

GRASP-GCN: Graph-Shape Prioritization for Neural Architecture Search under Distribution Shifts

Sofia Casarin¹, Oswald Lanz¹ Sergio Escalera^{2,3}

¹Free University of Bozen-Bolzano, Bolzano, Italy

²Computer Vision Center, Barcelona, Spain

³Universitat de Barcelona, Barcelona, Spain

{scasarin@unibz.it, lanz@inf.unibz.it, sergio@maia.ub.es

Abstract

Neural Architecture Search (NAS) methods have shown to output networks that largely outperform human-designed networks. However, conventional NAS methods have mostly tackled the single dataset scenario, incurring in a large computational cost as the procedure has to be run from scratch for every new dataset. In this work, we focus on predictor-based algorithms and propose a simple and efficient way of improving their prediction performance when dealing with data distribution shifts. We exploit the Kronecker-product on the randomly wired search-space and create a small NAS benchmark composed of networks trained over four different datasets. To improve the generalization abilities, we propose GRASP-GCN, a ranking Graph Convolutional Network that takes as additional input the shape of the layers of the neural networks. GRASP-GCN is trained with the not-at-convergence accuracies, and improves the state-of-the-art of 3.3 % for Cifar-10 and increasing moreover the generalization abilities under data distribution shift.

1. Introduction

Neural Architecture Search (NAS) has drawn large research attention due to its efficacy in automatically optimizing the architecture of Deep Neural Networks (DNNs), replacing the error-prone manual design which demands high expertise. As the NAS process can be very expensive many methods were proposed to save time or computation, following two main directions: i) reducing the time required to evaluate each searched architecture proposing a weight sharing mechanism ([2], [1], [14], [11], [19]), ii) using sample efficient algorithms so that only few architectures are evaluated ([22], [21], [9], [15]). Proxy task performance and Predictor-based algorithms follow the second approach. They estimate the performance of the DNN either as an approximation or a prediction based on lower fidelities, such as i) shorter training ([22], [21]), ii) training

on a subset of the data ([9]), iii) on lower-resolution images ([3]), or iv) with less filters per layer and less cells ([22], [15]). While these approximations reduce the computational cost, they also introduce a bias in the estimate as performance will typically be underestimated. Predictor-based algorithms follow the second approach, and train a proxy model that can infer the validation accuracy of DNNs directly from their network structure. During optimization, the proxy can be used to narrow down the number of architectures for which the true validation accuracy must be computed, which makes predictor-based algorithms sample efficient. Predictor-based algorithms have been proposed by [17]; [13] and [6]. Despite the success of these kinds of approaches, only few methods ([10], [7]) tackle the problem of sharing or re-using the predictor knowledge on different datasets. Most conventional NAS methods are indeed task-specific, requiring repeatedly training the model from scratch for each new dataset. Moreover, existing NAS benchmarks either i) provide architectures trained on a single dataset ([20]), ii) define a benchmark across different tasks but not datasets ([5]), iii) or do not provide the full training-log of architectures and define not-unique networks in their search space ([4]). This limit the possible studies that can be done on the datasets to correctly interpret possible predictor results. In this paper, we restrict the problem to predictor-based algorithms, which given a NAS benchmarks are extremely fast to train (~10 min on GeForce GTX 1080). We aim at answering this question: *can we re-use knowledge that a predictor has learned on one dataset and transfer it to get a more sample-efficient algorithm on another dataset?* We study the impact of distribution shifts on the predictor performance, by analysing the ranking of the architectures trained on commonly used datasets for the task of image classification, and propose simple yet effective solutions to address the shift. Specifically, we design a randomly wired search space, that quaintly exploits the Kronecker product to impose a Resnet-like structure and

create a dataset of 2000 architectures trained on Cifar-10, Cifar-100, Tiny-ImageNet and Fashion-MNIST datasets. We study the generalization abilities of the predictor when directly used on a new dataset without fine-tuning and propose to integrate the so called *vertex shapes* - the shapes each layer has given a different input size, and to adopt early stopping - training the predictor with not-converged accuracies. Our study shows which are the limitations of predictor based algorithms, and our simple approach improves of **3.7%** and **9.5%** (without and with distribution shift) with respect to the naïve approach. To summarize, our contributions are threefold:

