

Composing Recurrent Spiking Neural Networks using Locally-Recurrent Motifs and Risk-Mitigating Architectural Optimization

Wenrui Zhang 1 , Hejia Geng 1 and Peng Li 1,*

¹Department of Electrical and Computer Engineering, University of California, Santa Barbara, CA 93106

Correspondence*: Peng Li lip@ucsb.edu

ABSTRACT

3 In neural circuits, recurrent connectivity plays a crucial role in network function and stability. However, existing recurrent spiking neural networks (RSNNs) are often constructed by random 4 connections without optimization. While RSNNs can produce rich dynamics that are critical for memory formation and learning, systemic architectural optimization of RSNNs is still an open challenge. We aim to enable systematic design of large RSNNs via a new scalable RSNN architecture and automated architectural optimization. We compose RSNNs based on a layer architecture called Sparsely-Connected Recurrent Motif Layer (SC-ML) that consists of 9 multiple small recurrent motifs wired together by sparse lateral connections. The small size of 10 the motifs and sparse inter-motif connectivity leads to an RSNN architecture scalable to large 11 12 network sizes. We further propose a method called Hybrid Risk-Mitigating Architectural Search (HRMAS) to systematically optimize the topology of the proposed recurrent motifs and SC-ML layer architecture. HRMAS is an alternating two-step optimization process by which we mitigate the risk of network instability and performance degradation caused by architectural change by 15 introducing a novel biologically-inspired "self-repairing" mechanism through intrinsic plasticity. The intrinsic plasticity is introduced to the second step of each HRMAS iteration and acts as 17 unsupervised fast self-adaptation to structural and synaptic weight modifications introduced by 18 the first step during the RSNN architectural "evolution". To the best of the authors' knowledge, this 19 is the first work that performs systematic architectural optimization of RSNNs. Using one speech 20 and three neuromorphic datasets, we demonstrate the significant performance improvement 21 brought by the proposed automated architecture optimization over existing manually-designed 22 RSNNs.

- 24 Keywords: Brain Inspired Computing, Recurrent Spiking Neural Networks, Neural Architecture Search, Sparsely-Connected Recurrent
- 25 Motif Layer, Intrinsic Plasticity

1 INTRODUCTION

- 26 In the brain, recurrent connectivity is indispensable for maintaining dynamics, functions, and oscillations of
- 27 the network Buzsaki (2006). As a brain-inspired computational model, spiking neural networks (SNNs) are
- 28 well suited for processing spatiotemporal information (Maass, 1997). In particular, recurrent spiking neural

58 59

60

61

66

67

69

71

networks (RSNNs) can mimic microcircuits in the biological brain and induce rich behaviors that are critical 30 for memory formation and learning. Recurrence has been explored in conventional non-spiking artificial neural networks (ANNs) in terms of Long Short Term Memory (LSTM) (Hochreiter and Schmidhuber, 31 1997), Echo State Networks (ESN) (Jaeger, 2001), Deep RNNs (Graves et al., 2013), Gated Recurrent Units 32 (GRU) (Cho et al., 2014), and Legendre Memory Units (LMU) (Voelker et al., 2019). While recurrence 33 presents unique challenges and opportunities in the context of spiking neural networks, RSNNs are yet to 34 be well explored. 35

Most existing works on RSNNs adopt recurrent layers or reservoirs with randomly generated connections. 36 37 The Liquid State Machine (LSM) (Maass et al., 2002) is one of the most widely adopted RSNN architectures with one or multiple recurrent reservoirs and an output readout layer wired up using feedforward 38 39 synapses (Zhang et al., 2015; Wang and Li, 2016; Srinivasan et al., 2018). However, there is a lack of principled approaches for setting up the recurrent connections in reservoirs. Instead, ad-hoc randomly 40 generated wiring patterns are often adopted. Bellec et al. (2018) proposed an architecture called long short-41 term memory SNNs (LSNNs). The recurrent layer contains a regular spiking portion with both inhibitory 42 and excitatory spiking neurons and an adaptive neural population. Zhang and Li (2019b) proposed to 43 train deep RSNNs by a spike-train level backpropagation (BP) method. Maes et al. (2020) demonstrated 44 45 a new reservoir with multiple groups of excitatory neurons and a central group of inhibitory neurons. Furthermore, Zhang and Li (2020a) presented a recurrent structure named ScSr-SNNs in which recurrence 46 is simply formed by a self-recurrent connection to each neuron. However, the recurrent connections in 47 48 all of these works are either randomly generated with certain probabilities or simply constructed by selfrecurrent connections. Randomly generated or simple recurrent connections may not effectively optimize 49 RSNNs' performance. Recently, Pan et al. (2023) introduced a multi-objective Evolutionary Liquid State 50 Machine (ELSM) inspired by neuroevolution process. Chakraborty and Mukhopadhyay (2023) proposed 51 52 Heterogeneous recurrent spiking neural network (HRSNN), in which recurrent layers are composed of 53 heterogeneous neurons with different dynamics. Chen et al. (2023) introduced an intralayer-connected SNN and a hybrid training method combining probabilistic spike-timing dependent plasticity (STDP) with 55 BP. But their performance still has significant gaps. Systemic RSNN architecture design and optimization 56 remain as an open problem.

Neural architectural search (NAS), the process of automating the construction of non-spiking ANNs, has become prevalent recently after achieving state-of-the-art performance on various tasks (Elsken et al., 2019; Wistuba et al., 2019). Different types of strategies such as reinforcement learning (Zoph and Le, 2017), gradient-based optimization (Liu et al., 2018), and evolutionary algorithms (Real et al., 2019) have been proposed to find optimal architectures of traditional CNNs and RNNs. In contrast, the architectural optimization of SNNs has received little attention. Only recently, Tian et al. (2021) adopted a simulated 62 annealing algorithm to learn the optimal architecture hyperparameters of liquid state machine (LSM) 63 models through a three-step search. Similarly, a surrogate-assisted evolutionary search method was applied 64 in Zhou et al. (2020) to optimize the hyperparameters of LSM such as density, probability and distribution 65 of connections. However, both studies focused only on LSM for which hyperparameters indirectly affecting recurrent connections as opposed to specific connectivity patterns were optimized. Even after selecting the hyperparameters, the recurrence in the network remained randomly determined without any optimization. 68 Recently, Kim et al. (2022) explored a cell-based neural architecture search method on SNNs, but did not involve large-scale recurrent connections. Na et al. (2022) introduced a spike-aware NAS framework called AutoSNN to investigate the impact of architectural components on SNNs' performance and energy efficiency. Overall, NAS for RSNNs is still rarely explored.

91

92 93

94

95

96

97

98

99

100

101

104

73 This paper aims to enable systematic design of large recurrent spiking neural networks (RSNNs) via a 74 new scalable RSNN architecture and automated architectural optimization. RSNNs can create complex network dynamics both in time and space, which manifests itself as an opportunity for achieving great 75 76 learning capabilities and a challenge in practical realization. It is important to strike a balance between 77 theoretical computational power and architectural complexity. Firstly, we argue that composing RSNNs based on well-optimized building blocks small in size, or recurrent motifs, can lead to an architectural 78 solution scalable to large networks while achieving high performance. We assemble multiple recurrent 79 80 motifs into a layer architecture called Sparsely-Connected Recurrent Motif Layer (SC-ML). The motifs in each SC-ML share the same topology, defined by the size of the motif, i.e., the number of neurons, and 81 the recurrent connectivity pattern between the neurons. The motif topology is determined by the proposed architectural optimization while the weights within each motif may be tuned by standard backpropagation 83 training algorithms. Motifs in a recurrent SC-ML layer are wired together using sparse lateral connections 84 determined by imposing spatial connectivity constraints. As such, there exist two levels of structured recurrence: recurrence within each motif and recurrence between the motifs at the SC-ML level. The 86 fact that the motifs are small in size and that inter-motif connectivity is sparse alleviates the difficulty in 87 architectural optimization and training of these motifs and SC-ML. Furthermore, multiple SC-ML layers can be stacked and wired using additional feedforward weights to construct even larger recurrent networks. 89

Secondly, we demonstrate a method called Hybrid Risk-Mitigating Architectural Search (HRMAS) to optimize the proposed recurrent motifs and SC-ML layer architecture. HRMAS is an alternating two-step optimization process hybridizing bio-inspired intrinsic plasticity for mitigating the risk in architectural optimization. Facilitated by gradient-based methods (Liu et al., 2018; Zhang and Li, 2020b), the first step of optimization is formulated to optimize network architecture defined by the size of the motif, intra and inter-motif connectivity patterns, types of these connections, and the corresponding synaptic weight values, respectively.

While structural changes induced by the architectural-level optimization are essential for finding highperformance RSNNs, they may be misguided due to discontinuity in architectural search, and limited training data, hence leading to over-fitting. We mitigate the risk of network instability and performance degradation caused by architectural change by introducing a novel biologically-inspired "self-repairing" mechanism through intrinsic plasticity, which has the same spirit of homeostasis during neural development 102 (Tien and Kerschensteiner, 2018). The intrinsic plasticity is introduced in the second step of each HRMAS iteration and acts as unsupervised self-adaptation to mitigate the risks imposed by structural and synaptic 103 weight modifications introduced by the first step during the RSNN architectural "evolution".

