals.

TABLE III
ARCHITECTURE SEARCH STRATEGIES

Search Method	Works	Key Characteristics	Strengths	Limitations
Grid Search	HASNAS [22]	Systematically evaluates all hyperparameter combinations.	Parallelizable and simple to implement. Ensures the optimum architecture is found.	Computationally intensive. Impractical for large search spaces.
Random Search	SNASNet [113], Liu & Yi [20], Putra & Shafique [21]	Selects random combinations of hyperparameters from the search space and evaluates the resulting architecture performances.	Explores the search space broadly, finds good configurations faster, and often achieves near-optimal solutions with fewer evaluations. Simple to implement and parallelizable.	Lacks guided exploration. Does not adaptively focus on promising regions of the search space.
Evolutionary	Zhou et al. [13], AutoSNN [7], GAQ [16], Auto-Spikformer [17], NeuEvo [114], MixedSNN [19], ANAS [77], MSE-NAS [23], DESPINE [27], AutoST [28], SSTNAS [81]	Iteratively selects, changes, and reproduces neural networks to optimize a fitness value related to an objective like classification accuracy.	Does not require differentiable fitness functions. Suitable for multi-objective optimization.	Computationally expensive with more parameters.
Reinforcement Learning	ANCoEF [84]	Rewards positive behaviors and penalizes negative ones to derive efficient strategies toward a given goal. Balances exploration and exploitation.	Good for complex and high-dimensional search spaces. Custom reward functions align search with specific goals.	Computationally intensive.
Bayesian Optimization	SpikeExplorer [26]	Finds the optimum of expensive and poorly understood functions using a probabilistic model-based approach.	Requires fewer evaluations. Capable of finding global optima in complex, multimodal search spaces. Learns from prior evaluations to focus on promising regions.	Not well-suited for high-dimensional spaces. Computationally intensive and difficult to parallelize. Can converge to local optima in noisy settings.
Gradient-Based	SpikeDHS [15], Yan et al. [24], LitE-SNN [25], Che et al. [82]	Allows for continuous relaxation of the architecture space, enabling gradient-based optimization.	Fast exploration, scalable to large spaces. Efficient for complex architecture searches.	May under-explore the search space. Susceptible to local optima.
Curriculum	–	Begins with a simple neural network and progressively increases complexity by sampling and evaluating new candidate architectures.	Efficient and effective across different data types. Increases the likelihood of discovering optimal architectures.	Dependent on data quality. Implementation is complex. Limited generalization.
Progressive	Nijenhuis [14], Spike-NAS [12]	Gradually expands the search space, starting with simpler architectures.	Parallelizable and less computationally expensive. Handles multiple resource constraints. Improved search performance.	Shallow initial architectures may not reflect deeper architecture performance. Costly for large search spaces.

3) *Evolutionary Methods*: Evolutionary NAS methods (Figure 7) are rooted in early works [120]–[122] that mimic biological evolution by iteratively selecting, mutating, and recombining neural network configurations to optimize a fitness function—often accuracy, energy usage, or other objectives [123]. Unlike gradient-driven approaches, they do not require a differentiable objective, making them suitable for multi-objective and discrete design challenges. Genetic Algorithms (GAs) are a common technique that encode networks as strings or graphs, then repeatedly apply crossover and mutation to discover high-performing architectures. Historically, these methods have optimized weights (e.g., NEAT [72]), and more recent works extend them to structural and hyperparameter searches [124]–[126], sometimes using illumination algorithms like

MAP-Elites [127] for diverse solutions.

Despite high computational demands (e.g., AmoebaNet-A [128] required 3,150 GPU days, evolutionary algorithms can be accelerated through smaller training sets, prediction models, Lamarckian strategies [129], [130], or partial weight inheritance. Notable ANN HW-NAS efforts include TinyNAS [131], Once-for-all [132], HAT [133], NASCaps [134], and HURRICANE [135]. For SNNs, evolutionary approaches are similarly powerful for tasks requiring hardware, memory, and energy optimizations, as demonstrated by GAQ-SNN [16], DESPINE [27], MixedSNN [19], MSE-NAS [23], Auto-Spikformer [17], NeuEvo [18], and ANAS [77], among others. Their parallelizable, global search can reveal novel SNN designs in scenarios where gradient signals are limited, although balancing computation, efficient encoding, and hyperparameter

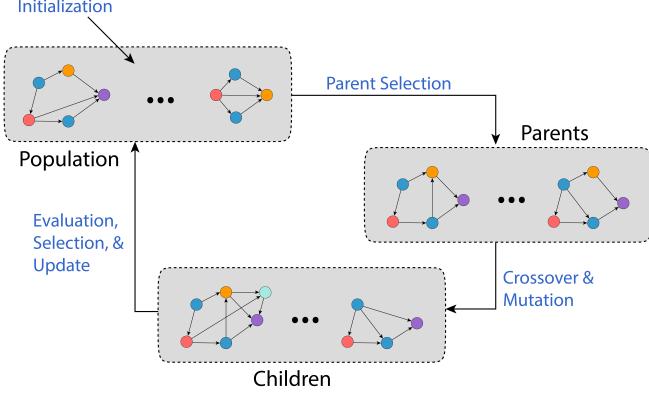


Fig. 7. Evolutionary Algorithm for NAS based on [136]. After initializing the population, the algorithm selects parents and then crosses or mutates them to generate children. It evaluates the adaptability, or the neural architecture performance, of the children and selects a group of individual neural architectures with the best performance from the children generated.

tuning remains a key challenge.

4) *Reinforcement Learning Methods*: Reinforcement learning (RL) methods, illustrated in Figure 8, frame NAS as a Markov Decision Process [137], where an agent explores network architectures (or hardware designs) by taking actions that yield rewards based on accuracy or efficiency. Early work by Zoph and Le [78] introduced RL-based controllers for NAS, but required extensive computational resources; subsequent techniques such as weight sharing [76] and learning curve prediction [138] reduced this overhead. RL is appealing for hardware-oriented NAS because controllers can dynamically adjust to multi-objective goals (e.g., latency, energy) [90], [94], [111], [139], [140].

In the SNNaS context, ANCoEF [84] demonstrates a multi-objective RL approach that simultaneously searches SNN architectures (via a supernet-based method [7]) and asynchronous hardware parameters (e.g., neuron connectivity, mapping, partitioning, arbitration). Its RL agent selects actions to maximize a reward reflecting accuracy, latency, energy consumption, and area utilization. By effectively modeling non-numerical design options (e.g., routing strategies) as discrete actions, ANCoEF outperforms the evolutionary algorithm-based ANAS [77] on N-MNIST with a $1.81\times$ lower energy-delay product (EDP) and $2.73\times$ less search time.

Despite these advances, RL in SNNaS remains relatively unexplored. Challenges include high computational demands, large action spaces, and potential local optima [128], [141], [142]. RL controllers can be sample-inefficient and may require extensive tuning to generalize across tasks. Future directions involve more robust environment designs, refined reward functions, and possible hybridization with evolutionary or Bayesian methods to mitigate sample intensity and improve scalability.

5) *Bayesian Optimization*: Bayesian Optimization (BO) adaptively searches expensive, complex objective functions by building a probabilistic surrogate model (e.g., a Gaussian process) and balancing exploration and exploitation via an acquisition function [143]. In SNNaS, BO typically encodes

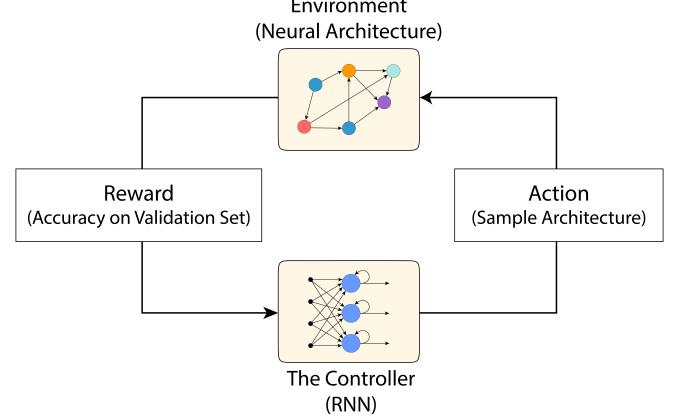


Fig. 8. Illustration of reinforcement learning for NAS based on [78].

architectural elements (e.g., neuron types, connectivity) and performance metrics (e.g., accuracy, sparsity) into a surrogate model that guides the next candidate selection, thereby minimizing costly simulation and training runs. The approach iteratively updates the surrogate model with newly evaluated architectures, refining predictions and uncertainties; careful tuning of hyperparameters, surrogate design, and acquisition functions remains essential. Despite potential computational overhead and risks of local optima, BO's data-efficient search can yield high-quality solutions, as shown by SpikeExplorer [26], which uses BO to co-optimize power, latency, area, and accuracy on FPGAs. This enables efficient exploration across multiple constraints, ultimately outperforming several accelerator baselines while preserving accuracy.