- We propose a new way of defining search spaces, that exploit the generality of randomly wired spaces but samples neural network efficiently through the Kronecker product and a criterion based on a desired skeleton, to obtain specific categories of neural networks.
- We analyse the generalization capabilities of naïve predictor based algorithms on two different scenarios, involving different latent data but the same observed data and different latent data with different observed data.
- We propose GRASP-GCN, which integrates the shapes of the layers of the neural networks as input to predictor, and trains with the accuracies of non-specialized neural networks.

2. Related Works

Different techniques were proposed to mitigate the large computation burden of NAS. These approaches primarily target the acceleration of either the evaluation or search modules within the NAS framework. The former accelerates the evaluation of each DNN, the latter increases the sample efficiency so that fewer architectures need to be evaluated for discovering a good network. Our work falls under the second category, as the predictor can be utilized to sample architectures that most likely perform well on a given task.

2.1. Single dataset predictor-based algorithms

Many predictor-based methods, that set the baseline for following works, focus on a single dataset. [17] train a regressor model on a small built dataset and select the top-K predicted architectures to train them from scratch. The proposed approach leads to a more than $20\times$ sample efficient algorithm with respect to standard used Evolution ones. [16] used graph neural network-based accuracy predictors and an iterative approach to estimate the accuracy of models. [17] propose a Graph-based neural Architecture Encoding Scheme, *i.e.* GATES, to improve the generalization abilities of performance predictors by modeling the information flow of the actual data processing of the architecture as the attributes of the input nodes. Neural Architecture Optimization, shortly NAO, [12]) reframes the NAS

problem as a continuous optimization problem. Through the use of a predictor that takes as input the continuous encoded representation of a neural network, NAO performs gradient based optimization in the continuous space to find the embedding of a new architecture with potentially better accuracy. [6] propose an efficient hardware-aware NAS method enabled by an accurate performance predictor based on Graph Convolutional Network (GCN). The authors show that the sample efficiency of predictor based NAS can be improved by considering binary relations of models and an iterative data selection strategy. Similarly to BRP-NAS, we employ a binary ranking GCN, but we extend the focus on multiple datasets and employ as labels the validation accuracy of non-specialized networks to improve the generalization abilities.

2.2. Transferable predictor-based algorithms

Among existing methods, relevant approaches to ours tackling the generalization problem across multiple dataset are MetaD2A by [10], and Arch-Graph by [7]. MetaD2A stochastically generates graphs from a dataset via a cross-modal latent space that is learned via amortized meta-learning. From the encoding of the dataset, obtained through a permutation invariant encoder set, a graph is decoded. A meta-predictor is then used to estimate and select the best architecture for a given dataset. Instead of using an encoder set, we propose a much simpler and more general solution that does not limit the approach to image-dataset due to the encoder-set, lacking the possibility to adapt it to video. We provide as additional input to the GCN the “vertex shapes”, which are strictly related to the shapes of the data. In Arch-Graph the generalization problem is addressed from the point of view of task generalization, rather than dataset generalization. The method predicts task-specific optimal architectures with respect to given task embeddings, by leveraging correlations across multiple tasks through their embeddings as a part of the predictor’s input for fast adaptation. Despite being sample efficient across many tasks, the method requires predictor-tuning on the new task/dataset. With respect to previous approach, our work tackles the problem from the point of view of the distribution shift in the dataset, proposing a simpler yet more general method that takes into consideration the specialization (or overfitting) of neural networks over datasets, the shape characteristics of the data and their effect on networks, and that does not require fine-tuning the predictor.

