

Research Article

Alexandr Udeneev, Petr Babkin, and Oleg Bakhteev

Surrogate assisted diversity estimation in neural ensemble search.

<https://doi.org/10.1515/sample-YYYY-XXXX>

Received Month DD, YYYY; revised Month DD, YYYY; accepted Month DD, YYYY

Abstract: Automated search for optimal neural architectures (NAS) is expensive, and extending it to ensembles (NES) can lead to exponential growth in computational cost due to the combinatorial number of architecture combinations. To address this, we propose a surrogate-based approach: each candidate architecture is encoded as a graph, its predictions on a held-out dataset form the training set for a surrogate model, and this model guides an efficient NES framework that selects both diverse and high-performing networks. Our final ensemble matches human-level accuracy on CIFAR-10 and outperforms other one-shot NES methods, demonstrating the practical effectiveness of the approach.

Keywords: NES, surrogate function, triplet loss.

1 Introduction

Neural network ensembles often demonstrate better accuracy compared to single models, especially in classification and regression tasks [1, 2]. This fact gives rise to the problem of constructing an efficient ensemble of models (NES) [3]. NES, in turn, relies on Neural Architecture Search (NAS) methods, which are extensively studied and applied to search for individual neural network architectures, such as evolutionary algorithms [4, 5], reinforcement learning [6–8], and Bayesian optimization [9, 10]. Selecting an optimal architecture for even a single model is a challenging task, particularly when considering data-specific constraints and computational limitations [11].

The simplest approach for ensemble construction is DeepEns [12]. It involves a random search for several architectures, which are then combined into an ensemble. Despite its simplicity in implementation and hyperparameter tuning, this method is computationally expensive. More sophisticated adaptation of NAS techniques are presented in some recent works [3, 13, 14], which are designed to efficiently combine multiple networks into an ensemble.

Our research also adapts ideas from NAS for NES, specifically utilisation a surrogate function [15–17]. Some modern NAS methods widely use surrogate functions to estimate architecture quality without requiring full model training. These functions significantly reduce computational costs, expanding the applicability of such methods. For example, in [15], evolutionary algorithms were proposed in combination with surrogate models for real-time semantic segmentation. In [17], a Surrogate-assisted Multiobjective Evolutionary-based Algorithm (SaMEA) is used for 3D medical image segmentation.

In this work, we propose a method for constructing neural network ensembles using a surrogate function that accounts for both model classification accuracy and architectural diversity. Diversity is crucial because ensembles consisting of similar models often fail to provide a significant performance gain. Two surrogate functions are used in the work: the first to represent the architecture in the latent space [18], the second to predict the accuracy of the architecture. Since a neural network architecture is represented as a graph, using a Graph Attention Network (GAT) [19] as a surrogate function [20] seems natural. To train it to predict model diversity, we use Triplet Loss [21], similar to [18]. We validate this approach on CIFAR-10, demonstrating the effectiveness of the surrogate function for predicting diversity and constructing ensembles.

We claim that ensembles constructed in this manner achieve state-of-the-art accuracy compared to one-shot NES algorithms, such as DeepEns [12].

Main Contributions:

- 1) We propose a method for encoding the DARTS [22] search space into a representation suitable for training a Graph Attention Network [19] (GAT), where graph nodes correspond to operations within the network.
- 2) We propose a way for training the surrogate function to predict the diversity of architectures.
- 3) We adapt surrogate functions for ensemble construction, taking into account both predictive performance and architectural diversity.

2 Problem statement

2.1 Neural Architecture Search

We define:

$$\mathcal{V} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\} \subset \mathbb{R}^d, \quad \mathcal{O} = \{o_1, \dots, o_K\}, \quad o_k : \mathbb{R}^d \rightarrow \mathbb{R}^d, \quad \mathcal{E} \subseteq \mathcal{V} \times \mathcal{V}$$

where \mathcal{V} is the set of node feature-vectors in \mathbb{R}^d , \mathcal{O} is a finite set of elementwise operations (e.g., convolutions, poolings), and \mathcal{E} is the set of directed edges between nodes.

We then define the space of feasible architectures as the set of all acyclic directed graphs over \mathcal{V} with operations on edges:

$$\mathcal{A} = \{(\mathcal{V}, E) \mid E \subseteq \mathcal{E}, (\mathcal{V}, E) \text{ is a DAG}, \forall (u \rightarrow v) \in E, \text{ assign } o_{u \rightarrow v} \in \mathcal{O}\}.$$

Equivalently, an architecture $\alpha \in \mathcal{A}$ can be represented by

$$(V, \{(u \rightarrow v, o_{u \rightarrow v}) \mid (u \rightarrow v) \in E\}),$$

where E induces an acyclic graph on \mathcal{V} and each edge $(u \rightarrow v)$ is labeled by an operation $o_{u \rightarrow v} \in \mathcal{O}$.

Denote \mathcal{L}_{train} and \mathcal{L}_{val} as the training and validation losses, respectively. The NAS problem can then be formulated as the search for an optimal architecture α^* that minimizes $\mathcal{L}_{val}(\alpha^*, \omega^*)$, under the constraint that the weights are obtained by minimizing the training loss:

$$\omega^* = \arg \min_{\omega \in \mathcal{W}} \mathcal{L}_{train}(\alpha^*, \omega)$$

This can be expressed as the following optimization problem:

$$\begin{aligned} & \min_{\alpha \in \mathcal{A}} \mathcal{L}_{val}(\omega^*(\alpha), \alpha) \\ \text{s.t. } & \omega^*(\alpha) = \arg \min_{\omega \in \mathcal{W}} \mathcal{L}_{train}(\omega, \alpha) \end{aligned} \tag{1}$$

The primary challenge in this optimization lies in the immense search space of possible architectures (e.g., in DARTS [22], it is approximately 10^{25}).