6) *Gradient-Based Methods*: Gradient-based NAS methods optimize architecture parameters via backpropagation through a continuous relaxation of discrete choices [144]. Early examples include DARTS [144], SNAS [113], ProxylessNAS [145], and FBNet [146], which reduced the memory overhead and search time by restricting continuous activation paths or using stochastic gradient updates. However, some studies, such as Bingham et al. [147], report that gradient descent may stagnate on a single architecture over many iterations, failing to exploit the broader search space as effectively as random or evolutionary methods. Despite these challenges, gradient-based strategies remain popular due to their efficiency once the continuous relaxation is set up, and they are frequently combined with techniques like inference latency constraints, as in HTAS [148].

For SNNs, gradient-based methods have spawned frameworks like SpikeDHS [15], which adapts DARTS for spiking neurons by addressing surrogate gradient noise and maintaining multiplication-free inference. Che et al. [82] further enhance SpikeDHS with differentiable surrogate gradient search (DGS) and temporal parameter search (TPS) to refine neuron dynamics, and these modules can be applied individually or jointly for added flexibility. LitE-SNN [25] extends SpikeDHS to incorporate spatial and temporal compression, and Yan et al. [24] propose a branchless supernet that encodes architectures with continuous parameters for gradient-based multi-

objective optimization. Collectively, these approaches have achieved state-of-the-art performance across various datasets by carefully balancing the high computational demands of differentiable searches with specialized SNN design considerations.

7) Progressive Neural Architecture Search: Progressive Neural Architecture Search (PNAS) [149] incrementally expands a base network by adding layers or connections, focusing computational effort on promising partial architectures. This “search by growing” strategy drastically reduces the exploration of unproductive designs and often delivers results significantly faster. PNAS can be five times faster than earlier RL-based methods [150] and calculate the results eight times faster overall [149]. Despite its efficiency and strong performance on benchmarks like CIFAR-10 and ImageNet, PNAS can prematurely prune architectures and incur high training costs for performance prediction [151]. In SNNaS, adaptations include the Sequential Greedy Architecture Search (SGAS) used by Nijenhuis [14], which did not surpass baseline SNN results, and SpikeNAS [12], a memory-aware PNAS variant that progressively refines cell-based SNN designs by pruning memory-intensive operations. These studies illustrate both the potential and the challenges of applying PNAS to spike-based architectures.

8) Curriculum Neural Architecture Search: CNAS [151] employs a staged, curriculum-inspired process to navigate complex NAS landscapes. It begins with a simple base network, then gradually adds operations or layers via mutation or component combination. Newly generated architectures are briefly trained and compared to the base; stronger performers supersede earlier models. By focusing resources on promising partial solutions, CNAS can achieve better convergence within limited budgets, splitting NAS into manageable subproblems that exploit learned knowledge from previous stages.

9) Hybrid Methods: Hybrid approaches blend multiple NAS strategies to combine their benefits, such as robust exploration from one method and efficient exploitation from another, thereby speeding searches and conserving resources. Examples include evolutionary algorithms coupled with reinforcement learning [152]–[154] and random search paired with local grid search. Although these hybrids can produce more balanced results and handle diverse tasks, they also introduce increased complexity and hyperparameter tuning. Careful design is crucial to avoid overlapping strengths and redundant computations, but when done well, hybrid NAS can outperform single-method solutions.

F. Non-differentiable Hardware Constraints

All of the search strategies mentioned above can be adapted to hardware-aware NAS (HW-NAS), which aims to optimize neural architectures to meet hardware-specific constraints such as memory, latency, and energy requirements. However, real-world hardware metrics such as bus widths, memory sizes, or specialized neuromorphic parameters are often discrete, making them non-differentiable. Although HW-NAS typically relies on differentiable loss functions for efficient optimization [93], these discrete choices complicate gradient-based approaches.

To address this challenge, methods like Gumbel-Softmax (used in FBNet [146]) introduce controlled noise to effectively “soften” discrete variables, while ProxylessNAS [145] either approximates the discrete search space with a continuous function or applies the REINFORCE algorithm [155] to handle binarized weights. These strategies facilitate backpropagation by approximating non-differentiable hardware constraints, allowing over-parameterized yet gradient-friendly models to incorporate crucial hardware considerations.

G. Performance Evaluation Metrics

Performance evaluation metrics are essential in SNNaS, as they guide the design of architectures that balance accuracy, latency, energy consumption, memory footprint, and area. While floating point operations (FLOPs) are commonly used, they can be misleading since equal FLOPs may yield varying latencies on different hardware [133], [135], [156]. Accuracy-focused but energy-aware NAS methods, such as MONAS [157] and NetAdapt [158], incorporate power consumption, and circuit area minimization has also been studied to reduce static power use [92]. For edge applications, limited memory budgets are critical, and latency or profiled energy often provide more reliable indicators than FLOPs alone. Balancing these hardware metrics alongside task performance ensures robust solutions for diverse deployment scenarios.

To estimate hardware cost, NAS methods use analytical formulas for quick approximations, albeit with less precision on complex designs; lookup tables (LUTs) leverage pre-measured data (e.g., layer-specific latency and energy in [159]); and simulation-based approaches, such as MnasNet [139], offer detail at the expense of time. Machine learning models predict hardware behavior from architecture-hardware mappings [92], [157], [160], [161], but rely on high-quality data. Hybrid approaches combine these methods. For instance, fast analytical or LUT-driven screening may precede more rigorous simulations or ML estimates for top-performing candidates.

H. Search Acceleration Strategies

Search acceleration strategies in SNNaS aim to mitigate the immense computational costs of training every candidate architecture from scratch [126], [128], [150]. One approach is to use lower fidelity estimates, which evaluate networks under reduced configurations, e.g., fewer epochs, surrogate models, smaller resolutions or channel counts. This approach significantly reduces search time but risks some loss in accuracy [76], [128], [144], [150], [162]–[166]. In SNNaS specifically, Nijenhuis [14] achieved a 45× speed-up by reducing both dataset size and time dimension. HotNAS [111] introduced a “Hot Start” strategy that begins with a strong base model, accelerating convergence. Learning Curve Extrapolation (LCE) [138] (illustrated in Figure 9) cuts computational costs by projecting final performance from partial training. Network morphism and weight inheritance rapidly adapt a trained model by making small architectural changes while retaining weights, as seen in works like Nijenhuis [14], improving efficiency and convergence.