We evaluate the proposed techniques on speech dataset TI46-Alpha (Liberman et al., 1991), neuromorphic 105 speech dataset N-TIDIGITS (Anumula et al., 2018), neuromorphic video dataset DVS-Gesture (Amir et al., 106 2017), and neuromorphic image dataset N-MNIST (Orchard et al., 2015). The SC-ML-based RSNNs 107 optimized by HRMAS achieve state-of-the-art performance on all four datasets. With the same network 108 size, automated network design via HRMAS outperforms existing RSNNs by up to 3.38% performance 109 improvement. 110

MATERIALS AND METHODS

2.1 **Spiking Neuron Model** 111

112 In this work, we adopt the leaky integrate-and-fire (LIF) neuron model Gerstner and Kistler (2002) which is one of the most popular neuron models for simulating SNNs. During the simulation, we use the fixed-step 113

- first-order Euler method to discretize the LIF model. In the rest of this paper, we only analyze an SNN in
- 115 the discretized form. Consider the input spike train from pre-synaptic neuron j: $s_j[t] = \sum_{t_j^{(f)}} \delta[t t_j^{(f)}],$
- where $t_i^{(f)}$ denotes a particular firing time of presynaptic neuron j. The incoming spikes are converted
- into an (unweighted) postsynaptic current (PSC) $a_i[t]$ through a synaptic model. We adopt the first-order
- 118 synaptic model Gerstner and Kistler (2002):

$$a_j[t] = (1 - \frac{1}{\tau_{syn}})a_j[t - 1] + s_j[t], \tag{1}$$

119 where τ_{syn} is the synaptic time constant. Then, the neuronal membrane voltage $u_i[t]$ of neuron i at time t is

120 given by

$$u_i^{-}[t] = (1 - \frac{1}{\tau})u_i[t - 1] + \frac{R}{\tau} \sum_j w_{ij} a_j[t], \tag{2}$$

121

125

$$u_i[t] = \begin{cases} 0, & \text{if } u_i^-[t] > V_{\text{th}} \\ u_i^-[t], & \text{otherwise} \end{cases}$$
 (3)

where R and τ are the resistance and time constant of the membrane, w_{ij} the synaptic weight from pre-synaptic neuron j to neuron i. Moreover, the firing output of the neuron is expressed as

$$s_i[t] = H\left(u_i[t] - V_{th}\right),\tag{4}$$

where V_{th} is the firing threshold and $H(\cdot)$ is the Heaviside step function.

2.2 Sparsely-Connected Recurrent Motif Layer (SC-ML)

Unlike the traditional non-spiking RNNs that are typically constructed with units like LSTM or GRU, the 126 127 structure of existing RSNNs is random without specific optimization, which hinders RSNNs' performance and prevents scaling to large networks. However, due to the complexity of recurrent connections and 128 dynamics of spiking neurons, the optimization of RSNNs weights is still an open problem. As shown in 129 Table 3, recurrent connections that are not carefully set up may hinder network performance. To solve 130 this problem, we first designed the SC-ML layer, which is composed of multiple sparsely-connected 131 recurrent motifs, where each motif consists of a group of recurrently connected spiking neurons, as shown 132 in Figure 1. The motifs in each SC-ML share the same topology, which is defined as the size of the 133 motif, i.e., the number of neurons, and the recurrent connectivity pattern between the neurons (excitatory, 134 135 inhibitory or non-existent). Within the motif, synaptic connections can be constructed between any two neurons including self-recurrent connections. Thus the problem of the recurrent layer optimization can be 136 simplified to that of learning the optimal motif and sparse inter-motif connectivity, alleviating the difficulty 137 in architectural optimization and allowing scalability to large networks. 138

This motif-based structure is motivated by both a biological and a computational perspective. First, from a biological point of view, there is evidence that the neocortex is not only organized in layered minicolumn structures but also into synaptically connected clusters of neurons within such structures (Perin et al., 2011; Ko et al., 2011). For example, the networks of pyramidal cells cluster into multiple groups of a few dozen neurons each. Second, we add onto the memory effects resulting from temporal integration of individual spiking neurons by introducing sparse intra or inter-motif connections. This corresponds to a scalable and biologically plausible RSNN architectural design space that closely mimics the microcircuits in the nervous

158

172

173174

175

176

177178

179

180

181

182 183

184 185

186

187

system. From a computational perspective, optimizing the connectivity of the basic building block, i.e., the motif, simplifies the problem of optimizing the connectivity of the whole recurrent layer. Furthermore, by constraining most recurrent connections inside the motifs and allowing a few lateral connections between neighboring motifs to exchange information across the SC-ML, the total number of recurrent connections is limited. This leads to a great deal of sparsity as observed in biological networks (Seeman et al., 2018).

Figure 1 presents an example of SC-ML with 12-neuron motifs. The lateral inter-motif connections can be introduced as the mutual connections between two corresponding neurons in neighboring motifs to ensure sparsity and reduce complexity. With the proposed SC-ML, one can easily stack multiple SC-MLs to form a multi-layer large RSNN using feedforward weights. Within a multi-layered network, information processing is facilitated through local processing of different motifs, communication of motif-level responses via inter-motif connections, and extraction and processing of higher-level features layer by layer.

2.3 Hybrid Risk-Mitigating Architectural Search (HRMAS)

spiking RNNs, where a substructure called cell is optimized by a search algorithm (Zoph and Le, 2017). Nevertheless, this NAS approach may not be the best fit for RSNNs. First, recurrence in the cell is only created by feeding previous hidden state back to the cell while connectivity inside the cell is feedforward. Second, the overall operations and connectivity found by the above NAS procedure do not go beyond an LSTM-like architecture. Finally, the considered combination operations and activation functions like addition and elementwise multiplication are not biologically plausible.

Neural architecture search (NAS) has been applied for architectural optimization of traditional non-

In order to extend NAS to a wider range of spiking RNNs, we introduce the Hybrid Risk-Mitigating Architectural Search (HRMAS). This framework systematically optimizes the motif topology and lateral connections of SC-ML. Each optimization iteration consists of two alternating steps, hybridizing gradient-based optimization and biologically-inspired intrinsic plasticity for robust NAS of RSNNs. We will introduce the overall idea of HRMAS in 2.3.1, the optimization problem of HRMAS in 2.3.2, the gradient-based optimization part in 2.3.3, and the bio-inspired optimization part in 2.3.4.

171 2.3.1 Hybrid Risk-Mitigating Architectural Search Framework

In HRMAS, all recurrent connections are categorized into three types: inhibitory, excitatory, and nonexistence. An inhibitory connection has a negative weight and is fixed without training in our current implementation. In the recurrent network, negative weights mainly provide the function of inhibitory stimulation. Here we follow the settings in previous research (Zhang and Li, 2020a, 2021) and adopt fixed negative weights. In experiments, fixed negative weights can reduce the optimization complexity without significant performance loss, while providing stable inhibitory connections. The weight of an excitatory connection is positive and trained by a backpropagation (BP) method. HRMAS is an alternating two-step optimization process, hybridizing architectural optimization with intrinsic plasticity (IP). The first step of each HRMAS optimization iteration optimizes the topology of the motif and inter-motif connectivity in SC-ML and the corresponding synaptic weights hierarchically. Specifically, the optimal number of neurons in the motif is optimized over a finite set of motif sizes. All possible intra-motif connections are considered and the type of each connection is optimized, which may lead to a sparser connectivity if the connection types of certain synapses are determined to be "non-existence". At the inter-motif level, a sparse motif-to-motif connectivity constraint is imposed: neurons in one motif are only allowed to be wired up with the corresponding neurons in the neighboring motifs as the Figure 1 shows. This locally connected topology will serve as a hard constraint in the subsequent optimization process. Inter-motif connections

199

200

201

202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

221

222

223

224

also fall under one of the three types ("inhibitory", "excitatory", "non-existence"). Hence, a greater level of sparsity is produced with the emergence of connections of type "non-existence". The second step in each 189 HRMAS iteration executes an unsupervised IP rule to stabilize the network function and mitigate potential 190 risks caused by architectural changes. 191

192 Figure 4 illustrates the incremental optimization strategy we adopt for the architectural parameters. Using the two-step optimization, initially all architectural parameters including motif size and connectivity are 193 194 optimized. After several training iterations, we choose the optimal motif size from a set of discrete options. As the most critical architectural parameter is set, we continue to optimize the remaining architectural 195 196 parameters defining connectivity, allowing fine-tuning of performance based on the chosen motif size.

Alternating Two-Step Optimization in HRMAS 2.3.2 197

The alternating two-step optimization in HRMAS is inspired by the evolution in neural development. As shown in Figure 2, neural circuits may experience weight changes through synaptic plasticity. Over a longer time scale, circuit architecture, i.e., connectivity, may evolve through learning and environmental changes. In addition, spontaneous firing behaviors of individual neurons may be adapted by intrinsic plasticity (IP). We are motivated by the important role of local IP mechanisms in stabilizing neuronal activity and coordinating structural changes to maintain proper circuit functions (Tien and Kerschensteiner, 2018). We view IP as a "fast-paced" self-adapting mechanism of individual neurons to react to and minimize the risks of weight and architectural modifications. As shown in Figure 3, we define the architectural parameters (motif size and intra/inter-motif connection types weights), synaptic weights, and intrinsic neuronal parameters as α , w, and β , respectively. Each HRMAS optimization iteration consists of two alternating steps. In the first step, we optimize α and w hierarchically based on gradient-based optimization using backpropagation (BP). In Figure 3, δ is the backpropagated error obtained via the employed BP method. In the second step, we use an unsupervised IP rule to adapt the intrinsic neuronal parameters of each neuron over a time window ("IP window") during which training examples are presented to the network. IP allows the neurons to respond to the weight and architectural changes introduced in the first step and mitigate possible risks caused by such changes. In Step 1 of the subsequent iteration, the error gradients w.r.t the synaptic weights and architectural parameters are computed based on the most recent values of β updated in the preceding iteration. In summary, the k-th HRMAS iteration solves a bi-level optimization problem:

$$\alpha^* = \arg\min_{\alpha} \mathcal{L}_{\text{valid}}(\alpha, w^*(\alpha), \beta^*)$$
 (5)

$$\alpha^* = \arg\min_{\alpha} \mathcal{L}_{\text{valid}}(\alpha, w^*(\alpha), \beta^*)$$
s.t.
$$\beta^* = \arg\min_{\beta} \mathcal{L}_{\text{ip}}(\alpha, w^*(\alpha), \beta_-^*),$$
s.t.
$$w^*(\alpha) = \arg\min_{w} \mathcal{L}_{\text{train}}(\alpha, w, \beta_-^*),$$
(6)

s.t.
$$w^*(\alpha) = \arg\min_{w} \mathcal{L}_{train}(\alpha, w, \beta_-^*),$$
 (7)

where \mathcal{L}_{valid} and \mathcal{L}_{train} are the loss functions defined based on the validation and training sets used to train α and w respectively; \mathcal{L}_{ip} is the local loss to be minimized by the IP rule as further discussed in Section 2.3.4; β_{-}^{*} are the intrinsic parameter values updated in the preceding (k-1)-th iteration; $w^{*}(\alpha)$ 219 220 denotes the optimal synaptic weights under the architecture specified by α .