2.2 Neural Ensemble Search

The primary objective of NES is to find an optimal ensemble of neural networks whose architectures lie within the NAS search space.

As before, we denote $\alpha \in \mathcal{A}$ as a network architecture and $\omega(\alpha)$ as its corresponding parameters. The action of this network on an input x is denoted by $f_\alpha(x, \omega(\alpha))$. Let $S \subset \mathcal{A}$ be a subset of architectures. Then, the NES problem can be formally described as follows:

$$\begin{aligned} \min_S \mathcal{L}_{val} \left(\frac{1}{|S|} \sum_{\alpha \in S} f_{\alpha}(x, \omega^*(\alpha)) \right) \\ \text{s.t. } \forall \alpha \in S : \omega^*(\alpha) = \arg \min_{\omega(\alpha)} \mathcal{L}_{train}(f_{\alpha}(x, \omega(\alpha))) \end{aligned} \quad (2)$$

Thus, in addition to searching over a vast number of architectures, we now also need to find the optimal ensemble composition.

3 Method

In this work, we consider the transformation of the architecture space proposed in DARTS [22] for application in Graph Attention Networks (GAT) (see Section 3.1). In Section 3.2, we present the architecture of the surrogate function and describe its working principle, while Section 3.3 provides a detailed discussion of the ensemble construction method based on this surrogate function.

3.1 Architecture Search Space

A critical aspect of our methodology involves constructing a training dataset for the surrogate model. While we adopt the DARTS framework [22], which defines normal and reduction cells, our approach diverges fundamentally in how these cells are generated. Specifically, instead of optimizing cells through continuous relaxation, we generate both cell types via discrete random sampling.

Each cell is represented as a directed acyclic graph (DAG) comprising n nodes and m edges, where every edge corresponds to an operation selected from the DARTS operation set \mathcal{O} . Formally, a cell architecture is defined as:

$$\alpha_{\text{cell}} = (G, \mathbf{o}), \quad G = (\mathcal{V}, \mathcal{E}), \quad \mathbf{o} = \{o_e\}_{e \in \mathcal{E}}, \quad o_e \in \mathcal{O}, \quad (3)$$

where \mathcal{V} denotes the set of n latent nodes, \mathcal{E} represents the m directed edges, and \mathcal{O} is the predefined operation space.

The full architecture space \mathcal{A} consists of pairs of such DAGs (denoted as normal and reduction cells):

$$\mathcal{A} = \{(G_1, \mathbf{o}_1) \times (G_2, \mathbf{o}_2) \mid G_1, G_2 \text{ are DAGs, } |\mathcal{V}_1| = |\mathcal{V}_2| = n, |\mathcal{E}| = m, \mathbf{o}_1, \mathbf{o}_2 \in \mathcal{O}^m\}. \quad (4)$$

To train the surrogate model, we generate N architectures, evaluate their performance on a fixed validation dataset, and compile the training dataset:

$$\mathcal{D}_{\text{train}} = \{(\alpha_i, \mathbf{y}_i, \text{acc}_i)\}_{i=1}^N, \quad (5)$$

where \mathbf{y}_i is the vector of predictions from architecture α_i on the validation set, and acc_i is its validation accuracy.

Unlike DARTS, which employs continuous relaxation and gradient-based optimization, our method relies on discrete sampling to explore the search space. This strategy enables broader diversity in architecture generation at the expense of computational efficiency.

Let $n = 5, m = 10$ and use DARTS search space with 7 operations, then the total number of unique cell architectures is computed as:

$$|\mathcal{A}_{\text{cell}}| = \binom{2}{2} \binom{3}{2} \binom{4}{2} \binom{5}{2} \cdot |\mathcal{O}|^8, \quad |\mathcal{O}| = 7, \quad (6)$$

resulting in a complete architecture space of:

$$|\mathcal{A}| = |\mathcal{A}_{\text{cell}}|^2 \approx 10^{18}. \quad (7)$$

This vast combinatorial space renders exhaustive search computationally intractable, necessitating efficient surrogate-guided exploration.

3.2 Surrogate Function

In order to construct the ensemble described in Section 3.3, we need to predict both the performance and the similarity of candidate architectures. However, due to the enormous size of the architecture space, obtaining these characteristics via full training is infeasible.

To address this, we employ surrogate functions—models that, given an architectural representation, predict key properties of neural networks. Formally, we define

$$f : \mathcal{A} \rightarrow \mathbb{R}^d,$$

where d is the dimension of the latent space. In particular, we consider two surrogates:

$$f_{\text{acc}} : \mathcal{A} \rightarrow [0, 1], \quad f_{\text{sim}} : \mathcal{A} \rightarrow \mathbb{R}^d.$$

Our surrogate models are graph neural networks based on the Graph Attention (GAT) mechanism. Each surrogate consists of N sequential GAT layers with residual connections and GraphNorm, followed by a global pooling aggregation.

For accuracy prediction f_{acc} , a final fully-connected layer with sigmoid activation is used.

For similarity embedding f_{sim} , the pre-classification output of the network defines a latent vector in \mathbb{R}^d . During training, we optimize f_{sim} using a triplet loss on architecture embeddings, ensuring that architectures with similar performance and structural properties are mapped to nearby points in the latent space, while dissimilar architectures are pushed apart.

Figure 1 illustrates the overall architecture: N GAT layers (each with residual connections), global pooling, and two separate output heads—sigmoid for f_{acc} and identity for f_{sim} .