TABLE IV
SPIKING NEURAL NETWORK ARCHITECTURE SEARCH PAPERS

Model	Search Space	Search Algorithm	Acceleration Strategy	Performance Evaluation
Zhou et al. (2020) [13]	Global	Evolutionary	Cheap Surrogate	Accuracy
Nijenhuis (2021) [14]	Cell-based	Progressive NAS	Weight Inheritance, Reduced Dataset, Lower Fidelity Estimates	Accuracy
AutoSNN (2022) [7]	Hierarchical	Evolutionary	One-Shot	Accuracy and Energy Efficiency
SpikeDHS (2022) [15]	Cell-based and Hierarchical	Gradient-Based	First-Order Gradient Approximation	Accuracy, Sparsity, Computation Cost, Energy Consumption
GAQ (2022) [16]	Global	Evolutionary	–	Accuracy and Memory Cost
SNASNet (2022) [8]	Cell-based	Random Search	NAS without Training	Accuracy
Auto-Spikformer (2023) [17]	Global	Evolutionary	One-Shot	Accuracy and Energy Consumption
NeuEvo (2023) [18]	Hierarchical	Evolutionary	–	Accuracy and Sparsity
MixedSNN (2023) [19]	Global	Evolutionary	NAS without Training	Accuracy and PBSE
ANAS (2023) [77]	Global	Evolutionary	Parallel Searching with Multi-Thread CanMore Simulator	Energy-Delay (Latency and Energy Consumption)
Liu & Yi (2024) [20]	Global	Random Search	NAS without Training	GMG (Accuracy), Energy Consumption, Model Parameters
SpikeNAS (2024) [12]	Cell-based	Progressive NAS	NAS without Training	Accuracy and Memory Budget
Putra & Shafique (2024) [21]	Cell-based	Random Search	NAS without Training	Accuracy and Memory Footprint
HASNAS (2024) [22]	Cell-based	Grid Search	NAS without Training	Accuracy, Memory, Area, Latency, and Energy Consumption
MSE-NAS (2024) [23]	Hierarchical	Evolutionary	NAS without Training	BIE (Neural Activity)
Yan et al. (2024) [24]	Global	Gradient-Based	One-Shot	Accuracy and Computational Cost (SynOps)
LitE-SNN (2024) [25]	Hierarchical	Gradient-Based	Weight Sharing	Accuracy, Memory, Computational Cost (SynOps)
SpikeExplorer (2024) [26]	Hierarchical	Bayesian Optimization	User-Set Constraints on Parameters	Accuracy, Area, Latency, Energy Consumption
DESPINE (2024) [27]	Hierarchical	Evolutionary	–	Accuracy, Energy Consumption (Spike Rate)
AutoST (2024) [28]	Global	Evolutionary	NAS without Training	FLOPs
Che et al. (2024) [82]	Cell-based and Hierarchical	Gradient-Based, DGS, TPS	Limited Search Space, First-Order Gradient Approximation, Merging Preprocessing	Accuracy, Sparsity, Computation Cost, Energy Consumption
SSTNAS (2024) [81]	Cell-based	Evolutionary	NAS without Training	Accuracy
ANCoEF (2024) [84]	Hierarchical	Reinforcement Learning	One-Shot	Accuracy, Hardware Performance, Energy Consumption, Area

Other methods reduce full-training overhead even more drastically. One-shot NAS (weight sharing), exemplified by ENAS [76], DARTS [144], and SNAS [113], trains a supernet whose subnets share weights. In the SNN domain, AutoSNN [7], ANCoEF [84], Auto-Spikformer [17], and Yan et al. [24] adopt this strategy to accelerate architecture search without separate training for each candidate. Moving further, zero-shot approaches predict performance without training at all, using metrics like activation overlap [167]–[171]. In SNNaS, Liu et al. [20] estimate Loihi-based architectures using gradient-based Gram matrices, while Kim et al. [8] introduce the Sparsity-Aware Hamming Distance for zero-shot SNN designs. Follow-up works like SpikeNAS [12], HASNAS [22],

and Li et al. [81] leverage similar metrics to achieve orders-of-magnitude faster searches, balancing computational efficiency with competitive accuracy.

I. Temporal Considerations

Temporal dynamics are central to SNNs' unique capabilities, yet many SNNaS approaches adapt ANN-centric methods that overlook time-related factors. Recent work has begun to address this gap. For instance, LitE-SNN [25] integrates both spatial and temporal compression, and MixedSNN [19] leverages Period-Based Spike Evaluation (PBSE) to capture critical temporal attributes. Che et al. [82] further introduce a spatial-temporal search methodology that uses Temporal Parameter

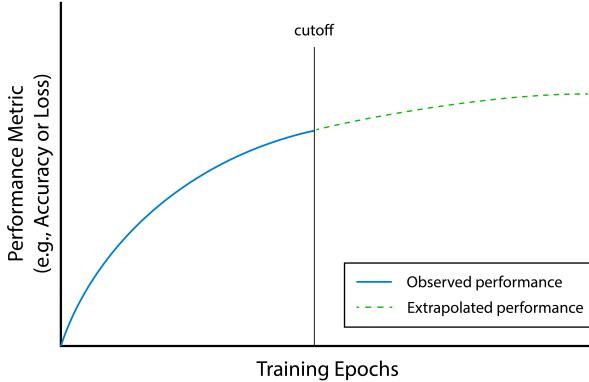


Fig. 9. Learning curve extrapolation illustrating the use of partial training data to predict final performance. The observed curve (solid blue) represents the model’s accuracy during early epochs, while the extrapolated curve (dashed green) predicts performance for the remaining training epochs. This method enables early termination of unpromising models, saving computational resources and time.

Search (TPS) to fine-tune different neurons’ time constants, optimizing both architecture and temporal behaviors. Such strategies signal an important transition toward a more holistic SNN design paradigm, embracing temporal as well as spatial dimensions to enhance neuromorphic computing.

IV. CHALLENGES AND FUTURE DIRECTIONS

Table IV outlines recent SNNaS studies surveyed in this work, highlighting their search spaces, search algorithms, acceleration strategies, and performance criteria. While these approaches have driven notable progress in the field, significant challenges remain regarding flexibility, interpretability, scalability, and broader applicability, especially beyond core tasks like image classification.

SNNaS faces persistent challenges regarding search space flexibility, interpretability, and reproducibility across diverse tasks. Human-designed or cell-based search spaces risk omitting genuinely novel architectures, while many methods remain tailored to specific datasets (often image classification) and do not naturally generalize to broader domains [172], [173]. Moreover, understanding *why* certain automatically discovered architectures perform well remains elusive, underscoring a need for deeper theoretical insights into SNN behaviors [174]. Resource-intensive search processes further restrict accessibility and reproducibility, despite partial solutions like one-shot NAS, weight sharing, or zero-cost proxies [73], [173]. These acceleration techniques can introduce biases—for instance, parameter sharing may degrade final model quality or misrank promising designs [117], [132], [175], [176].

Simultaneously optimizing architectures and hyperparameters could lessen human bias and yield more robust solutions, especially under few-shot learning or meta-learning conditions [174], [177]. There is also growing interest in architectures resilient to adversarial attacks [73], multi-task learning, and multi-objective optimization that accounts for both performance and resource efficiency [73]. In parallel, research on hardware constraints in SNNaS has remained limited, often defaulting to von Neumann architectures despite

memory walls and other performance bottlenecks [93]. Co-optimizing SNN architectures with emerging paradigms—such as in-memory computing, quantization, and pruning—could provide far more bio-plausible and efficient solutions [93], [178], [179]. Finally, most SNNaS methods still focus on a single, universal architecture rather than tailoring designs for specific hardware devices, suggesting future work in building device-targeted or co-explored solutions [93].

V. CONCLUSION

Spiking neural network architecture search (SNNaS) demands a tightly integrated hardware/software co-design strategy, as standard ANN-based NAS methods overlook spiking-specific challenges. Recent progress, encompassing hardware-aware design, explicit temporal considerations, and novel search algorithms such as surrogate gradients, Bayesian optimization, and hybrid techniques, demonstrates the growing sophistication in SNN-focused approaches. However, fundamental issues remain, such as bridging hardware-software integration, reducing computational costs, and improving energy efficiency. Addressing these challenges is essential for scaling SNNs to real-world tasks. Overcoming these challenges and deploying SNNs on dedicated neuromorphic hardware will be crucial to unlocking the full potential of spiking networks and driving neuromorphic computing forward.