In the implementation of HRMAS, architectural parameters and synaptic weights are optimized by the first step. The architectural parameters are defined as motif size and types of intra/inter-motif connections. The general architectural optimization is performed by generating architecture and evaluating the architecture by a standard training and validation process on data. The validation performance is used to

train the architectural parameters and generate a better structure. These steps are repeated until the optimal architecture is found. The first step of the k-th HRMAS iteration solves a bi-level optimization problem using BP:

$$min_{\alpha}\mathcal{L}_{valid}(\alpha, w^*(\alpha), \beta_{-}^*)$$
 (8)

s.t.
$$w^*(\alpha) = arg_w min \mathcal{L}_{train}(\alpha, w, \beta_-^*),$$
 (9)

where \mathcal{L}_{valid} and \mathcal{L}_{train} are the loss functions defined based on the validation and training sets used to train α and w respectively; β_{-}^{*} is the intrinsic parameter values updated in the preceding (k-1)-th iteration; $w^{*}(\alpha,\beta)$ denotes the optimal synaptic weights under the architecture specified by α . The second step of the k-th iteration solves the optimization problem below:

$$\beta^* = arg_{\beta}min\mathcal{L}_{ip}(\alpha^*, w^*, \beta) \tag{10}$$

232 \mathcal{L}_{ip} is the local loss to be minimized by the IP rule.

233 2.3.3 Gradient-based Optimization in HRMAS

234 2.3.3.1 Relaxing SC-ML layer's architectural parameters from discrete to continuous

Optimizing the weight and architectural parameters by solving the bi-level optimization problem of (5, 6, 7) can be computationally expensive. We adapt the recent method proposed in Liu et al. (2018) to reduce computational complexity by relaxing the discrete architectural parameters to continuous ones for efficient gradient-based optimization. Without loss of generality, we consider a multi-layered RSNN consisting of one or more SC-ML layers, where connections between layers are assumed to be feedforward. We focus on one SC-ML layer, as shown in Figure 5, to discuss the proposed gradient-based optimization.

The number of neurons in the SC-ML layer is fixed. The motif size is optimized such that each neuron is partitioned into a specific motif based on the chosen motif size. The largest white square in Figure 5 shows the layer-connectivity matrix of all intra-layer connections of the whole layer, where the dimension of the matrix corresponds to the neuron count of the layer. We superimpose three sets of smaller gray squares onto the layer-connectivity matrix, one for each of the three possible motif sizes of v_1 , v_2 , and v_3 considered. Choosing a particular motif size packs neurons in the layer into multiple motifs, and the corresponding gray squares illustrate the intra-motif connectivity introduced within the SC-ML layer.

The entry of the layer-connectivity matrix at row r and column i specifies the existence and nature of the connection from neuron r to neuron i. We consider multiple motif size and connection type choices during architectural search using continuous-valued parameterizations α^v and α^c_{ir} , respectively for each motif size v and connection type c. We relax the categorical choice of each motif size using a softmax over all possible options as $\hat{\alpha}^v$, and similarly relax the categorical choice of each connection type based on the corresponding motif size as $\hat{\alpha}^c_{ir}$:

$$\hat{\alpha}^{v} = \frac{exp(\alpha^{v})}{\sum_{v' \in \mathcal{V}} exp(\alpha^{v'})}, \quad \hat{\alpha}_{ir}^{c} = \frac{exp(\alpha_{ir}^{c})}{\sum_{c' \in \mathcal{C}} exp(\alpha_{ir}^{c'})}$$
(11)

Here, V and C are the set of all motif sizes and possible connection types, respectively; $\hat{\alpha}^v$ and $\hat{\alpha}^c_{ir}$ are the continuous-valued categorical choice of motif size v and connection type c, respectively, which can also be interpreted as the probability of selecting the corresponding motif size or connection type.

279

280

281

282

283

284

285

286

As in Figure 5, the synaptic weight of the connection from neuron r to neuron i is expressed as the summation of weights under all possible motif sizes and connection types weighted by the respective continuous-valued categorical choices (selection probabilities). In this paper, we use hat over the variable to denote the architectural parameter processed by softmax. Then, the task of architecture optimization is reduced to learn a set of continuous variables $\hat{\alpha} = \{\hat{\alpha}_{ir}^c, \hat{\alpha}^v\}$. With the continuous architectural parameters, a gradient-based method like BP is applicable to learn the recurrent connectivity.

Since IP rules are independent of the network architecture search problem, in following derivation, we do not express the IP method parameters β explicitly and express the term $\mathcal{L}_{valid}(\alpha, w^*(\alpha), \beta_-^*)$ as $\mathcal{L}_{valid}(\alpha, w^*(\alpha))$ for simplicity. In Liu et al. (2018), the bi-level optimization problem is simply approximated to a one-shot model to reduce the expensive computational cost of the inner optimization which can be expressed as

$$\nabla_{\hat{\alpha}} \mathcal{L}_{valid}(\hat{\alpha}, w^*(\hat{\alpha})) = \nabla_{\hat{\alpha}} \mathcal{L}_{valid}(\hat{\alpha}, w - \eta \nabla_w \mathcal{L}_{train}(w, \hat{\alpha})), \tag{12}$$

268 where η is the learning rate for a step of inner loop. Both the weights of the search network and the architectural parameters are trained by the BP method. The architectural gradient can be approximated by

$$\frac{d\mathcal{L}_{valid}}{d\hat{\alpha}}(\hat{\alpha}) = \nabla_{\hat{\alpha}}\mathcal{L}_{valid}(\hat{\alpha}, w^*) - \eta \nabla_{w}\mathcal{L}_{valid}(\hat{\alpha}, w^*) \nabla_{\hat{\alpha}, w}^2 \mathcal{L}_{train}(w^*, \hat{\alpha})). \tag{13}$$

The complexity is further reduced by using the finite difference approximation around $w^{\pm} = w \pm \epsilon \nabla_w \mathcal{L}_{valid}(\hat{\alpha}, w^*)$ for small perturbation ϵ to compute the gradient of $\nabla_{\hat{\alpha}} \mathcal{L}_{valid}(\hat{\alpha}, w^*)$. Finally the architectural updates in (13) can be calculated as

$$\frac{d\mathcal{L}_{valid}}{d\hat{\alpha}}(\hat{\alpha}) = \nabla_{\hat{\alpha}} \mathcal{L}_{valid}(\hat{\alpha}, w^*) - \frac{\eta}{2\epsilon} (\nabla_{\hat{\alpha}} \mathcal{L}_{train}(w^+, \hat{\alpha}) - \nabla_{\hat{\alpha}} \mathcal{L}_{train}(w^-, \hat{\alpha})). \tag{14}$$

273 2.3.3.2 Backpropagation via HRMAS framework

274 2.3.3.2.1 Integrating architectural parameterizations into the LIF model

Based on the leaky integrate-and-fire (LIF) neuron model in (3), the neuronal membrane voltage $u_i[t]$ of neuron i in the SC-ML layer at time t is given by integrating currents from all inter-layer inputs and intra-layer recurrent connections under all possible architectural parameterizations:

$$u_{i}[t] = (1 - \frac{1}{\tau})u_{i}[t - 1] + \frac{R}{\tau} \left(\sum_{j} w_{ij} a_{j}[t] + \sum_{v \in \mathcal{V}} (\hat{\alpha}^{v} \sum_{r}^{I_{i}^{v}} \sum_{c \in \mathcal{C}} (\hat{\alpha}_{ir}^{c} w_{ir}^{c} a_{r}[t - 1]))\right), \tag{15}$$

where R and τ are the resistance and time constant of the membrane, w_{ij} the synaptic weight from neuron j in the previous layer to neuron i, w_{ir}^c the recurrent weight from neuron r to neuron i of connection type c, and $a_j[t]$ the (unweighted) postsynaptic current (PSC) converted from spikes of neuron j through a synaptic model. To reduce clutter in the notation, we use I_i^v to denote the number of presynaptic connections afferent onto neuron i's input in the recurrent layer when choosing motif size v, which includes both inter and intra-motif connections and will be introduced in detail in the next paragraph. We further drop the explicit dependence of $\hat{\alpha}_{ir}^c$ on $\hat{\alpha}^v$. We assume feedforward connections have no time delay and recurrent connections have one time step delay. The response of neuron i obtained from recurrent connections is the summation of all the weighted recurrent inputs over the probabilities of connection types and motif sizes.

2.3.3.2.2 SC-ML's topology and scalability

288 In this section we formally describe the topology of SC-ML and discuss its scalability. I_i^v denote the number of presynaptic connections afferent onto neuron i's input in the recurrent layer when choosing 289 290 motif size v, and could be formally expressed as a union of inter-motif $I_{i,inter}^v$ and intra-motif $I_{i,intra}^v$ 291

neuron connections (We have omitted the superscript of connection type c for convenience):

$$I_i^v = I_{i,inter}^v \cup I_{i,intra}^v \tag{16}$$

Hence the recurrent input weight of neuron i be expressed by

$$w_{ir} = w_{ir}^{inter} \cup w_{ir}^{intra} \tag{17}$$

Let us consider a SC-ML layer with N neurons, divided into motif size = v, with N/v motifs within this layer. Let us denote the index of the motif by $k \in (0, 1, 2, ..., N/v - 1)$. Assuming the neuron i is located in the k_{th} motif (i.e.: $kv \le i \le kv + v - 1$), then the intra-layer recurrent connection into neuron i be expressed as

$$w_{ir}^{intra}$$
, where $r \in (kv, kv+1, kv+2, ..., kv+v-1)$

The inter-layer recurrent connection into neuron i be expressed as

$$w_{ir}^{inter}$$
, where $r \in (i - v, i + v)$

The figure expression is shown in Figure 1. The essence of SC-ML architecture design is to reduce the huge 293 search space of the recurrent matrix and improve optimization efficiency through biologically inspired 294 295 and carefully designed local recurrent connections as inductive bias. Hence, the SC-ML architecture can naturally adopt different inter and intra-motif topological connection patten across different layers, while 296 providing scalability. 297