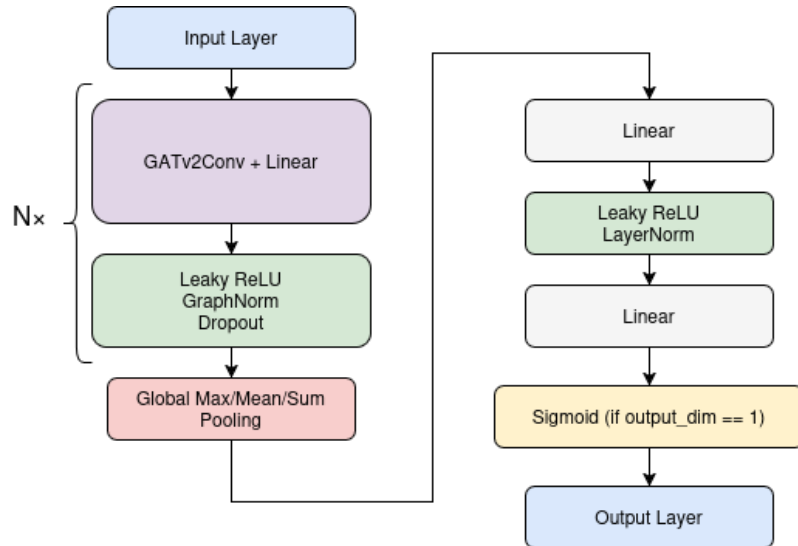


Fig. 1: Architecture of the surrogate functions: N GAT layers (each with residual + GraphNorm), global pooling, and output heads for accuracy and similarity embedding.

Currently, each architecture in the dataset is represented as a directed acyclic graph (DAG), where the nodes correspond to latent representations and the edges to operations, similar to the format presented in

NAS-Bench-201 [23]. For example, the DARTS architectures consist of connected normal and reduction cells (in a 2:1 ratio), as shown in Figure 4.

Subsequently, it is more convenient to adopt an alternative graph representation: the *nodes* represent operations and the *edges* their corresponding latent embeddings. In this way, the architecture is transformed into the NAS-Bench-101 format [24] (see Figure 5).

The conversion is carried out as follows:

1. Each edge of the original graph is transformed into a node in the new graph; the label of this node indicates the operation present on the corresponding edge in the original graph.
2. An oriented edge is drawn from the node with label o_i to the node with label o_j if, in the original graph, the operation o_i acts from a node x to a node y , and the operation o_j acts from node y to a node z (with x , y , and z being arbitrary nodes).

We encode the operations as one-hot vectors. Since our dataset consists of graphs (architectural encodings), we adopt a Graph Attention Network (GAT) as the surrogate function, as in [25].

For predicting accuracy, supervised learning can be employed; however, to train the model to predict similarity, we utilize the Triplet Loss [21]. As in [18], Triplet Loss is employed to learn the latent representations of architectures. Similarity between two models can be defined, for instance, by comparing the fraction of identical predictions or by computing the average distance between their output distributions using a suitable metric for multivariate distributions, such as a divergence or distance function in the probability space (e.g., Jensen–Shannon divergence or Hellinger distance). This allows us to construct a similarity matrix of size $N \times N$.

In this article, we construct a similarity matrix based on model responses on the validation dataset. The matrix is constructed as follows:

1. For each of the N models, predictions are computed on a fixed validation dataset consisting of K examples. Denote the predictions of model M_i by the vector

$$\mathbf{y}^{(i)} = \left(y_1^{(i)}, \dots, y_K^{(i)} \right).$$

2. For every pair of models (M_i, M_j) , where $i, j = 1, \dots, N$, the fraction of matching predictions is computed:

$$s_{ij} = \frac{1}{K} \sum_{k=1}^K \mathbb{I} \left(y_k^{(i)} = y_k^{(j)} \right),$$

where \mathbb{I} is the indicator function.

3. The values s_{ij} form the symmetric similarity matrix $\mathbf{S} \in [0, 1]^{N \times N}$, with each entry \mathbf{S}_{ij} representing the degree of similarity between models M_i and M_j .
4. To discretize the similarity for the purposes of Triplet Loss, the matrix \mathbf{S} is converted into a discrete matrix $\mathbf{M} \in \{-1, 0, 1\}^{N \times N}$ according to the following rule:

$$\mathbf{M}_{ij} = \begin{cases} 1, & \text{if } s_{ij} > q_p, \\ -1, & \text{if } s_{ij} < q_n, \\ 0, & \text{otherwise,} \end{cases}$$

where q_p and q_n are predetermined thresholds corresponding to the upper and lower quantiles of the distribution of s_{ij} values.

Algorithm 1 Training the surrogate model for architecture diversity

Input: f_{sim} : an untrained surrogate model; \mathcal{B} : a set of architectures of pretrained models; N : the number of architectures; \mathbf{M} : a discrete similarity matrix of size $N \times N$; n : the number of training epochs; **optimizer**: the optimization algorithm; m : the margin parameter for the Triplet Loss.