REFERENCES

- [1] S.-C. Liu and T. Delbrück, “Neuromorphic sensory systems,” *Current opinion in neurobiology*, vol. 20, no. 3, pp. 288–295, 2010.
- [2] M. Osswald, S.-H. Ieng, R. Benosman, and G. Indiveri, “A spiking neural network model of 3d perception for event-based neuromorphic stereo vision systems,” *Scientific reports*, vol. 7, no. 1, p. 40703, 2017.
- [3] L. Deng, Y. Wu, X. Hu, L. Liang, Y. Ding, G. Li, G. Zhao, P. Li, and Y. Xie, “Rethinking the performance comparison between snns and anns,” *Neural networks*, vol. 121, pp. 294–307, 2020.
- [4] Z. Sun, V. Cutsuridis, C. F. Caiafa, and J. Solé-Casals, “Brain simulation and spiking neural networks,” *Cognitive Computation*, vol. 15, no. 4, pp. 1103–1105, 2023.
- [5] J. D. Nunes, M. Carvalho, D. Carneiro, and J. S. Cardoso, “Spiking neural networks: A survey,” *IEEE Access*, vol. 10, pp. 60738–60764, 2022.
- [6] M. Pfeiffer and T. Pfeil, “Deep learning with spiking neurons: opportunities and challenges,” *Frontiers in neuroscience*, vol. 12, p. 409662, 2018.
- [7] B. Na, J. Mok, S. Park, D. Lee, H. Choe, and S. Yoon, “Autosnn: Towards energy-efficient spiking neural networks,” in *International Conference on Machine Learning*, pp. 16253–16269, PMLR, 2022.
- [8] Y. Kim, Y. Li, H. Park, Y. Venkatesha, and P. Panda, “Neural architecture search for spiking neural networks,” in *European Conference on Computer Vision*, pp. 36–56, Springer, 2022.
- [9] A. Tavanaei, M. Ghodrati, S. R. Kheradpisheh, T. Masquelier, and A. Maida, “Deep learning in spiking neural networks,” *Neural networks*, vol. 111, pp. 47–63, 2019.
- [10] K. Simonyan, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [11] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- [12] R. V. W. Putra and M. Shafique, “Spikenas: A fast memory-aware neural architecture search framework for spiking neural network systems,” *arXiv preprint arXiv:2402.11322*, 2024.
- [13] Y. Zhou, Y. Jin, and J. Ding, “Surrogate-assisted evolutionary search of spiking neural architectures in liquid state machines,” *Neurocomputing*, vol. 406, pp. 12–23, 2020.
- [14] J. Nijenhuis, “Using nas to improve accuracy of snns,” 2021.

- [15] K. Che, L. Leng, K. Zhang, J. Zhang, Q. Meng, J. Cheng, Q. Guo, and J. Liao, "Differentiable hierarchical and surrogate gradient search for spiking neural networks," *Advances in Neural Information Processing Systems*, vol. 35, pp. 24975–24990, 2022.
- [16] D.-A. Nguyen, X.-T. Tran, and F. Iacopi, "Gaq-snn: A genetic algorithm based quantization framework for deep spiking neural networks," in *2022 International Conference on IC Design and Technology (IC-CDT)*, pp. 93–96, 2022.
- [17] K. Che, Z. Zhou, Z. Ma, W. Fang, Y. Chen, S. Shen, L. Yuan, and Y. Tian, "Auto-spikformer: Spikformer architecture search," *arXiv preprint arXiv:2306.00807*, 2023.
- [18] G. Shen, D. Zhao, Y. Dong, and Y. Zeng, "Brain-inspired neural circuit evolution for spiking neural networks," *Proceedings of the National Academy of Sciences*, vol. 120, no. 39, p. e2218173120, 2023.
- [19] Z. Xie, Z. Liu, P. Chen, and J. Zhang, "Efficient spiking neural architecture search with mixed neuron models and variable thresholds," in *International Conference on Neural Information Processing*, pp. 466–481, Springer, 2023.
- [20] S. Liu and Y. Yi, "Unleashing energy-efficiency: Neural architecture search without training for spiking neural networks on loihi chip," in *2024 25th International Symposium on Quality Electronic Design (ISQED)*, pp. 1–7, IEEE, 2024.
- [21] R. V. W. Putra and M. Shafique, "A methodology for improving accuracy of embedded spiking neural networks through kernel size scaling," *arXiv preprint arXiv:2404.01685*, 2024.
- [22] R. Vidya Wicaksana Putra and M. Shafique, "Hasnas: A hardware-aware spiking neural architecture search framework for neuromorphic compute-in-memory systems," *arXiv e-prints*, pp. arXiv–2407, 2024.
- [23] W. Pan, F. Zhao, Z. Zhao, and Y. Zeng, "Brain-inspired evolutionary architectures for spiking neural networks," *IEEE Transactions on Artificial Intelligence*, 2024.
- [24] J. Yan, Q. Liu, M. Zhang, L. Feng, D. Ma, H. Li, and G. Pan, "Efficient spiking neural network design via neural architecture search," *Neural Networks*, p. 106172, 2024.
- [25] Q. Liu, J. Yan, M. Zhang, G. Pan, and H. Li, "Lite-snn: Designing lightweight and efficient spiking neural network through spatial-temporal compressive network search and joint optimization," *arXiv preprint arXiv:2401.14652*, 2024.
- [26] D. Padovano, A. Carpegnà, A. Savino, and S. Di Carlo, "Spikeexplorer: Hardware-oriented design space exploration for spiking neural networks on fpga," *Electronics*, vol. 13, no. 9, p. 1744, 2024.
- [27] B. Ajay, M. Rao, et al., "Despine: Nas generated deep evolutionary adaptive spiking network for low power edge computing applications," in *2024 25th International Symposium on Quality Electronic Design (ISQED)*, pp. 1–8, IEEE, 2024.
- [28] Z. Wang, Q. Zhao, J. Cui, X. Liu, and D. Xu, "Autost: Training-free neural architecture search for spiking transformers," in *ICASSP 2024-2024 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pp. 3455–3459, IEEE, 2024.
- [29] J. Teich, "Hardware/software codesign: The past, the present, and predicting the future," *Proceedings of the IEEE*, vol. 100, no. Special Centennial Issue, pp. 1411–1430, 2012.
- [30] W. Gerstner and W. M. Kistler, *Spiking neuron models: Single neurons, populations, plasticity*. Cambridge university press, 2002.
- [31] S. Ghosh-Dastidar and H. Adeli, "Spiking neural networks," *International journal of neural systems*, vol. 19 4, pp. 295–308, 2009.