Backpropagation in output layer 2.3.3.2.3

299 Through (15), the continuous architecture parameterizations influence the integration of input currents, and hence firing activities of neurons in all layers and affect the loss function defined at the output layer. 300 As such, the task of architecture optimization reduces to the one that learns the set of optimal continuous 301 302 variables $\hat{\alpha}^c$ and $\hat{\alpha}^v$. The final architecture is constructed by choosing the parameterizations with the highest 303 selection probabilities obtained from the optimization. During the learning, We define the loss function as

$$L = \sum_{k=0}^{T} E[t_k], \tag{18}$$

where T is the total time steps and $E[t_k]$ the loss at t_k . From (15) and (4), the membrane potential $u_i[t]$ of the neuron i at time t demonstrates contribution to all future fires and losses of the neuron through its PSC 305 $a_i[t]$. Therefore, the error gradient with respect to the presynaptic weight w_{ij} from neuron j to neuron i can 306

be defined as

$$\frac{\partial L}{\partial w_{ij}} = \sum_{k=0}^{T} \frac{\partial E[t_k]}{\partial w_{ij}} = \sum_{k=0}^{T} \sum_{m=0}^{k} \frac{\partial E[t_k]}{\partial u_i[t_m]} \frac{\partial u_i[t_m]}{\partial w_{ij}}$$

$$= \sum_{m=0}^{T} \frac{R}{\tau} a_j[t_m] \sum_{k=m}^{T} \frac{\partial E[t_k]}{\partial u_i[t_m]} = \sum_{m=0}^{T} \frac{R}{\tau} a_j[t_m] \delta_i[t_m],$$
(19)

where $\delta_i[t_m]$ denotes the error for neuron i at time t_m and is defined as:

$$\delta_i[t_m] = \sum_{k=m}^T \frac{\partial E[t_k]}{\partial u_i[t_m]} = \sum_{k=m}^T \frac{\partial E[t_k]}{\partial a_i[t_k]} \frac{\partial a_i[t_k]}{\partial u_i[t_m]}.$$
 (20)

In this work, the output layer is regular feedforward layer without recurrent connection. Therefore, the 309 weight w_{oj} of output neuron o is updated by

$$\frac{\partial L}{\partial w_{oj}} = \sum_{m=0}^{T} \frac{R}{\tau} a_j[t_m] \sum_{k=m}^{T} \frac{\partial E[t_k]}{\partial a_o[t_k]} \frac{\partial a_o[t_k]}{\partial u_o[t_m]},\tag{21}$$

where $\frac{\partial E[t_k]}{\partial a_o[t_k]}$ depends on the choice of the loss function.

2.3.3.2.4 Backpropagation in hidden layers

- 313 Now, we focus on the backpropagation in the recurrent hidden layer while the feedforward hidden layer
- case can be derived similarly. For a neuron i in SC-ML, in addition to the error signals from the next
- layer, the error backpropagated from the recurrent connections should also be taken into consideration. The
- backpropagated error can be calculated by:

$$\delta_{i}[t_{m}] = \sum_{k=m}^{T} \sum_{j=k}^{T} \frac{\partial a_{i}[t_{k}]}{\partial u_{i}[t_{m}]} \sum_{p=1}^{N_{p}} \left(\frac{\partial u_{p}[t_{k}]}{\partial a_{i}[t_{k}]} \frac{\partial E[t_{j}]}{\partial u_{p}[t_{k}]} \right) + \sum_{k=m}^{T} \sum_{j=k+1}^{T} \frac{\partial a_{i}[t_{k}]}{\partial u_{i}[t_{m}]} \sum_{r}^{N_{r}} \left(\frac{\partial u_{r}[t_{k}+1]}{\partial a_{i}[t_{k}]} \frac{\partial E[t_{j}]}{\partial u_{r}[t_{k}+1]} \right)$$

$$= \sum_{k=m}^{T} \frac{\partial a_{i}[t_{k}]}{\partial u_{i}[t_{m}]} \sum_{p=1}^{N} \left(\frac{R}{\tau} w_{pi} \delta_{p}[t_{k}] \right) + \sum_{k=m}^{T-1} \frac{\partial a_{i}^{(l)}[t_{k}]}{\partial u_{i}^{(l)}[t_{m}]} \sum_{v \in \mathcal{V}} \left(\hat{\alpha}^{v} \sum_{r}^{O_{i}^{v}} \sum_{c \in \mathcal{C}} \frac{R}{\tau} \hat{\alpha}_{ri}^{c} w_{ri}^{c} \delta_{r}[t_{k}+1] \right), \tag{22}$$

- where N_p and N_r are the number of neurons in the next layer and the number of neurons in this recurrent
- layer, respectively. δ_p and δ_r are the errors of the neuron p in the next layer and the error from the neuron r
- through the recurrent connection. O_i^v represents all the postsynaptic neurons of neuron i's outputs in the
- recurrent layer when choosing motif size v, which includes both inter and intra-motif connections. 320
- The key term in (22) is $\frac{\partial a[t_k]}{\partial u[t_m]}$ which reflects the effect of neuron's membrane potential on its output PSC. Due to the non-differentiable spiking events, it becomes the main difficulty for the BP of SNNs. 321
- 322
- Various approaches are proposed to handle this problem such as probability density function of spike 323
- state change Shrestha and Orchard (2018), surrogate gradient Neftci et al. (2019), and Temporal Spike 324
- Sequence Learning via Backpropagation (TSSL-BP) Zhang and Li (2020b). In our experiments, we adopt 325

336

337

338

339 340

341

342

343

344

345

346

347

348

349

350

351

352

353

354

355

356

357 358

359

360

the TSSL-BP method to calculated $\frac{\partial a[t_k]}{\partial u[t_m]}$. With the error backpropagated according to (22), the weights and architectural parameters can be updated by gradient descent as:

$$\Delta w_{ij} \propto \delta_{i}[t] \frac{R}{\tau} a_{j}[t], \quad \Delta \hat{\alpha}^{v} \propto \sum_{i}^{N_{r}} \delta_{i}[t] \frac{R}{\tau} \sum_{r}^{I_{i}^{v}} (\sum_{c \in \mathcal{C}} \hat{\alpha}_{ir}^{c} w_{ir}^{c} a_{r}[t-1]),$$

$$\Delta w_{ir}^{c} \propto \delta_{i}[t] \frac{R}{\tau} \sum_{v \in \mathcal{V}} (\hat{\alpha}^{v} \hat{\alpha}_{ir}^{c} a_{r}[t-1]), \quad \Delta \hat{\alpha}_{ir}^{c} \propto \delta_{i}[t] \frac{R}{\tau} \sum_{v \in \mathcal{V}} (\hat{\alpha}^{v} w_{ir}^{c} a_{r}[t-1]).$$
(23)

where $\delta_i[t]$ is the backpropagated error for neuron i at time t given in (22), N_r is the number of neurons in this recurrent layer, R and τ are the leaky resistance and membrane time constant, two intrinsic parameters adapted by the IP rule, $a_j[t]$ and $a_r[t]$ are the (unweighted) postsynaptic currents (PSCs) generated based on synpatic model by the presynaptic neuron j in the preceding layer and the r-th neuron in this recurrent layer, respectively.

333 2.3.4 Risk Minimizing Optimization with Intrinsic Plasticity

For architectural optimization of non-spiking RNNs, gradient-based methods are shown to be unstable in some cases due to misguided architectural changes and conversion from the optimized continuous-valued parameterization to a discrete architectural solution, hindering the final performance and demolishing the effectiveness of learning (Zela et al., 2019). Adaptive regularization which modifies the regularization strength (weight decay) guided by the largest eigenvalue of $\nabla_{\alpha}^2 \mathcal{L}_{valid}$ was proposed to address this problem (Zela et al., 2019). While this method shows promise for non-spiking RNNs, it is computationally intensive due to frequent expensive eigenvalue computation, severely limiting its scalability.

To address risks observed in architectural changes for RSNNs, we introduce a biologically-inspired risk-mitigation method. Biological circuits demonstrate that Intrinsic Plasticity (IP) is crucial in reducing such risks. IP is a self-regulating mechanism in biological neurons ensuring homeostasis and influencing neural circuit dynamics (Marder et al., 1996; Baddeley et al., 1997; Desai et al., 1999). IP is based on local neural firing activities and performs online adaptation with minimal additional computational overhead. It not only stabilizes neuronal activity but also coordinates connectivity and excitability changes across neurons to stabilize circuits (Maffei and Fontanini, 2009; Tien and Kerschensteiner, 2018). IP has been applied in spiking neural networks for locally regulating neuron activity (Lazar et al., 2007; Bellec et al., 2018). In Zhang et al. (2019), the application of IP mechanism significantly improves computational performance in terms of learning speed, accuracy, and robustness to input variations and noise. Fourati et al. (2020) proposes a deep echo state network that utilizes intrinsic plasticity to drive reservoir neuron activities to follow a desired Gaussian distribution, enabling the learning of discriminative EEG representations and demonstrating its effectiveness on emotion recognition benchmarks. Zhang et al. (2020) proposes a novel IP learning rule based on a soft-reset spiking neuron model, which ensures the neuron's membrane potential is mathematically continuous and differentiable. Experimental results demonstrate that the proposed IP rule can effectively improve the classification accuracy, inference speed, and noise robustness. Zhang et al. (2021) proposes input-driven and self-driven intrinsic IP learning rules for spiking convolutional neural networks (SCNNs), where IP updates occur only when a neuron receives input spikes or generates an output spike, respectively. Experiments show that the event-driven IP rules significantly reduce IP update operations and accelerate convergence while maintaining accuracy.

362

363

364

365

366

367

368

Drawing from these findings, we make use of IP for mitigating the risk of RSNN architectural modifications in this work. Our HRMAS framework integrates the IP rule into the architectural optimization, applied in the second step of each iteration. We adopt the SpiKL-IP rule (Zhang and Li, 2019a) for all recurrent neurons during architecture optimization. SpiKL-IP adapts the intrinsic parameters of a spiking neuron while minimizing the KL-divergence from the output firing rate distribution to a targeted exponential distribution. It both maintains a level of network activity and maximizes the information transfer for each neuron. We adapt leaky resistance and membrane time constant of each neuron using SpiKL-IP which effectively solves the optimization problem in (6) in an online manner. Specifically:

$$\Delta R = \frac{2y\tau V_{th} - W - V_{th} - \frac{1}{\mu}\tau V_{th}y^2}{RW}, \quad \Delta \tau = \frac{-1 + \frac{y}{\mu}}{\tau}, \quad W = \frac{V_{th}}{e^{\frac{1}{\tau y}} - 1}, \tag{24}$$

where μ is the desired mean firing rate, y the average firing rate of the neuron. Similar to biological neurons, we use the intracellular calcium concentration $\phi[t]$ as a good indicator of the averaged firing activity and y can be expressed with the time constant of calcium concentration τ_{cal} as

$$\phi_i[t] = (1 - \frac{1}{\tau_{cal}})\phi_i[t - 1] + s_i[t], \quad y_i[t] = \frac{\phi_i[t]}{\tau_{cal}}.$$
 (25)