Output: Trained surrogate model f_{sim}

```

1 for  $i \leftarrow 1$  to  $n$  do
2   for  $j \leftarrow 1$  to  $N$  do
3      $\mathcal{P}_j \leftarrow \{k \mid \mathbf{M}[j, k] = 1\}$  // Sample a positive example
4      $k_p \leftarrow \text{UniformSample}(\mathcal{P}_j)$ 
5      $\mathcal{N}_j \leftarrow \{k \mid \mathbf{M}[j, k] = -1\}$  // Sample a negative example
6      $k_n \leftarrow \text{UniformSample}(\mathcal{N}_j)$ 
7     // Compute embeddings
8      $\mathbf{e}_a \leftarrow f_{sim}(\mathcal{B}[j])$ 
9      $\mathbf{e}_p \leftarrow f_{sim}(\mathcal{B}[k_p])$ 
10     $\mathbf{e}_n \leftarrow f_{sim}(\mathcal{B}[k_n])$ 
11    // Compute the triplet loss
12     $\mathcal{L} \leftarrow \max(0, \|\mathbf{e}_a - \mathbf{e}_p\|_2^2 - \|\mathbf{e}_a - \mathbf{e}_n\|_2^2 + m)$ 
13    // Optimization step
14    optimizer.zero_grad()
15     $\mathcal{L}.backward()$ 
16    optimizer.step()
17  end
18 end
19 return  $f_{sim}$ 

```

The training procedure of f_{sim} (Algorithm 1) proceeds as follows. For each epoch and for each anchor architecture, we first sample one positive and one negative example uniformly from the corresponding rows of the similarity matrix \mathbf{M} . Next, the anchor, positive, and negative architectures are embedded via f_{sim} , and a triplet loss with margin m is computed:

$$\mathcal{L} = \max\left(0, \|\mathbf{e}_a - \mathbf{e}_p\|_2^2 - \|\mathbf{e}_a - \mathbf{e}_n\|_2^2 + m\right).$$

Finally, we perform the standard backward-propagation step—zeroing gradients, backpropagating \mathcal{L} , and updating the model parameters with the chosen optimizer. This optimization scheme encourages the model to map similar architectures closer together and dissimilar ones at least margin m apart in the latent space.

3.3 Ensemble construction

The algorithm 2 begins by initializing three sets: \mathcal{A}_{pool} , which stores candidate architectures whose predicted accuracy exceeds a predefined threshold α ; \mathcal{E}_{pool} , which holds the corresponding diversity embeddings; and \mathcal{S}_{pool} , which contains the predicted accuracy scores obtained from surrogate models.

At each iteration, a batch of N candidate architectures, denoted by \mathcal{A}_{cand} , is generated. For each architecture, both the predicted accuracy and a representation in the latent diversity space are computed using the surrogate functions f_{acc} and f_{div} , respectively. Only those architectures with predicted accuracy above the threshold α are added to the pool.

Once the number of architectures in \mathcal{A}_{pool} reaches or exceeds M , the diversity embeddings \mathcal{E}_{pool} are clustered into K groups using the KMeans algorithm. For each cluster center, the architecture whose embedding lies closest to the center (in Euclidean distance) is selected. The resulting K architectures form the final ensemble, aiming to balance both predictive performance and diversity.

Algorithm 2 Ensemble Construction

Input: K : ensemble size; M : minimum pool size ($M \geq K$); N : number of architectures generated per iteration; α : accuracy threshold; f_{acc} : surrogate model for accuracy prediction; f_{div} : surrogate model for diversity embeddings;

Output: $\mathcal{A}_{\text{best}}$: set of K selected architectures

```

1  $\mathcal{A}_{\text{pool}} \leftarrow \emptyset$  // Initialize architecture pool
2  $\mathcal{E}_{\text{pool}} \leftarrow \emptyset$  // Initialize diversity embedding pool
3  $\mathcal{S}_{\text{pool}} \leftarrow \emptyset$  // Initialize pool of predicted accuracy scores
4 while  $|\mathcal{A}_{\text{pool}}| < M$  do
5    $\mathcal{A}_{\text{cand}} \leftarrow \text{GenerateArchs}(N)$  // Generate N random candidate architectures
6    $s \leftarrow f_{\text{acc}}(\mathcal{A}_{\text{cand}})$  // Predict accuracy for candidates
7    $e \leftarrow f_{\text{div}}(\mathcal{A}_{\text{cand}})$  // Predict diversity embeddings for candidates
8   for  $(a, s, e) \in \text{zip}(\mathcal{A}_{\text{cand}}, s, e)$  do
9     if  $s \geq \alpha$  then
10       // Updating pools
11        $\mathcal{A}_{\text{pool}} \leftarrow \mathcal{A}_{\text{pool}} \cup \{a\}$ 
12        $\mathcal{E}_{\text{pool}} \leftarrow \mathcal{E}_{\text{pool}} \cup \{e\}$ 
13        $\mathcal{S}_{\text{pool}} \leftarrow \mathcal{S}_{\text{pool}} \cup \{s\}$ 
14     end
15   end
16  $\mathcal{C} \leftarrow \text{KMeans}(\mathcal{E}_{\text{pool}}, K)$  // Cluster diversity embeddings into K groups
17  $\mathcal{A}_{\text{best}} \leftarrow \emptyset$  // Initialize set of best architectures
18 foreach  $c \in \mathcal{C}$  do
19    $\mu_c \leftarrow c.\text{center}$ 
20    $i^* \leftarrow \arg \min_i \|\mathcal{E}_{\text{pool}}[i] - \mu_c\|_2$  // Find index of embedding closest to cluster center
21    $\mathcal{A}_{\text{best}} \leftarrow \mathcal{A}_{\text{best}} \cup \{\mathcal{A}_{\text{pool}}[i^*]\}$  // Add selected architecture to best set
22   Remove index  $i^*$  from  $\mathcal{A}_{\text{pool}}$ ,  $\mathcal{E}_{\text{pool}}$ , and  $\mathcal{S}_{\text{pool}}$  // Remove selected index from pools
23 end
24 return  $\mathcal{A}_{\text{best}}$ 