- [32] C. Koch and I. Segev, *Methods in neuronal modeling: from ions to networks*. MIT press, 1998.
- [33] A. L. Hodgkin and A. F. Huxley, "A quantitative description of membrane current and its application to conduction and excitation in nerve," *The Journal of physiology*, vol. 117, no. 4, p. 500, 1952.
- [34] E. M. Izhikevich, "Simple model of spiking neurons," *IEEE Transactions on neural networks*, vol. 14, no. 6, pp. 1569–1572, 2003.
- [35] A. Taherkhani, A. Belatreche, Y. Li, G. Cosma, L. P. Maguire, and T. M. McGinnity, "A review of learning in biologically plausible spiking neural networks," *Neural Networks*, vol. 122, pp. 253–272, 2020.
- [36] A. Almomani, M. Alauthman, M. Alweshah, O. Dorgham, and F. Albalas, "A comparative study on spiking neural network encoding schema: implemented with cloud computing," *Cluster Computing*, vol. 22, pp. 419–433, 2019.
- [37] I. Aliyev, K. Svoboda, T. Adegbija, and J.-M. Fellous, "Sparsity-aware hardware-software co-design of spiking neural networks: An overview," in *2024 IEEE 17th International Symposium on Embedded Multicore/Many-core Systems-on-Chip (MCSoc)*, pp. 413–420, IEEE, 2024.
- [38] W. Guo, M. E. Fouad, A. M. Eltawil, and K. N. Salama, "Neural coding in spiking neural networks: A comparative study for robust neuromorphic systems," *Frontiers in Neuroscience*, vol. 15, p. 638474, 2021.
- [39] S. Park, S. Kim, B. Na, and S. Yoon, "T2fsnn: Deep spiking neural networks with time-to-first-spike coding," in *2020 57th ACM/IEEE Design Automation Conference (DAC)*, pp. 1–6, IEEE, 2020.
- [40] J. Kim, H. Kim, S. Huh, J. Lee, and K. Choi, "Deep neural networks with weighted spikes," *Neurocomputing*, vol. 311, pp. 373–386, 2018.
- [41] E. M. Izhikevich, N. S. Desai, E. C. Walcott, and F. C. Hoppensteadt, "Bursts as a unit of neural information: selective communication via resonance," *Trends in neurosciences*, vol. 26, no. 3, pp. 161–167, 2003.
- [42] S. Sheik, "Spike based information processing in spiking neural networks," in *Proceedings of the 4th International Conference on Applications in Nonlinear Dynamics (ICAND 2016) 5*, pp. 177–188, Springer, 2017.
- [43] B. Rueckauer and S.-C. Liu, "Temporal pattern coding in deep spiking neural networks," in *2021 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, IEEE, 2021.
- [44] L. F. Abbott and S. B. Nelson, "Synaptic plasticity: taming the beast," *Nature neuroscience*, vol. 3, no. 11, pp. 1178–1183, 2000.
- [45] A. Shrestha and A. Mahmood, "Review of deep learning algorithms and architectures," *IEEE access*, vol. 7, pp. 53040–53065, 2019.
- [46] T. Masquelier and S. J. Thorpe, "Unsupervised learning of visual features through spike timing dependent plasticity," *PLoS computational biology*, vol. 3, no. 2, p. e31, 2007.
- [47] M. Beyeler, N. D. Dutt, and J. L. Krichmar, "Categorization and decision-making in a neurobiologically plausible spiking network using a stdp-like learning rule," *Neural Networks*, vol. 48, pp. 109–124, 2013.
- [48] A. Tavanei, T. Masquelier, and A. S. Maida, "Acquisition of visual features through probabilistic spike-timing-dependent plasticity," in *2016 International Joint Conference on Neural Networks (IJCNN)*, pp. 307–314, IEEE, 2016.
- [49] R. Güting, "To spike, or when to spike?," *Current opinion in neurobiology*, vol. 25, pp. 134–139, 2014.
- [50] R.-M. Memmesheimer, R. Rubin, B. P. Ölveczky, and H. Sompolinsky, "Learning precisely timed spikes," *Neuron*, vol. 82, no. 4, pp. 925–938, 2014.
- [51] N. Anwani and B. Rajendran, "Normad-normalized approximate descent based supervised learning rule for spiking neurons," in *2015 international joint conference on neural networks (IJCNN)*, pp. 1–8, IEEE, 2015.
- [52] A. Gilra and W. Gerstner, "Predicting non-linear dynamics by stable local learning in a recurrent spiking neural network," *Elife*, vol. 6, p. e28295, 2017.
- [53] W. Nicola and C. Clopath, "Supervised learning in spiking neural networks with force training," *Nature communications*, vol. 8, no. 1, p. 2208, 2017.
- [54] E. O. Neftci, H. Mostafa, and F. Zenke, "Surrogate gradient learning in spiking neural networks: Bringing the power of gradient-based optimization to spiking neural networks," *IEEE Signal Processing Magazine*, vol. 36, no. 6, pp. 51–63, 2019.
- [55] D. Huh and T. J. Sejnowski, "Gradient descent for spiking neural networks," *Advances in neural information processing systems*, vol. 31, 2018.
- [56] J. H. Lee, T. Delbrück, and M. Pfeiffer, "Training deep spiking neural networks using backpropagation," *Frontiers in neuroscience*, vol. 10, p. 508, 2016.
- [57] S. M. Bohte, J. N. Kok, and J. A. La Poutré, "Spikeprop: backpropagation for networks of spiking neurons.,," in *ESANN*, vol. 48, pp. 419–424, Bruges, 2000.
- [58] Y. Wu, L. Deng, G. Li, J. Zhu, and L. Shi, "Spatio-temporal backpropagation for training high-performance spiking neural networks," *Frontiers in neuroscience*, vol. 12, p. 331, 2018.
- [59] N. Rathi and K. Roy, "Diet-snn: A low-latency spiking neural network with direct input encoding and leakage and threshold optimization," *IEEE Transactions on Neural Networks and Learning Systems*, 2021.
- [60] I. Aliyev and T. Adegbija, "Fine-tuning surrogate gradient learning for optimal hardware performance in spiking neural networks," in *2024 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pp. 1–2, IEEE, 2024.
- [61] Y. Wang, M. Zhang, Y. Chen, and H. Qu, "Signed neuron with memory: Towards simple, accurate and high-efficient ann-snn conversion.,," in *IJCAI*, pp. 2501–2508, 2022.
- [62] B. Rueckauer, I.-A. Lungu, Y. Hu, M. Pfeiffer, and S.-C. Liu, "Conversion of continuous-valued deep networks to efficient event-driven