We explicitly express the neuronal parameters R and τ of neuron i tuned through time as $R_i[t]$ and $\tau_i[t]$, since they are adjusted by the IP rule at each time step. They are updated by

$$R_i[t] = R_i[t-1] - \gamma \Delta R_i, \quad \tau_i[t] = \tau_i[t-1] - \gamma \Delta \tau_i, \tag{26}$$

where γ is the learning rate of the SpiKL-IP rule. By including time-variant neuronal parameters R and τ into (22) and (23), the one time step architectural parameter and weight updates change to

$$\delta_{i}[t_{m}] = \sum_{k=m}^{T} \frac{\partial a_{i}[t_{k}]}{\partial u_{i}[t_{m}]} \sum_{p=1}^{N} \left(\frac{R_{p}[t_{k}]}{\tau_{p}[t_{k}]} w_{pi} \delta_{p}[t_{k}]\right)
+ \sum_{k=m}^{T-1} \frac{\partial a_{i}^{(l)}[t_{k}]}{\partial u_{i}^{(l)}[t_{m}]} \sum_{v \in \mathcal{V}} \left(\hat{\alpha}^{v} \sum_{r}^{O_{i}^{v}} \sum_{c \in \mathcal{C}} \frac{R_{r}[t_{k}+1]}{\tau_{r}[t_{k}+1]} \hat{\alpha}_{ri}^{c} w_{ri}^{c} \delta_{r}[t_{k}+1]\right)$$
(27)

$$\Delta w_{ij} \propto \delta_{i}[t] \frac{R_{i}[t]}{\tau_{i}[t]} a_{j}[t], \quad \Delta \hat{\alpha}^{v} \propto \sum_{i}^{N_{r}} \delta_{i}[t] \frac{R_{i}[t]}{\tau_{i}[t]} \sum_{r}^{I_{v}^{v}} (\sum_{c \in \mathcal{C}} \hat{\alpha}_{ir}^{c} w_{ir}^{c} a_{r}[t-1]),$$

$$\Delta w_{ir}^{c} \propto \delta_{i}[t] \frac{R_{i}[t]}{\tau_{i}[t]} \sum_{v \in \mathcal{V}} (\hat{\alpha}^{v} \hat{\alpha}_{ir}^{c} a_{r}[t-1]), \quad \Delta \hat{\alpha}_{ir} \propto \delta_{i}[t] \frac{R_{i}[t]}{\tau_{i}[t]} \sum_{v \in \mathcal{V}} (\hat{\alpha}^{v} w_{ir}^{c} a_{r}[t-1]),$$

$$(28)$$

The proposed alternating two-step optimization of HRMAS is summarized in Algorithm 1. Architectural parameters α includes size of motif, and type of motif connections. They are optimized separately in two consecutive stages. We express here only a formal unification, for the sake of clarity in the architecture search problem.

Algorithm 1 Hybrid Risk-Mitigating Architectural Search

```
Initialize weights w, intrinsic parameters \beta, architectural parameters \alpha, and correspondingly \hat{\alpha}. 

repeat

Update \hat{\alpha} by

\eta_1 \nabla_{\hat{\alpha}} \mathcal{L}_{valid}(\hat{\alpha}, w - \eta_2 \nabla_w \mathcal{L}_{train}(\hat{\alpha}, w, \beta));

Update w by \eta_2 \nabla_w \mathcal{L}_{train}(\hat{\alpha}, w, \beta);
\beta \longleftarrow \text{SpiKL-IP}(\hat{\alpha}, w);

until converged
```

3 RESULTS

The proposed HRMAS optimized RSNNs with the SC-ML layer architecture and five motif size options 380 are evaluated on speech dataset TI46-Alpha (Liberman et al., 1991), neuromorphic speech dataset N-381 TIDIGITS (Anumula et al., 2018), neuromorphic video dataset DVS-Gesture (Amir et al., 2017), and 382 neuromorphic image dataset N-MNIST (Orchard et al., 2015). The performances are compared with 383 recently reported state-of-the-art manually designed architectures of SNNs and ANNs such as feedforward 384 SNNs, RSNNs, LSM, and LSTM. For the proposed work, the architectural parameters are optimized by 385 HRMAS with the weights trained on a training set and architectural parameters learned on a validation set 386 387 as shown in Algorithm 1. The accuracy of each HRMAS optimized network is evaluated on a separate testing set with all weights reinitialized. Table 2 shows all results. 388

3.1 Experimental Settings

390 3.1.1 Dataset

389

The proposed HRMAS framework with SC-ML is evaluated on speech dataset TI46-Alpha Liberman et al. (1991), neuromorphic speech dataset N-TIDIGITS Anumula et al. (2018), neuromorphic video dataset DVS-Gesture Amir et al. (2017), and neuromorphic image dataset N-MNIST Orchard et al. (2015). The performances are compared with several existing results on different structures of SNNs and ANNs such as feedforward SNNs, RSNNs, Liquid State Machine(LSM), LSTM, and so on.

396 3.1.2 Loss Function

For the BP method used in this work, the loss function can be defined by any errors that measure the 397 398 distance between the actual outputs and the desired outputs. In our experiments, since hundreds of time steps are required for simulating speech and neuromorphic inputs, we choose the accumulated output 399 PSCs to define the error which is similar to the firing count used in many existing works Jin et al. (2018); 400 401 Shrestha and Orchard (2018). We suppose the simulation time steps for a sample is T. In addition, for neuron o of the output layer, we define the desired output as $d_o = (d_o[t_0], d_o[t_1], \dots, d_o[t_N])$ and real output 402 as $a_o = (a_o[t_0], a_o[t_1], \dots, a_o[t_N])$ and d_o is manually determined. Therefore, the loss is determined by the 403 404 square error of the outputs

$$L = \sum_{k=1}^{T} E[t_k] = \sum_{k=1}^{T} \sum_{o}^{N^{(out)}} \frac{1}{2} (d_o[t_k] - a_o[t_k])^2$$
(29)

where $N^{(out)}$ is the number of neurons in the output layer and $E[t_k]$ is the error at time step t_k . is simply defined by the averaged loss through all the time steps:

$$E[t_k] \triangleq \sum_{o}^{N^{(out)}} \frac{1}{2} (d_o[t_k] - a_o[t_k])^2.$$
 (30)

With the loss function defined above, the error δ can be calculated for each layer according to (22). We use

- 08 a manually specified target output sequence to calculate the loss. Typically, we want neurons in a target
- 409 class to fire at every timestep with spike train output: (1,1, ...,1), while neurons in other classes are silenced
- 410 with spike train output: (0,0,...,0). Loss is then calculated by comparing the target spike train's PSC d_0
- 411 with the actual spike train's PSC a_o in 29.

412 3.1.3 Network architecture and hyperparameters

- In the SNNs of the experiments, the fully connected weights between layers are initialized by the He
- 414 Normal initialization proposed in He et al. (2015). The recurrent weights of excitatory connections are
- 415 initialized to 0.2 and tuned by the BP method. The weights of inhibitory connections are initialized to -2
- and fixed. The simulation step size is set to 1 ms. The parameters like thresholds and learning rate are
- 417 empirically tuned. No synaptic delay is applied for feedforward connections while recurrent connections
- 418 have 1 time step delay. No refractory period, normalization, or dropout is used. Adam Kingma and Ba
- 419 (2014) is adopted as the optimizer. The mean and standard deviation (std) of the accuracy reported is
- 420 obtained by repeating the experiments five times.
- Table 1 lists the typical constant values of parameters adopted in our experiments for each dataset. The
- 422 SC-ML size denotes the number of neurons in the SC-ML. In our experiments, each network contains one
- 423 SC-ML as the hidden layer. In addition, five motif sizes are predetermined before the experiment. The
- 424 HRMAS framework optimizes the motif size from one of the five options.

425 3.1.4 Traing process

- 426 Our experiments contain two phases. In the first phase, the weights are trained via the training set while
- 427 the validation set is used to optimize architectural parameters. In the second phase, the motif topology and
- 428 type of lateral connections are fixed after obtaining the optimal architecture. All the weights of the network
- 429 are reinitialized. Then, the new network is trained on the training set and tested on the testing set. The test
- 430 performance is reported in the paper. In addition, since all the datasets adopted in this paper only contain
- 431 training sets and testing sets, our strategy is to divide the training set. In the first phase, the training set is
- 432 equally divided into a training subset and a validation subset. Then, the architecture is optimized on these
- 433 subsets. In the second phase, since all the weights are reinitialized, we can train the weights with the full
- 434 training set and test on the testing set. Note that the testing set is only used for the final evaluation.

435 3.2 Performance of HRMAS

- Table 2 shows the results on the TI46-Alpha dataset. In order to verify the performance of the HRMAS
- 437 algorithm, we conducted five experiments on each dataset (using different initialization seeds) and recorded
- 438 the highest accuracy, average accuracy and standard deviation. The HRMAS-optimized RSNN has one
- 439 hidden SC-ML layer with 800 neurons, and outperforms all other models while achieving 96.44% accuracy
- 440 with mean of 96.08% and standard deviation (std) of 0.27% on the testing set. The proposed RSNN

- outperforms the LSM model in Wijesinghe et al. (2019) by 18.44%. It also outperforms the larger multi-
- 442 layered RSNN with more tunable parameters in Zhang and Li (2019b) trained by the spike-train level BP
- 443 (ST-RSBP) by 3.1%. Recently, Zhang and Li (2020a) demonstrated improved performances from manually
- designed RNNs with self-recurrent connections trained using the same TSSL-BP method. Our automated
- 445 HRMAS architectural search also produces better performing networks.
- We also show that a HRMAS-optimized RSNN with a 400-neuron SC-ML layer outperforms several
- state-of-the-art results on the N-TIDIGITS dataset (Zhang and Li, 2019b), achieving 94.66% testing
- 448 accuracy (mean: 94.27%, std: 0.35%). Our RSNN has more than a 3% performance gain over the widely
- 449 adopted recurrent structures of ANNs, the GRU and LSTM. It also significantly outperforms a feedforward
- 450 SNN with the same hyperparameters, achieving an accuracy improvement of almost 9.82%, demonstrating
- 451 the potential of automated architectural optimization.
- On DVS-Gesture and N-MNIST, our method achieves accuracies of 90.28% (mean: 88.40%, std: 1.71%)
- and 98.72% (mean: 98.60%, std: 0.08%), respectively. Table2 compares a HRMAS-optimized RSNN with
- 454 models including feedforward SNNs trained by TSSL-BP (Zhang and Li, 2020b) or STBP (Wu et al.,
- 455 2018) with the same size, and non-spiking ANNs vanilla LSTM (He et al., 2020). Note that although
- 456 our RSNN and the LSTM model have the same number of units in the recurrent layer, the LSTM model
- 457 has a much greater number of tunable parameters and a improved rate-coding-inspired loss function. Our
- 458 HRMAS-optimized model surpasses all other models. For a more intuitive understanding, Figure 6 presents
- 459 two examples of the motif topology optimized by HRMAS: motif sizes 2 in options [2, 4, 8, 16, 32] for the
- 460 N-MNIST dataset and motif size 16 in options [5, 10, 16, 25, 40] for the TI-Alpha dataset. We also shows in
- 461 the Figure 7 the weight matrix of the RSNN with SC-ML optimized by the HRMAS method. The original
- 462 fully connected recurrent matrix size is 800*800. We set the search space of motif size to [2,4,8,16,32]. In
- 463 five random experiments, the HRMAS optimization method always gave the search results of motif-size=2,
- 464 with similar inter/intra motif topology. This limits the huge recurrent matrix to a highly sparse band matrix
- 465 with non-zero values only near the diagonal, greatly reducing the search space, parameter amount, and
- 466 optimization difficulty.