```

The selected models are presented on Figure 2

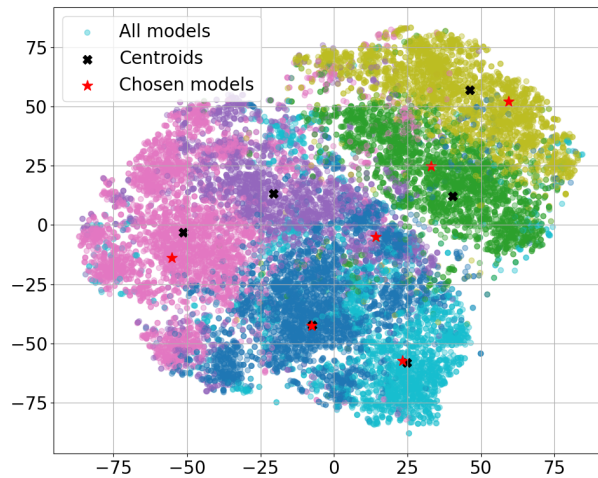


Fig. 2: Selected models for ensemble.

4 Computational Experiment

In this section we present the experimental results as well as the metrics used for comparison. In Section 4.1, we describe the dataset collection procedure. Section 4.2 details the training setup for the surrogate functions. Finally, in Section 4.3, we compare our proposed method against DeepEnsemble [12] and Random Search [26].

4.1 Dataset Construction

The models in our dataset were generated according to the following algorithm:

1. Split the CIFAR-10 dataset into training and validation subsets with an 80%/20% ratio.
2. Sample “normal” and “reduction” cells from the DARTS search space, ensuring each node has exactly two incoming edges from previous nodes.
3. Assemble a DARTS architecture consisting of two normal cells followed by one reduction cell.
4. Train each sampled architecture on the training subset.
5. For each trained model, record:
 - its architectural description,
 - its predictions on the validation set,
 - its validation accuracy.

Table 1 summarizes the training hyperparameters for all sampled models:

Tab. 1: Training Hyperparameters for Dataset Models

Hyperparameter	Value
Number of epochs	80
Optimizer	Adam
Learning rate	0.001
Training set size	48 000
Validation set size	12 000
Batch size	256

4.2 Training of the Surrogate Functions

In our experiment each cell contains $n = 5$ nodes, with two outgoing edges per node assigned random operations from \mathcal{O} , yielding $m = 10$ edges per cell. The operation set \mathcal{O} includes:

- Separable convolutions (kernel sizes: 3×3 , 5×5)
- Dilated separable convolutions (kernel sizes: 3×3 , 5×5)
- 3×3 max pooling
- 3×3 average pooling
- Identity operation

We employed a graph attention network (GAT) with four convolutional layers and two fully connected layers at the output, i.e. N=4 (see Figure 1). In both surrogate training regimes, the learning rate was halved every 5 epochs.

Full models are constructed by stacking normal and reduction cells in a 2:1 ratio, starting from a predefined number of initial channels. Each model combines one normal and one reduction cell sampled independently from \mathcal{A} .

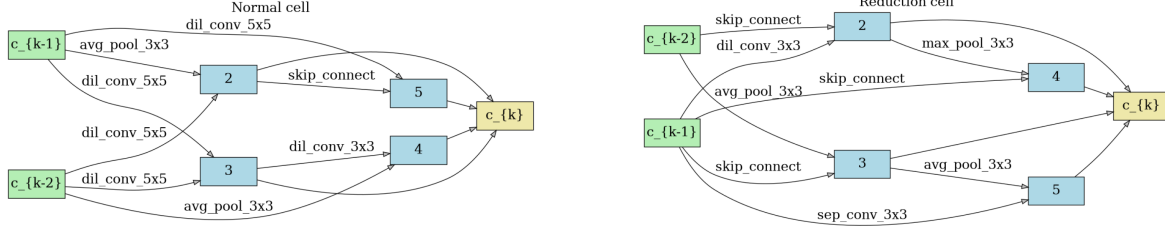


Fig. 3: Illustration of normal and reduction cell structures used in our architecture search space. The normal cell (left) preserves spatial dimensions while the reduction cell (right) halves them.

Table 2 lists the training hyperparameters for both surrogate models.

In the diversity regime, the GAT surrogate has an output dimensionality of 128, employs 4 attention heads and a low dropout rate of 10%, and is trained for 30 epochs with an initial learning rate of 0.001 scheduled by a cosine annealing scheduler down to a final learning rate of 1×10^{-6} . For accuracy prediction, the surrogate outputs a single scalar via sigmoid activation, uses 16 attention heads and a higher dropout rate of 40%, and is trained for 50 epochs with an initial learning rate of 0.01 under the same cosine scheduler ending at 1×10^{-6} . Both models are optimized with Adam on the same train/validation split of 1040/260 instances.

Tab. 2: Training Hyperparameters for Surrogate Models

Hyperparameter	Diversity Surrogate	Accuracy Surrogate
Number of epochs	30	50
Optimizer	Adam	Adam
Learning rate	0.001	0.01
Dropout	0.1	0.4
Heads	4	16
Training set size	1040	1040
Validation set size	260	260

4.3 Experimental Results on Ensemble Construction

We compare our surrogate-assisted NES ensemble construction method against DARTS [26] and a pure Random Search baseline [26]. Our ensemble comprises six models. Training was performed for 100 epochs with a batch size of 96 images, using the Adam optimizer (learning rate 0.025, cosine learning scheduler, weight decay 10^{-4}). The total training time was 11 hours on NVIDIA P100 GPU.

Table 3 presents the test accuracies of the compared methods.

Tab. 3: Comparison of ensemble test accuracies.