3. Methods

Our goal is to obtain a predictor that generalizes well across different datasets without need to re-train or finetune. To this end, we sample architectures from our search-space and train them over 4 image classification datasets (3.1) and pro-

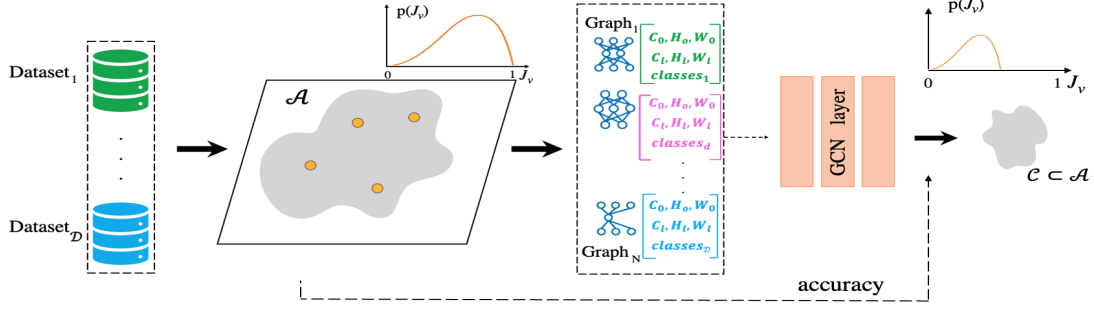


Figure 1. Architectures are sampled from the search space and trained over 4 datasets. The structure of DNNs is given as input with the shapes of the layers to a ranking GCN, which given the accuracy learns to rank DNNs so that the search space is narrowed down.

pose GRASP-GCN (3.2), which trains a ranking predictor with an additional input consisting in the vertex shapes.

3.1. Search-space definition

NAS is formalized as a bi-level optimization problem:

$$\begin{aligned} \mathcal{A}^* &= \operatorname{argmin}_{\mathcal{A}} J_v(\mathcal{A}, \mathbf{w}^*) \\ \text{s.t. } \mathbf{w}^* &= \operatorname{argmin}_{\mathbf{w}} J_t(\mathcal{A}, \mathbf{w}), \end{aligned} \quad (1)$$

where \mathcal{A} describes the architecture, \mathbf{w} are the trainable weights of the considered DNN, and J_v and J_t are the validation and training loss, respectively. In our work we define $\mathcal{A} = (\mathbf{A}, \mathbf{X})$, with $\mathbf{A}^{N \times N}$ and $\mathbf{X}^{N \times D}$ encoding the connections in the graph and the types of layers, respectively. The dimensions of the matrices are related to the N number of nodes (layers) and the D number of input features, *i.e.* the number of layer types allowed. We sample \mathcal{A} from our search space composed of all feed-forward convolutional networks, belonging to the randomly-wired search space [18] and exploit the Kronecker product, a trick that allows us balancing flexibility and efficiency, avoiding the sampling of big random matrices. We generate resnet-like architectures miming evolutionary sampling. Two random matrices $\mathbf{R}_1^{4 \times 4}$ and $\mathbf{R}_2^{4 \times 4}$ are indeed sampled from the search space and multiplied with two so-called skeleton matrices $\mathbf{K}_1^{4 \times 4}$ and $\mathbf{K}_2^{4 \times 4}$:

$$\begin{aligned} \mathbf{A} &= \underbrace{\begin{bmatrix} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \end{bmatrix}}_{\mathbf{K}_1} \otimes \begin{bmatrix} R_1 \\ R_1 \\ R_1 \\ R_1 \end{bmatrix} \\ &+ \underbrace{\begin{bmatrix} 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \\ 0 & 0 & 0 & 1 \\ 0 & 0 & 0 & 0 \end{bmatrix}}_{\mathbf{K}_2} \otimes \begin{bmatrix} R_2 \\ R_2 \\ R_2 \\ R_2 \end{bmatrix} \\ &= \begin{bmatrix} R_1 & R_2 & & \\ & R_1 & R_2 & \\ & & R_1 & R_2 \\ & & & R_1 \end{bmatrix} \end{aligned} \quad (2)$$