3.3 Ablation Analysis

468 3.3.1 Ablation experiments of proposed components

We conduct ablation studies on the RSNN optimized by HRMAS for the TI46-Alpha dataset to reveal the contributions of various proposed techniques. When all proposed techniques are included, the HRMAS-

contributions of various proposed techniques. When all proposed techniques are included, the HRMASoptimized RSNN achieves 96.44% accuracy. In Table 3, removing of the IP rule from the second step of

- 472 the HRMAS optimization iteration visibly degrades the performance, showing the efficacy of intrinsic
- 473 plasticity for mitigating risks of architectural changes. A similar performance degradation is observed when
- 474 the sparse inter-motif connections are excluded from the SC-ML layer architecture. Without imposing a
- 475 structure in the hidden layer by using motifs as a basic building block, HRMAS can optimize all possible
- 476 connectivity types of the large set of 800 hidden neurons. However, this creates a large and highly complex
- 477 architectural search space, rendering a tremendous performance drop. Finally, we compare the HRMAS
- 478 model with an RSNN of a fixed architecture with full recurrent connectivity in the hidden layer. The
- application of the BP method is able to train the latter model since no architectural (motifs or connection
- 480 types) optimization is involved. However, albeit its significantly increased model complexity due to dense
- 481 connections, this model has a large performance drop in comparison with the RSNN fully optimized by
- 482 HRMAS. We provide additional data including ablation experiments, the computational resources required
- 483 by our method, and the IP rule's effect on performance during optimizing process in Section 3.3.1.

503

504

505 506

507

508

509

Ablation experiments of network parameters 484

We provided additional ablation experiments in Table 4, including: random search in the search space as 485 the baseline for HRMAS, effect of IP rule, HRMAS perfermence on larger network. Experimental results 486 show that: the HRMAS method shows consistent superiority (around 2%) over the random search baseline; 487 IP rule brings stable performance improvement (around 1.3%); our method can be efficiently extended to 488 networks with more neurons while providing good performance. 489

3.4 IP rule's effect on performance during optimizing process 490

We plotted the performance curve of the network optimization process on the TI46-alpha dataset. Figure 8 and Figure 9 show the loss and accuracy on the validation set respectively. The solid line and shading show 492 the mean and standard deviation of the 5 experiments. We conducted experiments with IP rule turned on 493 and off. We use green text to mark each phase of architecture optimization. 494

The experimental results show: 1. When the network architecture changes drastically, such as iteration = 495 750 (the cell size search ends and the connection type search starts), and iteration = 1500 (the connection 496 type search ends, the network is discretized and fine-tuned), the network There will be a slight performance 497 degradation. But it can be quickly improved to a higher level by the next stage of training. 2.It can be found 498 that IP method brings two benefits: improved network performance and a more stable training process. The 499 500 figures showed that the red solid line (mean) has always performed better than the blue solid line without the IP method; at the same time, the red shadow (standard deviation) has always been narrower than the 501 blue shadow, which means a more stable network architecture search process.

The effect of IP rules is mainly to stabilize the performance loss caused by architecture changes when the network architecture undergoes huge changes. Therefore, we found that at 750 epoch, that is, the network plays the most significant role from searching for cell size to searching for connection type: the loss distribution with ip rules (red shading) is much smaller than the loss distribution without ip rules (blue shading). At epoch 1500, since the fine-tuning phase does not involve drastic architectural changes, the role of IP rules is relatively limited.

The computational resources required for HRMAS

The proposed HRMAS bi-level optimization process is similar to DARTS, so the overall computational 510 complexity is similar to DARTS; the IP method is an localized unsupervised learning method and does not 511 constitute significant computational consumption. Furthermore, our proposed SC-ML topology greatly 512 reduces the search space. Specifically, as the Figure 7 shows, the HRMAS optimization of a SC-ML layer, 513 with n neurons, a motif size of s and n/s motifs, reduces the parameters that need to be optimized for 514 the recurrent connection matrix from $O(n^2)$ to O(sn): O(n/s) inter-motif connections + $O(n/s * s^2)$ 515 intra-motif connections + O(n) neuron hyperparameters. Generally, $s \ll n$, which reduces the parameter 516 space of recurrent connections to linear growth with the neuron numbers, allowing our algorithm scale 517 well. Specifically, a complete training process of a RSNN with 800 neurons hidden layer for TI46-alpha 519 dataset, including 150 epoch for cell size search, 150 epoch for connection type search and 400 epoch for finetune, takes 4 hours on single NVIDIA GeForce RTX 3090 GPU.

CONCLUSION

We present an RSNN architecture based on SC-ML layers composed of multiple recurrent motifs with 521 sparse inter-motif connections as a solution to constructing large recurrent spiking neural models. We

- 523 further propose the automated architectural optimization framework HRMAS hybridizing the "evolution" of
- 524 the architectural parameters and corresponding synaptic weights based on backpropagation and biologically-
- 525 inspired mitigation of risks of architectural changes using intrinsic plasticity. We show that HRMAS-
- 526 optimized RSNNs impressively improve performance on four datasets over the previously reported state-of-
- 527 the-art RSNNs and SNNs. Notably, our HRMAS framework can be easily extended to more flexible network
- 528 architectures, optimizing sparse and scalable RSNN architectures. By sharing the PyTorch implementation
- 529 of our HRMAS framework, this work aims to foster advancements in high-performance RSNNs for both
- 530 general-purpose and dedicated neuromorphic computing platforms, potentially inspiring innovative designs
- 531 in brain-inspired recurrent spiking neural models and their energy-efficient deployment.

CONFLICT OF INTEREST STATEMENT

- 532 The authors declare that the research was conducted in the absence of any commercial or financial
- 533 relationships that could be construed as a potential conflict of interest.

AUTHOR CONTRIBUTIONS

- 534 WZ and PL developed the theoretical approach for HRMAS. WZ and HG implemented HRMAS and
- related learning rules and performed the simulation studies. WZ, HG and PL wrote the paper.

FUNDING

- 536 The author(s) declare that financial support was received for the research, authorship, and/or publication
- 537 of this article. This work supported by the National Science Foundation (NSF). Any opinions, findings,
- conclusions or recommendations expressed in this material are those of the authors and do not necessarily
- 539 reflect the views of NSF, UCSB and their contractors. This material is based upon work supported by the
- 540 National Science Foundation under Grant Nos. 1948201 and 2310170.

REFERENCES

- 541 Amir, A., Taba, B., Berg, D., Melano, T., McKinstry, J., Di Nolfo, C., et al. (2017). A low power, fully
- event-based gesture recognition system. In Proceedings of the IEEE Conference on Computer Vision
- and Pattern Recognition. 7243–7252
- 544 Anumula, J., Neil, D., Delbruck, T., and Liu, S.-C. (2018). Feature representations for neuromorphic audio
- spike streams. Frontiers in neuroscience 12, 23
- 546 Baddeley, R., Abbott, L. F., Booth, M. C., Sengpiel, F., Freeman, T., Wakeman, E. A., et al. (1997).
- Responses of neurons in primary and inferior temporal visual cortices to natural scenes. Proceedings of
- the Royal Society of London B: Biological Sciences 264, 1775–1783
- 549 Bellec, G., Salaj, D., Subramoney, A., Legenstein, R., and Maass, W. (2018). Long short-term memory
- and learning-to-learn in networks of spiking neurons. In Advances in Neural Information Processing
- 551 Systems. 787–797
- 552 Buzsaki, G. (2006). Rhythms of the Brain (Oxford University Press)
- 553 Chakraborty, B. and Mukhopadhyay, S. (2023). Heterogeneous recurrent spiking neural network for
- spatio-temporal classification. Frontiers in Neuroscience 17, 994517
- 555 Chen, L., Li, X., Zhu, Y., Wang, H., Li, J., Liu, Y., et al. (2023). Intralayer-connected spiking neural
- network with hybrid training using backpropagation and probabilistic spike-timing dependent plasticity.
- 557 International Journal of Intelligent Systems 2023