- nets for image classification," *Frontiers in neuroscience*, vol. 11, p. 682, 2017.
- [63] S. Deng and S. Gu, "Optimal conversion of conventional artificial neural networks to spiking neural networks," *arXiv preprint arXiv:2103.00476*, 2021.
- [64] N.-D. Ho and I.-J. Chang, "Tcl: an ann-to-snn conversion with trainable clipping layers," in *2021 58th ACM/IEEE Design Automation Conference (DAC)*, pp. 793–798, 2021.
- [65] G.-q. Bi and M.-m. Poo, "Synaptic modifications in cultured hippocampal neurons: dependence on spike timing, synaptic strength, and postsynaptic cell type," *Journal of neuroscience*, vol. 18, no. 24, pp. 10464–10472, 1998.
- [66] A. Vigneron and J. Martinet, "A critical survey of stdp in spiking neural networks for pattern recognition," in *2020 international joint conference on neural networks (ijcnn)*, pp. 1–9, IEEE, 2020.
- [67] P. Panda, S. A. Aketi, and K. Roy, "Toward scalable, efficient, and accurate deep spiking neural networks with backward residual connections, stochastic softmax, and hybridization," *Frontiers in Neuroscience*, vol. 14, p. 653, 2020.
- [68] F. Liu, W. Zhao, Y. Chen, Z. Wang, T. Yang, and L. Jiang, "Sstdp: Supervised spike timing dependent plasticity for efficient spiking neural network training," *Frontiers in Neuroscience*, vol. 15, p. 756876, 2021.
- [69] S. Valadez-Godinez, H. Sossa, and R. Santiago-Montero, "The step size impact on the computational cost of spiking neuron simulation," in *2017 Computing Conference*, pp. 722–728, IEEE, 2017.
- [70] B. Cramer, Y. Stradmann, J. Schemmel, and F. Zenke, "The heidelberg spiking data sets for the systematic evaluation of spiking neural networks," *IEEE Transactions on Neural Networks and Learning Systems*, vol. 33, no. 7, pp. 2744–2757, 2020.
- [71] M. Bouvier, A. Valentian, T. Mesquida, F. Rummens, M. Reyboz, E. Vianello, and E. Beigne, "Spiking neural networks hardware implementations and challenges: A survey," *ACM Journal on Emerging Technologies in Computing Systems (JETC)*, vol. 15, no. 2, pp. 1–35, 2019.
- [72] K. O. Stanley and R. Miikkulainen, "Evolving neural networks through augmenting topologies," *Evolutionary computation*, vol. 10, no. 2, pp. 99–127, 2002.
- [73] T. Elsken, J. H. Metzen, and F. Hutter, "Neural architecture search: A survey," *The Journal of Machine Learning Research*, vol. 20, no. 1, pp. 1997–2017, 2019.
- [74] S. Deng, H. Lin, Y. Li, and S. Gu, "Surrogate module learning: Reduce the gradient error accumulation in training spiking neural networks," in *International Conference on Machine Learning*, pp. 7645–7657, PMLR, 2023.
- [75] A. Javanshir, T. T. Nguyen, M. P. Mahmud, and A. Z. Kouzani, "Advancements in algorithms and neuromorphic hardware for spiking neural networks," *Neural Computation*, vol. 34, no. 6, pp. 1289–1328, 2022.
- [76] H. Pham, M. Guan, B. Zoph, Q. Le, and J. Dean, "Efficient neural architecture search via parameters sharing," in *International conference on machine learning*, pp. 4095–4104, PMLR, 2018.
- [77] J. Zhang, J. Zhang, D. Huo, and H. Chen, "Anas: Asynchronous neuromorphic hardware architecture search based on a system-level simulator," in *2023 60th ACM/IEEE Design Automation Conference (DAC)*, pp. 1–6, IEEE, 2023.
- [78] B. Zoph and Q. V. Le, "Neural architecture search with reinforcement learning," *arXiv preprint arXiv:1611.01578*, 2016.
- [79] B. Baker, O. Gupta, N. Naik, and R. Raskar, "Designing neural network architectures using reinforcement learning," *arXiv preprint arXiv:1611.02167*, 2016.
- [80] X. Dong and Y. Yang, "Nas-bench-201: Extending the scope of reproducible neural architecture search," *arXiv preprint arXiv:2001.00326*, 2020.
- [81] W. Li, Z. Zhu, S. Shao, Y. Lu, and A. Song, "Spiking spatio-temporal neural architecture search for eeg-based emotion recognition," *IEEE Transactions on Instrumentation and Measurement*, 2024.
- [82] K. Che, Z. Zhou, L. Yuan, J. Zhang, Y. Tian, and L. Leng, "Spatial-temporal search for spiking neural networks," *arXiv preprint arXiv:2410.18580*, 2024.
- [83] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," 2015.
- [84] J. Zhang, X. Zhang, J. Huang, J. Zhang, and H. Chen, "Ancoef: Asynchronous neuromorphic algorithm/hardware co-exploration framework with a fully asynchronous simulator," *arXiv preprint arXiv:2411.06059*, 2024.
- [85] H. Liu, K. Simonyan, O. Vinyals, C. Fernando, and K. Kavukcuoglu, "Hierarchical representations for efficient architecture search," *arXiv preprint arXiv:1711.00436*, 2017.
- [86] A. Brock, T. Lim, J. M. Ritchie, and N. Weston, "Smash: one-shot model architecture search through hypernetworks," *arXiv preprint arXiv:1708.05344*, 2017.
- [87] G. Hinton, O. Vinyals, and J. Dean, "Distilling the knowledge in a neural network," *arXiv preprint arXiv:1503.02531*, 2015.
- [88] H. Qiu, M. Ning, Z. Song, W. Fang, Y. Chen, T. Sun, Z. Ma, L. Yuan, and Y. Tian, "Self-architectural knowledge distillation for spiking neural networks," *Neural Networks*, p. 106475, 2024.
- [89] J. Gou, B. Yu, S. J. Maybank, and D. Tao, "Knowledge distillation: A survey," *International Journal of Computer Vision*, vol. 129, no. 6, pp. 1789–1819, 2021.
- [90] M. S. Abdelfattah, L. Dudziak, T. Chau, R. Lee, H. Kim, and N. D. Lane, "Best of both worlds: Automl codesign of a cnn and its hardware accelerator," in *2020 57th ACM/IEEE Design Automation Conference (DAC)*, pp. 1–6, IEEE, 2020.
- [91] W. Chen, Y. Wang, S. Yang, C. Liu, and L. Zhang, "You only search once: A fast automation framework for single-stage dnn/accelerator co-design," in *2020 Design, Automation & Test in Europe Conference & Exhibition (DATE)*, pp. 1283–1286, IEEE, 2020.
- [92] L. Yang, Z. Yan, M. Li, H. Kwon, L. Lai, T. Krishna, V. Chandra, W. Jiang, and Y. Shi, "Co-exploration of neural architectures and heterogeneous asic accelerator designs targeting multiple tasks," in *2020 57th ACM/IEEE Design Automation Conference (DAC)*, pp. 1–6, IEEE, 2020.
- [93] H. Benmeziane, K. E. Maghraoui, H. Ouarnoughi, S. Niar, M. Wistuba, and N. Wang, "A comprehensive survey on hardware-aware neural architecture search," *arXiv preprint arXiv:2101.09336*, 2021.
- [94] X. Zhang, W. Jiang, Y. Shi, and J. Hu, "When neural architecture search meets hardware implementation: from hardware awareness to co-design," in *2019 IEEE Computer Society Annual Symposium on VLSI (ISVLSI)*, pp. 25–30, IEEE, 2019.
- [95] N. Rathj, I. Chakraborty, A. Kosta, A. Sengupta, A. Ankit, P. Panda, and K. Roy, "Exploring neuromorphic computing based on spiking neural networks: Algorithms to hardware," *ACM Computing Surveys*, vol. 55, no. 12, pp. 1–49, 2023.
- [96] P. Qu, Y. Zhang, X. Fei, and W. Zheng, "High performance simulation of spiking neural network on gpgpus," *IEEE Transactions on Parallel and Distributed Systems*, vol. 31, no. 11, pp. 2510–2523, 2020.
- [97] E. Yavuz, J. Turner, and T. Nowotny, "Genn: a code generation framework for accelerated brain simulations," *Scientific reports*, vol. 6, no. 1, p. 18854, 2016.
- [98] M. Beyeler, K. D. Carlson, T.-S. Chou, N. Dutt, and J. L. Krichmar, "Carlsim 3: A user-friendly and highly optimized library for the creation of neurobiologically detailed spiking neural networks," in *2015 International Joint Conference on Neural Networks (IJCNN)*, pp. 1–8, IEEE, 2015.
- [99] J. Han, Z. Li, W. Zheng, and Y. Zhang, "Hardware implementation of spiking neural networks on fpga," *Tsinghua Science and Technology*, vol. 25, no. 4, pp. 479–486, 2020.
- [100] D. Pani, P. Meloni, G. Tuveri, F. Palumbo, P. Massobrio, and L. Raffo, "An fpga platform for real-time simulation of spiking neuronal networks," *Frontiers in neuroscience*, vol. 11, p. 90, 2017.
- [101] Y. Liu, Y. Chen, W. Ye, and Y. Gui, "Fpga-nhp: A general fpga-based neuromorphic hardware acceleration platform with high speed and low power," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 69, no. 6, pp. 2553–2566, 2022.
- [102] X. Ju, B. Fang, R. Yan, X. Xu, and H. Tang, "An fpga implementation of deep spiking neural networks for low-power and fast classification," *Neural computation*, vol. 32, no. 1, pp. 182–204, 2020.
- [103] D. Ma, J. Shen, Z. Gu, M. Zhang, X. Zhu, X. Xu, Q. Xu, Y. Shen, and G. Pan, "Darwin: A neuromorphic hardware co-processor based on spiking neural networks," *Journal of systems architecture*, vol. 77, pp. 43–51, 2017.
- [104] S. Li, Z. Zhang, R. Mao, J. Xiao, L. Chang, and J. Zhou, "A fast and energy-efficient snn processor with adaptive clock/event-driven computation scheme and online learning," *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 68, no. 4, pp. 1543–1552, 2021.
- [105] M. Davies, N. Srinivasa, T.-H. Lin, G. Chinya, Y. Cao, S. H. Choday, G. Dimou, P. Joshi, N. Imam, S. Jain, et al., "Loihi: A neuromorphic manycore processor with on-chip learning," *Ieee Micro*, vol. 38, no. 1, pp. 82–99, 2018.
- [106] F. Akopyan, J. Sawada, A. Cassidy, R. Alvarez-Icaza, J. Arthur, P. Merolla, N. Imam, Y. Nakamura, P. Datta, G.-J. Nam, et al., "Truenorth: Design and tool flow of a 65 mw 1 million neuron