In Eq. 2, \mathbf{K}_1 which is multiplied by \mathbf{R}_1 defines the feed-forward structure, while \mathbf{K}_2 , with the off-diagonal values,

defines the shortcuts. The *input* and the *output* layers are finally added to the generated matrix $\mathbf{A}^{16 \times 16}$, leading to a maximum dimension of 18×18 . As shown, the Kronecker product \otimes allows repeating \mathbf{R}_1 and \mathbf{R}_2 structures, limiting the randomly wired search space in a meaningful way. Moreover, it generates easily scalable networks, a key advantage as proven by [22], by stacking multiple blocks. This design choice was intentional to focus on analyzing the true impact of the dataset itself. Moreover, it does not represent a limitation as the only requirement to properly train a GCN predictor is to have both very-well and very badly-performing DNNs. Table 1 summarized our candidate layers which constitute our $\mathbf{X}^{18 \times 9}$.

3.2. Ranking GCN

Graph Convolutional Networks are DNN architectures that extract multi-scaled localized spatial features to extract highly expressive representations of graphs, dealing with the difficulty of “localized convolution” filters in non-Euclidean domains. GCNs performs a convolution looking for essential vertices and edges with the goal of learning the features of the graph. It takes as input: (i) a feature description X_i for every node i summarized in a feature matrix $\mathbf{X}^{N \times D}$ where N is the number of nodes, D the number of input features; (ii) a representative description of the graph structure summarized in the adjacency matrix $\mathbf{A}^{N \times N}$. For the classification task, the GCN produces a node-level output $\mathbf{H}^{N \times F}$, where F is the number of output features per node. The GCN outputs $h(\mathbf{A}, \mathbf{X})$ which is the concatenation of each \mathbf{H}^l layer mapping done, as described in Eq. 3:

$$h(\mathbf{A}, \mathbf{X}) = (H^L \circ H^{L-1} \circ \dots \circ H^l \circ \dots \circ H^1)(\mathbf{A}, \mathbf{X}) \quad (3)$$

where $l = 1, 2, \dots, L$, is the number of layers in the GCN, and \mathbf{H}^l is given by the propagation rule (eq. 4) defined by in [8]:

$$H^{(l+1)} = \sigma(\tilde{D}^{-\frac{1}{2}} \tilde{A} \tilde{D}^{-\frac{1}{2}} H^{(l)} W^{(l)}) \quad (4)$$

with $\tilde{A} = A + I$, (I identity matrix), $\tilde{D}_{ii} = \sum_j \tilde{A}_{ij}$ is the degree matrix, $W^{(l)}$ is the layer-specific trainable weight

Layer	stem conv3×3	conv3×3	conv3×3×d	conv3×3×h	conv3×3s2	conv3×3s2×d	conv3×3s2×h
Channels	64	Same	Doubles	Halves	Same	Doubles	Halves
Stride	-	-	-	-	2	2	2

Table 1. Layers in the search space. Same / Double / Halves refer to the channels of the parent nodes.