- 558 Cho, K., van Merrienboer, B., Gülçehre, Ç., Bahdanau, D., Bougares, F., Schwenk, H., et al. (2014).
- Learning phrase representations using rnn encoder-decoder for statistical machine translation. In
- 560 EMNLP
- 561 Desai, N. S., Rutherford, L. C., and Turrigiano, G. G. (1999). Plasticity in the intrinsic excitability of
- 562 cortical pyramidal neurons. Nature neuroscience 2, 515
- 563 Elsken, T., Metzen, J. H., Hutter, F., et al. (2019). Neural architecture search: A survey. J. Mach. Learn.
- 564 Res. 20, 1–21
- 565 Fourati, R., Ammar, B., Jin, Y., and Alimi, A. M. (2020). Eeg feature learning with intrinsic plasticity based
- deep echo state network. In 2020 international joint conference on neural networks (IJCNN) (IEEE),
- 567 1–8
- 568 Gerstner, W. and Kistler, W. M. (2002). Spiking neuron models: Single neurons, populations, plasticity
- 569 (Cambridge university press)
- 570 Graves, A., Mohamed, A.-r., and Hinton, G. (2013). Speech recognition with deep recurrent neural
- networks. In 2013 IEEE international conference on acoustics, speech and signal processing (IEEE),
- 572 6645–6649
- 573 He, K., Zhang, X., Ren, S., and Sun, J. (2015). Delving deep into rectifiers: Surpassing human-level per-
- formance on imagenet classification. In <u>Proceedings of the IEEE international conference on computer</u>
- 575 vision. 1026–1034
- 576 He, W., Wu, Y., Deng, L., Li, G., Wang, H., Tian, Y., et al. (2020). Comparing snns and rnns on
- 577 neuromorphic vision datasets: Similarities and differences. Neural Networks 132, 108–120
- 578 Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. Neural computation 9, 1735–1780
- 579 Jaeger, H. (2001). The "echo state" approach to analysing and training recurrent neural networks-with an
- erratum note. Bonn, Germany: German National Research Center for Information Technology GMD
- Technical Report 148, 13
- 582 Jin, Y., Zhang, W., and Li, P. (2018). Hybrid macro/micro level backpropagation for training deep spiking
- neural networks. Advances in neural information processing systems 31, 7005–7015
- 584 [Dataset] Kim, Y., Li, Y., Park, H., Venkatesha, Y., and Panda, P. (2022). Neural architecture search for
- 585 spiking neural networks
- 586 Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. arXiv preprint
- 587 arXiv:1412.6980
- 588 Ko, H., Hofer, S. B., Pichler, B., Buchanan, K. A., Sjöström, P. J., and Mrsic-Flogel, T. D. (2011).
- Functional specificity of local synaptic connections in neocortical networks. Nature 473, 87–91
- 590 Lazar, A., Pipa, G., and Triesch, J. (2007). Fading memory and time series prediction in recurrent networks
- with different forms of plasticity. Neural Networks 20, 312–322
- 592 [Dataset] Liberman, M., Amsler, R., Church, K., Fox, E., Hafner, C., Klavans, J., et al. (1991). TI 46-word
- 593 LDC93S9
- 594 Liu, H., Simonyan, K., and Yang, Y. (2018). Darts: Differentiable architecture search. In International
- 595 Conference on Learning Representations
- 596 Maass, W. (1997). Networks of spiking neurons: the third generation of neural network models. Neural
- 597 networks 10, 1659–1671
- 598 Maass, W., Natschläger, T., and Markram, H. (2002). Real-time computing without stable states: A new
- framework for neural computation based on perturbations. Neural computation 14, 2531–2560
- 600 Maes, A., Barahona, M., and Clopath, C. (2020). Learning spatiotemporal signals using a recurrent spiking
- network that discretizes time. PLoS computational biology 16, e1007606

- Maffei, A. and Fontanini, A. (2009). Network homeostasis: a matter of coordination. <u>Current opinion in</u> neurobiology 19, 168–173
- Marder, E., Abbott, L., Turrigiano, G. G., Liu, Z., and Golowasch, J. (1996). Memory from the dynamics of intrinsic membrane currents. Proceedings of the national academy of sciences 93, 13481–13486
- 606 [Dataset] Na, B., Mok, J., Park, S., Lee, D., Choe, H., and Yoon, S. (2022). Autosnn: Towards energy-efficient spiking neural networks
- Neftci, E. O., Mostafa, H., and Zenke, F. (2019). Surrogate gradient learning in spiking neural networks:
- Bringing the power of gradient-based optimization to spiking neural networks. <u>IEEE Signal Processing</u>
- 610 Magazine 36, 51–63
- 611 Orchard, G., Jayawant, A., Cohen, G. K., and Thakor, N. (2015). Converting static image datasets to
- spiking neuromorphic datasets using saccades. Frontiers in neuroscience 9, 437
- 613 [Dataset] Pan, W., Zhao, F., Han, B., Dong, Y., and Zeng, Y. (2023). Emergence of brain-inspired 614 small-world spiking neural network through neuroevolution
- Perin, R., Berger, T. K., and Markram, H. (2011). A synaptic organizing principle for cortical neuronal groups. Proceedings of the National Academy of Sciences 108, 5419–5424
- 617 Real, E., Aggarwal, A., Huang, Y., and Le, Q. V. (2019). Regularized evolution for image classifier
- architecture search. In Proceedings of the aaai conference on artificial intelligence. vol. 33, 4780–4789
- 619 Seeman, S. C., Campagnola, L., Davoudian, P. A., Hoggarth, A., Hage, T. A., Bosma-Moody, A., et al.
- 620 (2018). Sparse recurrent excitatory connectivity in the microcircuit of the adult mouse and human cortex.
- 621 Elife 7, e37349
- 622 Shrestha, S. B. and Orchard, G. (2018). Slayer: Spike layer error reassignment in time. In Advances in
- Neural Information Processing Systems. 1412–1421
- 624 Srinivasan, G., Panda, P., and Roy, K. (2018). Spilinc: Spiking liquid-ensemble computing for unsupervised
- speech and image recognition. <u>Frontiers in neuroscience</u> 12
- 626 Tian, S., Qu, L., Wang, L., Hu, K., Li, N., and Xu, W. (2021). A neural architecture search based framework
- for liquid state machine design. Neurocomputing 443, 174–182
- 628 Tien, N.-W. and Kerschensteiner, D. (2018). Homeostatic plasticity in neural development. <u>Neural</u>
- 629 development 13, 1–7
- 630 Voelker, A., Kajić, I., and Eliasmith, C. (2019). Legendre memory units: Continuous-time representation
- in recurrent neural networks. In Advances in Neural Information Processing Systems. 15570–15579
- 632 Wang, Q. and Li, P. (2016). D-lsm: Deep liquid state machine with unsupervised recurrent reservoir tuning.
- In 2016 23rd International Conference on Pattern Recognition (ICPR) (IEEE), 2652–2657
- 634 Wijesinghe, P., Srinivasan, G., Panda, P., and Roy, K. (2019). Analysis of liquid ensembles for enhancing
- the performance and accuracy of liquid state machines. Frontiers in Neuroscience 13, 504
- 636 Wistuba, M., Rawat, A., and Pedapati, T. (2019). A survey on neural architecture search. <u>arXiv preprint</u>
- 637 arXiv:1905.01392
- 638 Wu, Y., Deng, L., Li, G., Zhu, J., and Shi, L. (2018). Spatio-temporal backpropagation for training
- high-performance spiking neural networks. Frontiers in neuroscience 12, 331
- 640 Zela, A., Elsken, T., Saikia, T., Marrakchi, Y., Brox, T., and Hutter, F. (2019). Understanding and
- robustifying differentiable architecture search. In International Conference on Learning Representations
- 642 Zhang, A., Gao, Y., Niu, Y., Li, X., and Chen, Q. (2020). Intrinsic plasticity for online unsupervised
- learning based on soft-reset spiking neuron model. IEEE Transactions on Cognitive and Developmental
- 644 Systems 15, 337–347
- 645 Zhang, A., Li, X., Gao, Y., and Niu, Y. (2021). Event-driven intrinsic plasticity for spiking convolutional
- neural networks. IEEE Transactions on Neural Networks and Learning Systems 33, 1986–1995

- Zhang, A., Zhou, H., Li, X., and Zhu, W. (2019). Fast and robust learning in spiking feed-forward neural networks based on intrinsic plasticity mechanism. Neurocomputing 365, 102–112
- Zhang, W. and Li, P. (2019a). Information-theoretic intrinsic plasticity for online unsupervised learning in spiking neural networks. Frontiers in neuroscience 13, 31
- Zhang, W. and Li, P. (2019b). Spike-train level backpropagation for training deep recurrent spiking neural
 networks. In Advances in Neural Information Processing Systems. 7800–7811
- Zhang, W. and Li, P. (2020a). Skip-connected self-recurrent spiking neural networks with joint intrinsic
 parameter and synaptic weight training. arXiv preprint arXiv:2010.12691
- Zhang, W. and Li, P. (2020b). Temporal spike sequence learning via backpropagation for deep spiking
 neural networks. Advances in Neural Information Processing Systems 33
- Zhang, W. and Li, P. (2021). Spiking neural networks with laterally-inhibited self-recurrent units. In <u>2021</u>
 International Joint Conference on Neural Networks (IJCNN) (IEEE), 1–8
- Zhang, Y., Li, P., Jin, Y., and Choe, Y. (2015). A digital liquid state machine with biologically inspired
 learning and its application to speech recognition. <u>IEEE transactions on neural networks and learning</u>
 systems 26, 2635–2649
- Zhou, Y., Jin, Y., and Ding, J. (2020). Surrogate-assisted evolutionary search of spiking neural architectures
 in liquid state machines. Neurocomputing 406, 12–23
- Zoph, B. and Le, Q. V. (2017). Neural architecture search with reinforcement learning. In <u>5th International</u>
 Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference
 Track Proceedings (OpenReview.net)

Table 1. P	arameters	settings.
------------	-----------	-----------

Parameter	TI46-Alpha	N-TIDIGITS	DvsGesture	N-MNIST
$ au_m$	16 ms	64 ms	64 ms	16 ms
$ au_s$	8 ms	8 ms	8 ms	8 ms
$ au_{cal}$	16 ms	16 ms	16 ms	16 ms
learning rate	0.0005	0.0005	0.0001	0.0005
Batch Size	50	50	20	50
Time steps	100	300	400	100
Epochs for searching	300	200	60	30
Epochs for testing	400	400	150	100
SC-ML size	800	800	512	512
Motif size options	[5, 10, 1	6, 25, 40]	[2, 4, 8,	16, 32]

Table 2. Accuracy on TI46-Alpha, N-TIDIGITS, DVS-Gesture and N-MNIST. HeNHeS result is from (Chakraborty and Mukhopadhyay, 2023)

Dataset	Network Structure	Learning Rule	Hidden Layers	Best
TI46-Alpha	LSM (Wijesinghe et al., 2019)	Non-spiking BP	2000	78%
	RSNN (Zhang and Li, 2019b)	ST-RSBP	400 - 400 - 400	93.35%
	Sr-SNN (Zhang and Li, 2020a)	TSSL-BP	400 - 400 - 400	94.62%
	This work	TSSL-BP	800	96.44%
N-TIDIGITS	GRU (Anumula et al., 2018)	Non-spiking BP	200 - 200 - 100	90.90%
	Phase LSTM (Anumula et al., 2018)	Non-spiking BP	250 - 250	91.25%
	RSNN (Zhang and Li, 2019b)	ST-RSBP	400 - 400 - 400	93.90%
	Feedforward SNN	TSSL-BP	400	84.84%
	This work	TSSL-BP	400	94.66%
DVS-Gesture	Feedforward SNN (He et al., 2020)	STBP	P4 - 512	87.50%
	LSTM (He et al., 2020)	Non-spiking BP	P4 - 512	88.19%
	HeNHeS	STDP	500	90.15%
	Feedforward SNN	TSSL-BP	P4 - 512	88.19%
	This work	TSSL-BP	P4 - 512	90.28%
N-MNIST	Feedforward SNN (He et al., 2020)	STBP	512	98.19%
	RNN (He et al., 2020)	Non-spiking BP	512	98.15%
	LSTM (He et al., 2020)	Non-spiking BP	512	98.69%
	ELSM(Pan et al., 2023)	Non-spiking BP	8000	97.23%
	This work	TSSL-BP	512	98.72%

Table 3. Ablation studies of HRMAS on TI46-Alpha

Setting	Accuracy
Full HRMAS	96.44%
Without IP	95.20%
Without motif	88.35%
Without inter-motif connections	95.73%
Fully connected RSNN	94.10%

Table 4. Test Accuracy on TI46-Alpha, obtained by repeating 5 times with different random seeds, including: HRMAS perfermence on Larger network, effect of IP rule, random search as the baseline architecture.