- programmable neurosynaptic chip.” *IEEE transactions on computer-aided design of integrated circuits and systems*, vol. 34, no. 10, pp. 1537–1557, 2015.
- [107] H.-Y. Hsieh and K.-T. Tang, “Vlsi implementation of a bio-inspired olfactory spiking neural network,” *IEEE transactions on neural networks and learning systems*, vol. 23, no. 7, pp. 1065–1073, 2012.
- [108] J. Schemmel, D. Brüderle, A. Grübl, M. Hock, K. Meier, and S. Millner, “A wafer-scale neuromorphic hardware system for large-scale neural modeling,” in *2010 IEEE International Symposium on Circuits and Systems (ISCAS)*, pp. 1947–1950, IEEE, 2010.
- [109] C. Gao and D. Hammerstrom, “Cortical models onto cmos and cmos—architectures and performance/price,” *IEEE Transactions on Circuits and Systems I: Regular Papers*, vol. 54, no. 11, pp. 2502–2515, 2007.
- [110] E. Painkras, L. A. Plana, J. Garside, S. Temple, F. Galluppi, C. Patterson, D. R. Lester, A. D. Brown, and S. B. Furber, “Spinnaker: A 1-w 18-core system-on-chip for massively-parallel neural network simulation,” *IEEE Journal of Solid-State Circuits*, vol. 48, no. 8, pp. 1943–1953, 2013.
- [111] W. Jiang, L. Yang, S. Dasgupta, J. Hu, and Y. Shi, “Standing on the shoulders of giants: Hardware and neural architecture co-search with hot start,” *IEEE Transactions on Computer-Aided Design of Integrated Circuits and Systems*, vol. 39, no. 11, pp. 4154–4165, 2020.
- [112] E. Strubell, A. Ganesh, and A. McCallum, “Energy and policy considerations for deep learning in NLP,” *CoRR*, vol. abs/1906.02243, 2019.
- [113] S. Xie, H. Zheng, C. Liu, and L. Lin, “Snas: stochastic neural architecture search,” *arXiv preprint arXiv:1812.09926*, 2018.
- [114] S. Shen, R. Zhang, C. Wang, R. Huang, A. Tuerhong, Q. Guo, Z. Lu, J. Zhang, and L. Leng, “Evolutionary spiking neural networks: a survey,” *Journal of Membrane Computing*, pp. 1–12, 2024.
- [115] P. Liashchynskyi and P. Liashchynskyi, “Grid search, random search, genetic algorithm: a big comparison for nas,” *arXiv preprint arXiv:1912.06059*, 2019.
- [116] J. Bergstra and Y. Bengio, “Random search for hyper-parameter optimization.,” *Journal of machine learning research*, vol. 13, no. 2, 2012.
- [117] K. Yu, C. Sciuto, M. Jaggi, C. Musat, and M. Salzmann, “Evaluating the search phase of neural architecture search,” *arXiv preprint arXiv:1902.08142*, 2019.
- [118] L. Li and A. Talwalkar, “Random search and reproducibility for neural architecture search,” in *Uncertainty in artificial intelligence*, pp. 367–377, PMLR, 2020.
- [119] I. Radenovic, R. P. Kosaraju, R. Girshick, K. He, and P. Dollár, “Designing network design spaces,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10428–10436, 2020.
- [120] G. F. Miller, P. M. Todd, and S. U. Hegde, “Designing neural networks using genetic algorithms.,” in *ICGA*, vol. 89, pp. 379–384, 1989.
- [121] A. Guha, S. A. Harp, and T. Samad, “Genetic synthesis of neural networks,” *Honeywell Corporate Syst. Development Division, Tech. Rep. CSDD-88-I4852-CC-1*, 1988.
- [122] P. Todd, “Evolutionary methods for connectionist architectures,” *Psychology Dept. Stanford University, unpublished Manuscript*, 1988.
- [123] Y. Liu, Y. Sun, B. Xue, M. Zhang, G. G. Yen, and K. C. Tan, “A survey on evolutionary neural architecture search,” *IEEE transactions on neural networks and learning systems*, 2021.
- [124] A. Shrestha and A. Mahmood, “Optimizing deep neural network architecture with enhanced genetic algorithm,” in *2019 18th IEEE International Conference On Machine Learning And Applications (ICMLA)*, pp. 1365–1370, IEEE, 2019.
- [125] S. R. Young, D. C. Rose, T. P. Karnowski, S.-H. Lim, and R. M. Patton, “Optimizing deep learning hyper-parameters through an evolutionary algorithm,” in *Proceedings of the workshop on machine learning in high-performance computing environments*, pp. 1–5, 2015.
- [126] E. Real, S. Moore, A. Selle, S. Saxena, Y. L. Suematsu, J. Tan, Q. V. Le, and A. Kurakin, “Large-scale evolution of image classifiers,” in *International conference on machine learning*, pp. 2902–2911, PMLR, 2017.
- [127] J.-B. Mouret and J. Clune, “Illuminating search spaces by mapping elites,” *arXiv preprint arXiv:1504.04909*, 2015.
- [128] E. Real, A. Aggarwal, Y. Huang, and Q. V. Le, “Regularized evolution for image classifier architecture search,” in *Proceedings of the aaai conference on artificial intelligence*, vol. 33, pp. 4780–4789, 2019.
- [129] J. Prellberg and O. Kramer, “Lamarckian evolution of convolutional neural networks,” in *Parallel Problem Solving from Nature—PPSN XV: 15th International Conference, Coimbra, Portugal, September 8–12, 2018, Proceedings, Part II 15*, pp. 424–435, Springer, 2018.
- [130] T. Elsken, J. H. Metzen, and F. Hutter, “Efficient multi-objective neural architecture search via lamarckian evolution,” *arXiv preprint arXiv:1804.09081*, 2018.
- [131] J. Lin, W.-M. Chen, Y. Lin, C. Gan, S. Han, et al., “Mcunet: Tiny deep learning on iot devices,” *Advances in Neural Information Processing Systems*, vol. 33, pp. 11711–11722, 2020.
- [132] H. Cai, C. Gan, T. Wang, Z. Zhang, and S. Han, “Once-for-all: Train one network and specialize it for efficient deployment,” *arXiv preprint arXiv:1908.09791*, 2019.
- [133] H. Wang, Z. Wu, Z. Liu, H. Cai, L. Zhu, C. Gan, and S. Han, “Hat: Hardware-aware transformers for efficient natural language processing,” *arXiv preprint arXiv:2005.14187*, 2020.
- [134] A. Marchisio, A. Massa, V. Mrazek, B. Bussolino, M. Martina, and M. Shafique, “Nascaps: A framework for neural architecture search to optimize the accuracy and hardware efficiency of convolutional capsule networks,” in *Proceedings of the 39th International Conference on Computer-Aided Design*, pp. 1–9, 2020.
- [135] L. L. Zhang, Y. Yang, Y. Jiang, W. Zhu, and Y. Liu, “Fast hardware-aware neural architecture search,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pp. 692–693, 2020.
- [136] P. Ren, Y. Xiao, X. Chang, P.-Y. Huang, Z. Li, X. Chen, and X. Wang, “A comprehensive survey of neural architecture search: Challenges and solutions,” *ACM Computing Surveys (CSUR)*, vol. 54, no. 4, pp. 1–34, 2021.
- [137] R. S. Sutton and A. G. Barto, *Reinforcement learning: An introduction*. MIT press, 2018.
- [138] B. Baker, O. Gupta, R. Raskar, and N. Naik, “Accelerating neural architecture search using performance prediction,” *arXiv preprint arXiv:1705.10823*, 2017.
- [139] M. Tan, B. Chen, R. Pang, V. Vasudevan, M. Sandler, A. Howard, and Q. V. Le, “Mnasnet: Platform-aware neural architecture search for mobile,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 2820–2828, 2019.
- [140] Q. Lu, W. Jiang, X. Xu, Y. Shi, and J. Hu, “On neural architecture search for resource-constrained hardware platforms,” *arXiv preprint arXiv:1911.00105*, 2019.
- [141] H. Cai, T. Chen, W. Zhang, Y. Yu, and J. Wang, “Efficient architecture search by network transformation,” in *Proceedings of the AAAI conference on artificial intelligence*, vol. 32, 2018.
- [142] C. Ying, A. Klein, E. Christiansen, E. Real, K. Murphy, and F. Hutter, “Nas-bench-101: Towards reproducible neural architecture search,” in *International conference on machine learning*, pp. 7105–7114, PMLR, 2019.
- [143] B. Shahriari, K. Swersky, Z. Wang, R. P. Adams, and N. De Freitas, “Taking the human out of the loop: A review of bayesian optimization,” *Proceedings of the IEEE*, vol. 104, no. 1, pp. 148–175, 2015.
- [144] H. Liu, K. Simonyan, and Y. Yang, “Darts: Differentiable architecture search,” *arXiv preprint arXiv:1806.09055*, 2018.
- [145] H. Cai, L. Zhu, and S. Han, “Proxlessnas: Direct neural architecture search on target task and hardware,” *CoRR*, vol. abs/1812.00332, 2018.
- [146] B. Wu, X. Dai, P. Zhang, Y. Wang, F. Sun, Y. Wu, Y. Tian, P. Vajda, Y. Jia, and K. Keutzer, “Fbnets: Hardware-aware efficient convnet design via differentiable neural architecture search,” in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 10734–10742, 2019.
- [147] G. Bingham, “Ineffectiveness of gradient-based neural architecture search,”
- [148] Y. Jiang, X. Wang, and W. Zhu, “Hardware-aware transformable architecture search with efficient search space,” in *2020 IEEE International Conference on Multimedia and Expo (ICME)*, pp. 1–6, IEEE, 2020.
- [149] C. Liu, B. Zoph, M. Neumann, J. Shlens, W. Hua, L.-J. Li, L. Fei-Fei, A. Yuille, J. Huang, and K. Murphy, “Progressive neural architecture search,” in *Proceedings of the European conference on computer vision (ECCV)*, pp. 19–34, 2018.
- [150] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le, “Learning transferable architectures for scalable image recognition,” in *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 8697–8710, 2018.
- [151] Y. Guo, Y. Chen, Y. Zheng, P. Zhao, J. Chen, J. Huang, and M. Tan, “Breaking the curse of space explosion: Towards efficient nas with curriculum search,” in *International Conference on Machine Learning*, pp. 3822–3831, PMLR, 2020.
- [152] Y. Chen, G. Meng, Q. Zhang, S. Xiang, C. Huang, L. Mu, and X. Wang, “Reinforced evolutionary neural architecture search,” *arXiv preprint arXiv:1808.00193*, 2018.