Model	Test Accuracy (%)
Surrogate-assisted NES (ours)	94.2
DeepEns (DARTS) [26]	97.38
Random Search [26]	97.29

5 Conclusion

In this study, we demonstrated that surrogate diversity functions can be effectively learned from embedding distances: we observed strong negative correlations between the Euclidean distance of embeddings and the proportion of identical responses in the similarity matrix (Pearson's $r = -0.39$; Spearman's $\rho = -0.37$). Building on this finding, we developed an ensemble construction algorithm that simultaneously leverages surrogate accuracy and surrogate similarity functions to select complementary models.

The primary limitation of our approach lies in the substantial computational cost of training a large number of base models to generate the required dataset for surrogate training.

For future work, it would be worthwhile to investigate the square root of the Jensen–Shannon divergence as an alternative similarity metric. Unlike simple response-proportion measures, a similarity matrix based on Jensen–Shannon distance may yield a validation loss that continues to decrease, potentially improving the quality of diversity estimates. Additionally, integrating these empirical results within the theoretical framework of the bias–variance tradeoff could provide deeper insights into the interplay between model accuracy and diversity in ensemble performance.

References

- [1] P.N. Suganthan Ye Ren, Le Zhang. Ensemble classification and regression—recent developments, applications and future directions [review article]. *IEEE Computational Intelligence Magazine*, 11(1):41–53, Feb 2016. ISSN 1556-603X. 10.1109/mci.2015.2471235. URL <https://doi.org/10.1109/mci.2015.2471235>.
- [2] Lars Kai Hansen and Peter Salamon. Neural network ensembles. *IEEE Trans. Pattern Anal. Mach. Intell.*, 12(10): 993–1001, 1990. 10.1109/34.58871.
- [3] Sheheryar Zaidi, Arber Zela, Thomas Elsken, Chris C. Holmes, Frank Hutter, and Yee Whye Teh. Neural ensemble search for uncertainty estimation and dataset shift. In Marc'Aurelio Ranzato, Alina Beygelzimer, Yann N. Dauphin, Percy Liang, and Jennifer Wortman Vaughan, editors, *Advances in Neural Information Processing Systems 34: Annual Conference on Neural Information Processing Systems 2021, NeurIPS 2021, December 6-14, 2021, virtual*, pages 7898–7911, 2021. URL <https://proceedings.neurips.cc/paper/2021/hash/41a6fd31aa2e75c3c6d427db3d17ea80-Abstract.html>.
- [4] Esteban Real, Sherry Moore, Andrew Selle, Saurabh Saxena, Yutaka Leon Suematsu, Jie Tan, Quoc V Le, and Alexey Kurakin. Large-scale evolution of image classifiers. In *International conference on machine learning*, pages 2902–2911. PMLR, 2017.
- [5] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In *Proceedings of the aaai conference on artificial intelligence*, volume 33, pages 4780–4789, 2019.
- [6] Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. In *5th International Conference on Learning Representations, ICLR 2017, Toulon, France, April 24-26, 2017, Conference Track Proceedings*. OpenReview.net, 2017. URL <https://openreview.net/forum?id=r1Ue8Hcxg>.
- [7] Sirui Xie, Huihui Zheng, Chunxiao Liu, and Liang Lin. Snas: stochastic neural architecture search. *arXiv preprint arXiv:1812.09926*, 2018.
- [8] Yuqiao Liu, Yanan Sun, Bing Xue, Mengjie Zhang, Gary G. Yen, and Kay Chen Tan. A survey on evolutionary neural architecture search. *IEEE Trans. Neural Networks Learn. Syst.*, 34(2):550–570, 2023. 10.1109/TNNLS.2021.3100554.
- [9] Haifeng Jin, Qingquan Song, and Xia Hu. Auto-keras: An efficient neural architecture search system. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 1946–1956, 2019.
- [10] Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabas Poczos, and Eric P Xing. Neural architecture search with bayesian optimisation and optimal transport. *Advances in neural information processing systems*, 31, 2018.
- [11] Ankit Kumar Saroj Kumar Pandey Neeraj varshney Teekam Singh Chetan Swarup, Kamred Udham Singh. Brain tumor detection using cnn, alexnet amp; googlenet ensembling learning approaches. *Electronic Research Archive*, 31(5): 2900–2924, 2023. ISSN 2688-1594. 10.3934/era.2023146. URL <https://doi.org/10.3934/era.2023146>.
- [12] Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. *Advances in neural information processing systems*, 30, 2017.
- [13] Yao Shu, Yizhou Chen, Zhongxiang Dai, and Bryan Kian Hsiang Low. Neural ensemble search via bayesian sampling. In James Cussens and Kun Zhang, editors, *Proceedings of the Thirty-Eighth Conference on Uncertainty in Artificial Intelligence*, volume 180 of *Proceedings of Machine Learning Research*, pages 1803–1812. PMLR, 01–05 Aug 2022. URL <https://proceedings.mlr.press/v180/shu22a.html>.