Arch Optimization	Learning Rule	SC-ML Sizes	Best	Mean	Std
HRMAS (with IP)	TSSL-BP	800	96.44%	96.08%	0.27 %
HRMAS (with IP)	TSSL-BP	1600	96.26%	-	-
HRMAS (with IP)	TSSL-BP	2400	96.58%	-	-
HRMAS (with IP)	TSSL-BP	3200	96.45%	-	-
HRMAS (w/o IP)	TSSL-BP	800	95.17%	94.74%	0.32%
Random	TSSL-BP	800	94.47%	94.18%	0.30%

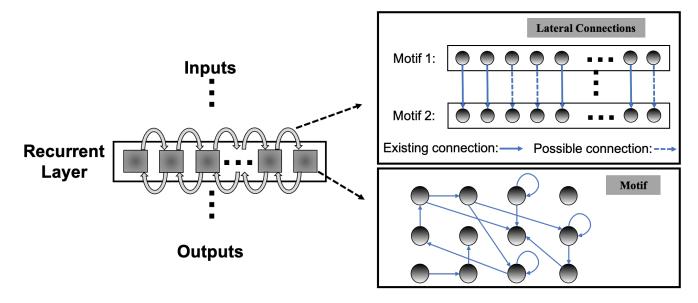


Figure 1. Sparsely-Connected Recurrent Motif Layer.

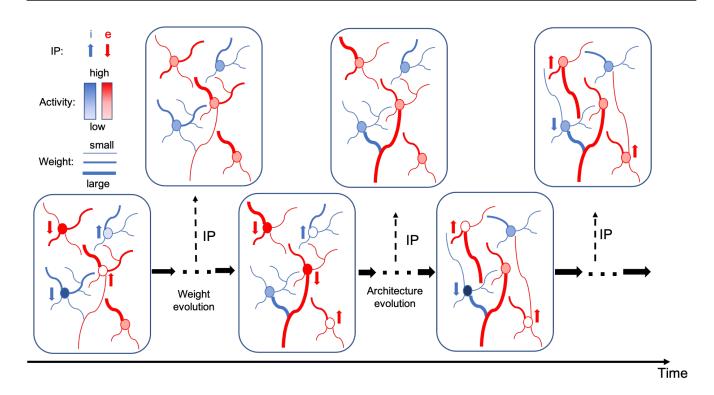


Figure 2. Evolution in neural development.

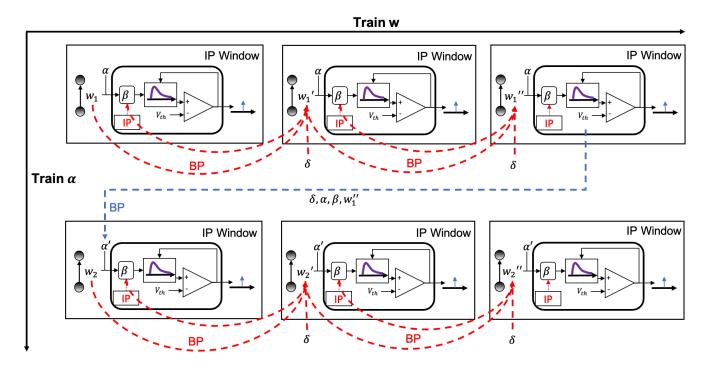


Figure 3. Proposed HRMAS.

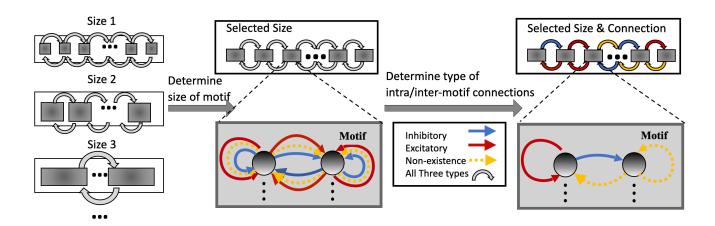


Figure 4. Architectural optimization in HRMAS.

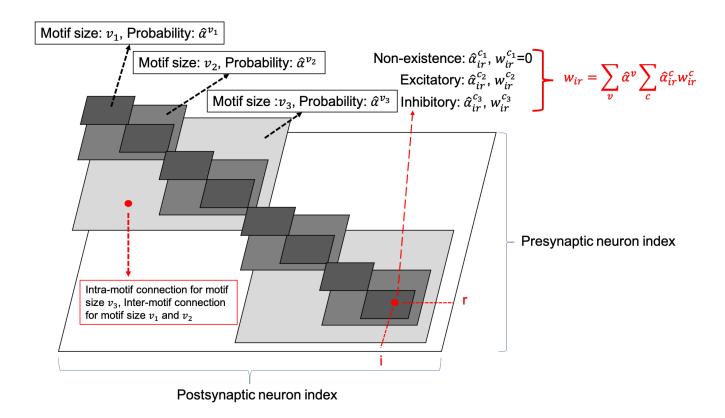
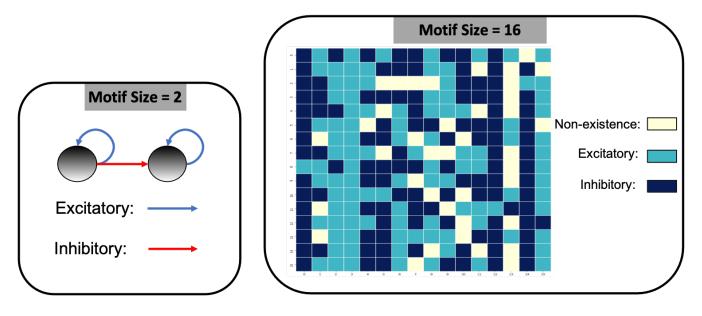


Figure 5. SC-ML with relaxed architectural parameters.



Optimized motif for N-MNIST

Optimized motif for TI46-Alpha

Figure 6. Optimized motif topologies.

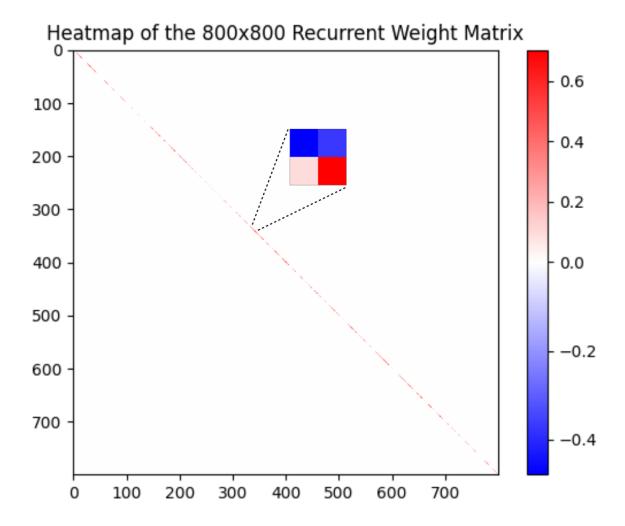


Figure 7. Recurrent Weight Matrix after optimization by HRMAS.

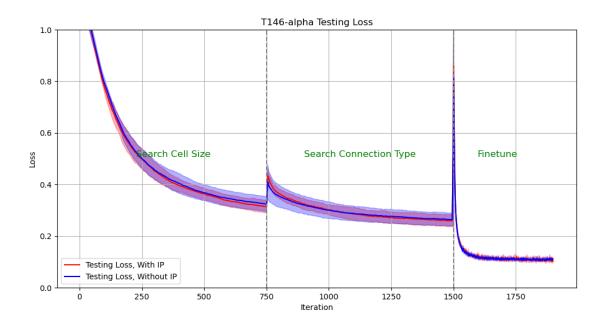


Figure 8. Test Loss in architectural optimization in HRMAS. The solid line and shading show the mean and standard deviation of the 5 experiments. We conducted experiments with ip rule turned on (red) and off (blue). It can be found that IP method brings two benefits: improved network performance and a more stable training process. The figures showed that the red solid line (mean) has lower loss than the blue solid line without the IP method; at the same time, the red shadow (standard deviation) has always been narrower than the blue shadow, which means a more stable network architecture search process.

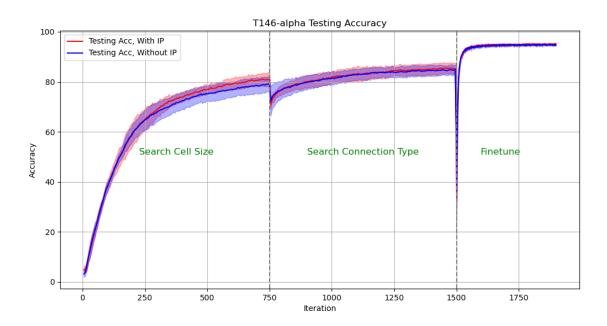


Figure 9. Test Accuracy in architectural optimization in HRMAS. The solid line and shading show the mean and standard deviation of the 5 experiments. We conducted experiments with ip rule turned on (red) and off (blue). It can be found that IP method brings two benefits: improved network performance and a more stable training process. The figures showed that the red solid line (mean) has higher accuracy than the blue solid line without the IP method; at the same time, the red shadow (standard deviation) has always been narrower than the blue shadow, which means a more stable network architecture search process.