- [153] Y. Chen, G. Meng, Q. Zhang, S. Xiang, C. Huang, L. Mu, and X. Wang, "Renas: Reinforced evolutionary neural architecture search," in *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pp. 4787–4796, 2019.
- [154] K. Maziarz, M. Tan, A. Khorlin, M. Georgiev, and A. Gesmundo, "Evolutionary-neural hybrid agents for architecture search," *arXiv preprint arXiv:1811.09828*, 2018.
- [155] R. J. Williams, "Simple statistical gradient-following algorithms for connectionist reinforcement learning," *Machine learning*, vol. 8, pp. 229–256, 1992.
- [156] H. Bouzidi, H. Ouarnoughi, S. Niar, and A. A. E. Cadi, "Performance prediction for convolutional neural networks in edge devices," *arXiv preprint arXiv:2010.11297*, 2020.
- [157] C.-H. Hsu, S.-H. Chang, J.-H. Liang, H.-P. Chou, C.-H. Liu, S.-C. Chang, J.-Y. Pan, Y.-T. Chen, W. Wei, and D.-C. Juan, "Monas: Multi-objective neural architecture search using reinforcement learning," *arXiv preprint arXiv:1806.10332*, 2018.
- [158] T.-J. Yang, A. Howard, B. Chen, X. Zhang, A. Go, M. Sandler, V. Sze, and H. Adam, "Netadapt: Platform-aware neural network adaptation for mobile applications," in *Proceedings of the European Conference on Computer Vision (ECCV)*, pp. 285–300, 2018.
- [159] T. Wang, K. Wang, H. Cai, J. Lin, Z. Liu, H. Wang, Y. Lin, and S. Han, "Apq: Joint search for network architecture, pruning and quantization policy," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 2078–2087, 2020.
- [160] A. Sood, B. Elder, B. Herta, C. Xue, C. Bekas, A. C. I. Malossi, D. Saha, F. Scheidegger, G. Venkataraman, G. Thomas, *et al.*, "Neunets: An automated synthesis engine for neural network design," *arXiv preprint arXiv:1901.06261*, 2019.
- [161] H. Cai, L. Zhu, and S. Han, "Proxylessnas: Direct neural architecture search on target task and hardware," *arXiv preprint arXiv:1812.00332*, 2018.
- [162] S. Falkner, A. Klein, and F. Hutter, "Bohb: Robust and efficient hyperparameter optimization at scale," in *International conference on machine learning*, pp. 1437–1446, PMLR, 2018.
- [163] L. Li, K. G. Jamieson, G. DeSalvo, A. Rostamizadeh, and A. Talwalkar, "Hyperband: Bandit-based configuration evaluation for hyperparameter optimization," in *ICLR (Poster)*, p. 53, 2017.
- [164] F. Runge, D. Stoll, S. Falkner, and F. Hutter, "Learning to design rna," *arXiv preprint arXiv:1812.11951*, 2018.
- [165] A. Zela, A. Klein, S. Falkner, and F. Hutter, "Towards automated deep learning: Efficient joint neural architecture and hyperparameter search," *arXiv preprint arXiv:1807.06906*, 2018.
- [166] D. Zhou, X. Zhou, W. Zhang, C. C. Loy, S. Yi, X. Zhang, and W. Ouyang, "Econas: Finding proxies for economical neural architecture search," in *Proceedings of the IEEE/CVF Conference on computer vision and pattern recognition*, pp. 11396–11404, 2020.
- [167] J. Mellor, J. Turner, A. Storkey, and E. J. Crowley, "Neural architecture search without training," in *International Conference on Machine Learning*, pp. 7588–7598, PMLR, 2021.
- [168] W. Chen, X. Gong, and Z. Wang, "Neural architecture search on imagenet in four gpu hours: A theoretically inspired perspective," *arXiv preprint arXiv:2102.11535*, 2021.
- [169] J. Xu, L. Zhao, J. Lin, R. Gao, X. Sun, and H. Yang, "Knas: green neural architecture search," in *International Conference on Machine Learning*, pp. 11613–11625, PMLR, 2021.
- [170] Z. Sun, M. Lin, X. Sun, Z. Tan, H. Li, and R. Jin, "Mae-det: Revisiting maximum entropy principle in zero-shot nas for efficient object detection," *arXiv preprint arXiv:2111.13336*, 2021.
- [171] V. Lopes, S. Alirezaee, and L. A. Alexandre, "Epe-nas: Efficient performance estimation without training for neural architecture search," in *International conference on artificial neural networks*, pp. 552–563, Springer, 2021.
- [172] Y. Jaafra, J. L. Laurent, A. Deruyver, and M. S. Naceur, "Reinforcement learning for neural architecture search: A review," *Image and Vision Computing*, vol. 89, pp. 57–66, 2019.
- [173] C. White, M. Safari, R. Sukthanker, B. Ru, T. Elsken, A. Zela, D. Dey, and F. Hutter, "Neural architecture search: Insights from 1000 papers," *arXiv preprint arXiv:2301.08727*, 2023.
- [174] X. He, K. Zhao, and X. Chu, "Automl: A survey of the state-of-the-art," *Knowledge-based systems*, vol. 212, p. 106622, 2021.
- [175] C. Li, J. Peng, L. Yuan, G. Wang, X. Liang, L. Lin, and X. Chang, "Block-wisely supervised neural architecture search with knowledge distillation," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 1989–1998, 2020.
- [176] M. Zhang, H. Li, S. Pan, X. Chang, and S. Su, "Overcoming multi-model forgetting in one-shot nas with diversity maximization," in *Proceedings of the ieee/cvf conference on computer vision and pattern recognition*, pp. 7809–7818, 2020.
- [177] C. Liu, P. Dollár, K. He, R. Girshick, A. Yuille, and S. Xie, "Are labels necessary for neural architecture search?," in *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pp. 798–813, Springer, 2020.
- [178] W. Jiang, Q. Lou, Z. Yan, L. Yang, J. Hu, X. S. Hu, and Y. Shi, "Device-circuit-architecture co-exploration for computing-in-memory neural accelerators," *IEEE Transactions on Computers*, vol. 70, no. 4, pp. 595–605, 2020.
- [179] I. Aliyev, J. Lopez, and T. Adegbija, "Exploring the sparsity-quantization interplay on a novel hybrid snn event-driven architecture," 2025.