- [14] Haibin Ling Minghao Chen, Jianlong Fu. One-shot neural ensemble architecture search by diversity-guided search space shrinking. *2021 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 16525–16534, Jun 2021. 10.1109/cvpr46437.2021.01626. URL <https://doi.org/10.1109/cvpr46437.2021.01626>.
- [15] Zhichao Lu, Ran Cheng, Shihua Huang, Haoming Zhang, Changxiao Qiu, and Fan Yang. Surrogate-assisted multi-objective neural architecture search for real-time semantic segmentation. *CoRR*, abs/2208.06820, 2022. 10.48550/ARXIV.2208.06820.
- [16] Zhichao Lu, Kalyanmoy Deb, Erik D. Goodman, Wolfgang Banzhaf, and Vishnu Naresh Boddeti. Nsganetv2: Evolutionary multi-objective surrogate-assisted neural architecture search. In Andrea Vedaldi, Horst Bischof, Thomas Brox, and Jan-Michael Frahm, editors, *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part I*, volume 12346 of *Lecture Notes in Computer Science*, pages 35–51. Springer, 2020. 10.1007/978-3-030-58452-8_3.
- [17] Maria G. Baldeon Calisto and Susana K. Lai-Yuen. Emonas-net: Efficient multiobjective neural architecture search using surrogate-assisted evolutionary algorithm for 3d medical image segmentation. *Artif. Intell. Medicine*, 119:102154, 2021. 10.1016/J.ARTMED.2021.102154.
- [18] Ferrante Neri Yu Xue, Zhenman Zhang. Similarity surrogate-assisted evolutionary neural architecture search with dual encoding strategy. *Electronic Research Archive*, 32(2):1017–1043, 2024. ISSN 2688-1594. 10.3934/era.2024050. URL <https://doi.org/10.3934/era.2024050>.
- [19] Petar Velickovic, Guillem Cucurull, Arantxa Casanova, Adriana Romero, Pietro Liò, and Yoshua Bengio. Graph attention networks. *CoRR*, abs/1710.10903, 2017. 10.48550/arxiv.1710.10903. URL <http://arxiv.org/abs/1710.10903>.
- [20] Wei Wen, Hanxiao Liu, Yiran Chen, Hai Li, Gabriel Bender, and Pieter-Jan Kindermans. Neural predictor for neural architecture search. In *European conference on computer vision*, pages 660–676. Springer, 2020.
- [21] Florian Schroff, Dmitry Kalenichenko, and James Philbin. Facenet: A unified embedding for face recognition and clustering. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 815–823, 2015.
- [22] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: differentiable architecture search. *CoRR*, abs/1806.09055, 2018. 10.48550/arxiv.1806.09055. URL <http://arxiv.org/abs/1806.09055>.
- [23] Xuanyi Dong and Yi Yang. Nas-bench-201: Extending the scope of reproducible neural architecture search. *arXiv preprint arXiv:2001.00326*, 2020.
- [24] Chris Ying, Aaron Klein, Esteban Real, Eric Christiansen, Kevin Murphy, and Frank Hutter. Nas-bench-101: Towards reproducible neural architecture search. *CoRR*, abs/1902.09635, 2019. 10.48550/arxiv.1902.09635. URL <http://arxiv.org/abs/1902.09635>.
- [25] Wei Wen, Hanxiao Liu, Hai Li, Yiran Chen, Gabriel Bender, and Pieter-Jan Kindermans. Neural predictor for neural architecture search. *CoRR*, abs/1912.00848, 2019. 10.48550/arxiv.1912.00848. URL <http://arxiv.org/abs/1912.00848>.
- [26] Liam Li and Ameet Talwalkar. Random search and reproducibility for neural architecture search. In *Uncertainty in artificial intelligence*, pages 367–377. PMLR, 2020.

A NAS-Bench formats

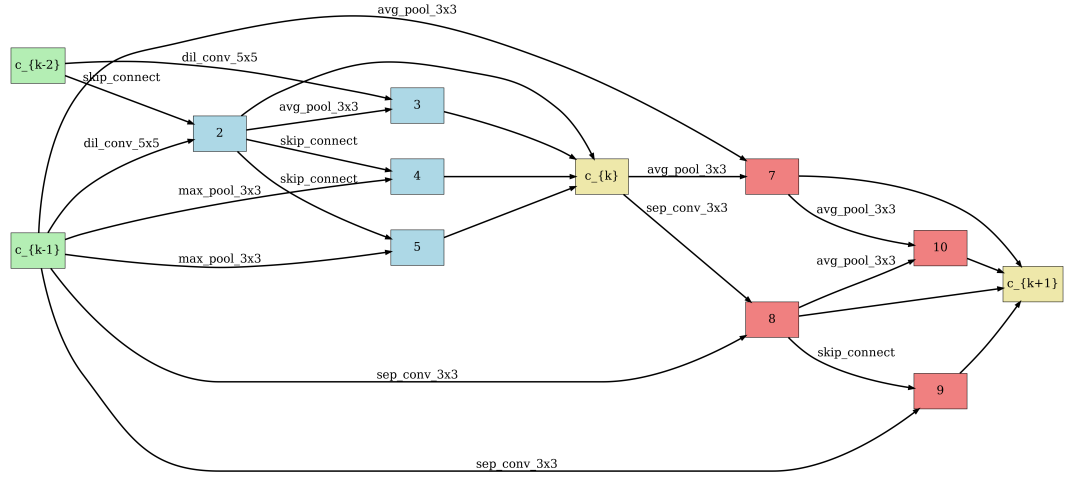


Fig. 4: Combined normal and reduced cells. The red vertices belong to the reduction cell; the blue vertices belong to the normal cell.

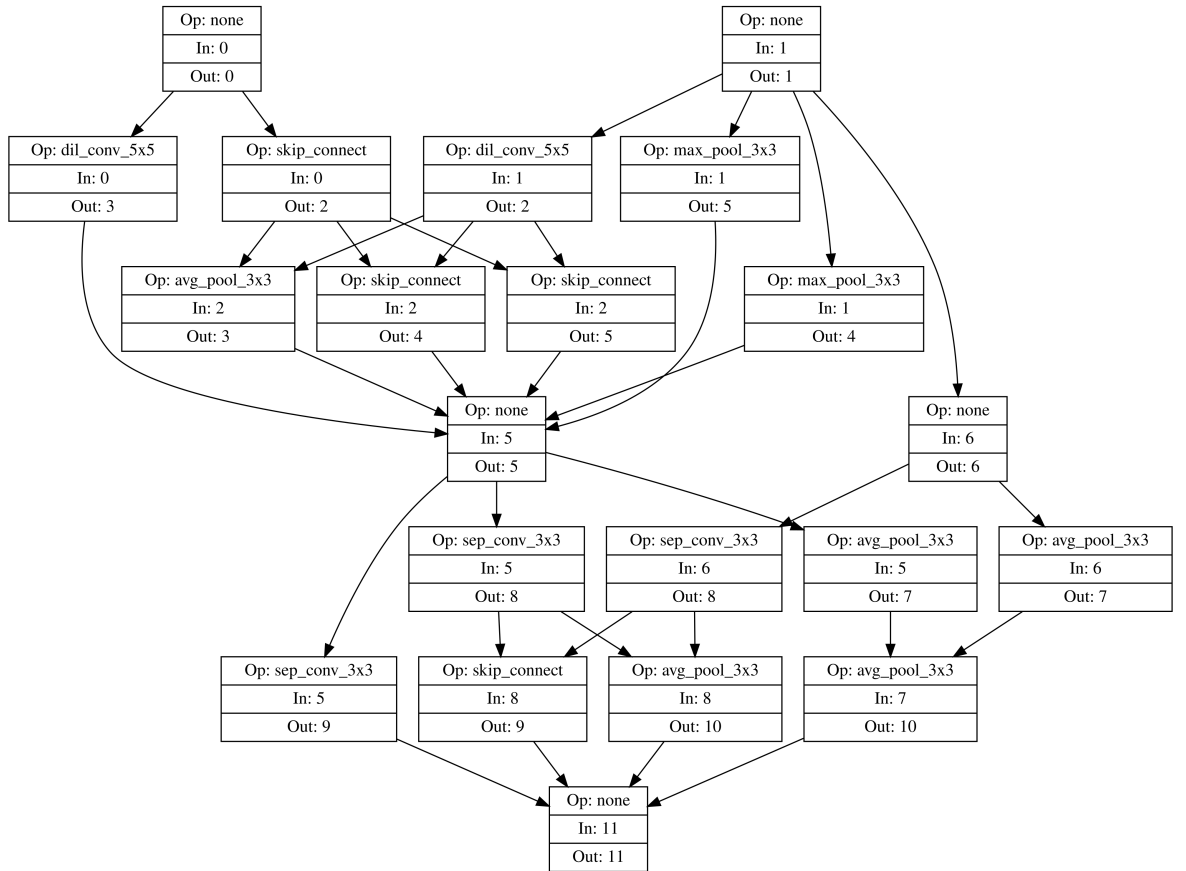


Fig. 5: Conversion of an architecture to the NAS-Bench-101 format.

B Distributions of models in dataset

The main distributions of trained models you can see on Figure 6 and Figure 7:

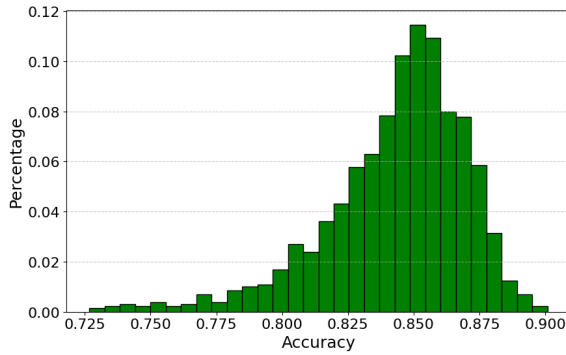


Fig. 6: Distribution of model accuracies

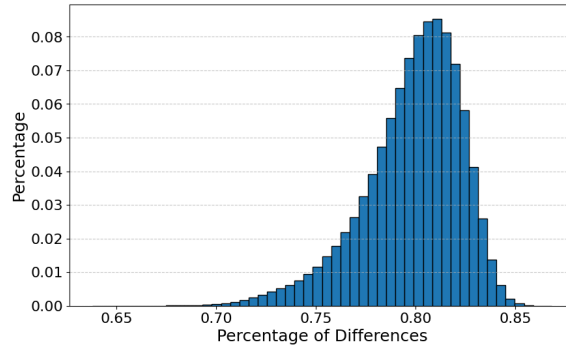


Fig. 7: Distribution of inter-model diversity

C Training curve for surrogate fuctionons

You can see the process of training accuracy surrogate function on Figure 9 and similarity surrogate function on Figure 8

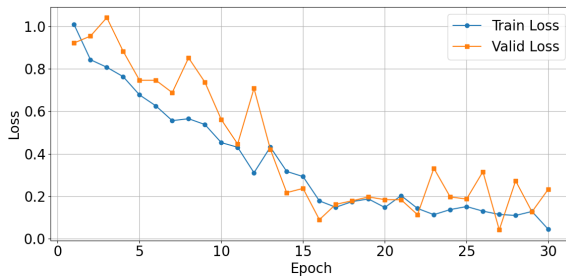


Fig. 8: Training curve of the surrogate model for diversity

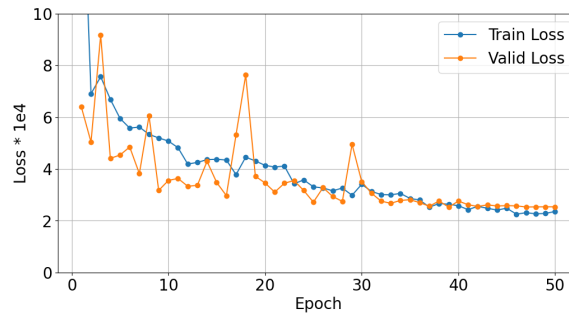


Fig. 9: Training curve of the surrogate model for accuracy