matrix, σ is the ReLU activation function, $H^{(l)} \in \mathbb{R}_{N \times D}$ is the matrix of activations in the $l^{(th)}$ layer, and $H^{(l+1)}$ is the output to the next layer. Graph-level outputs can then be modeled by introducing some form of pooling operation. In our work we exploit a GCN with 3 layers and a classification head, and following BRP-NAS we exploit a Ranking GCN, which learns a binary relation that focuses on the prediction of ranking. Indeed, as previously observed by [6] *i*) accuracy prediction is not necessarily required to produce faithful estimates (in the absolute sense) as long as the predicted accuracy preserves the ranking of the models; *ii*) any antisymmetric, transitive and convex binary relation produces a linear ordering of its domain, implying that NAS could be solved by learning binary relations, where $O(n^2)$ training samples can be used from \underline{n} measurements. Given the architecture’s search space, a ranking network predicts *how likely* any network in the search space reaches a higher accuracy than the current best. Fig. 2 shows how a GCN can be used as a ranking network: two architectures \mathcal{A}_1 and \mathcal{A}_2 are fed to the GCN. The GCN outputs two graph encodings h_1 and h_2 , which are then concatenated into the vector \mathbf{h} . A softmax function $\sigma(\cdot)$ is applied, obtaining the ranking probabilities which are then compared with the target \mathbf{t} . Given for example the tuple $(acc_{i,\mathcal{A}_j}, acc_{i,\mathcal{A}_k})$ where acc_{i,\mathcal{A}_j} is the accuracy of architecture \mathcal{A}_j and acc_{i,\mathcal{A}_k} is the accuracy of \mathcal{A}_k , both over dataset i , if $acc_{i,\mathcal{A}_j} > acc_{i,\mathcal{A}_k}$, then the target vector takes as values $\mathbf{t} = [1 \ 0]$. Our goal is to maximize the log probability that \mathcal{A}_j is better than \mathcal{A}_k , and therefore the predictor is trained with the loss in Eq. 5:

$$J = -\mathbf{t}^\top \cdot \ln(\sigma(\mathbf{y})). \quad (5)$$

As shown in Fig. 2, $\mathcal{A}_{j,k}$, which gives the information about the structure of the graphs, and the feature vector \mathcal{X} , which defines the layers in the DNN architecture. Therefore, given a DNN architecture that is characterized by a bottleneck, the predictor won’t capture the implications of a severe reduction of the feature maps. Clearly, the same DNN architecture would have very different performance depending on the input shape and the predictor could not see the difference if given with \mathcal{A} and \mathcal{X} only. We therefore add the dimensions characterizing each DNN layer. This info shortly called “*vertex shapes*” and is concatenated to \mathcal{X} that has now dimensions 18×12 . We normalize the vertex shapes with respect to the maximum dimensions in our dataset and encode them as float numbers. We did not choose one hot encoding as it implicitly loses any ordering and distance knowledge, *i.e.* we don’t know if shape (3,32,32) is closer to (3,64,64) or to (4096,3,3).

4. Implementation details

Our hyper-parameters are summarized in Tab. 2 which displays those used for training the architectures sampled from our search-space and in Tab 3 for training our GRASP-GCN. All hyper-parameters were optimized using Optuna framework¹. We moreover provide the pseudo-code for the creation of our NAS dataset (Algorithm 1), and for the functioning of predictor-based algorithms within the NAS framework (Algorithm 2).

Algorithm 1 Dataset creation

```

1: if Cifar10 then
2:   for iteration = 1, ..., 2000 do
3:     sample  $\mathcal{A} = \mathbf{A}, \mathbf{X}$ 
4:     train  $\mathcal{A}$ 
5:     save values in hash( $\mathcal{A}$ )
6:   end for
7: end if
8: if FashionMNIST or Cifar100 or TinyImageNET then
9:   for iteration = 1, ..., 2000 do
10:    random sample  $\mathcal{A}$  from Cifar10 dataframe
11:    new directory = hash( $\mathcal{A}$ )
12:    while new directory exist do
13:      random sample  $\mathcal{A}$  from Cifar10 dataframe
14:      new directory = hash( $\mathcal{A}$ )
15:    end while
16:    train  $\mathcal{A}$ 
17:    save values in new directory
18:   end for
19: end if

```

5. Experiments

In this section, we introduce the used ranking measures (Sec. 6), we investigate the statistics of our NAS dataset composed of (architectures, accuracy) pairs to provide ground-truth values for the predictor (Sec. 6.1). We show the performance of our predictor under distribution shift (Sec. 6.2) and further validate the usage of the vertex shapes as additional input and of not converged accuracies through ablations (Sec. 6.3).

6. Ranking measures

Several measures are commonly used to assess the quality of a ranking. We focus on: 1) NDCG@k, 2) Precision@k and 3) Kendall’s τ .

¹<https://optuna.readthedocs.io/en/stable/>

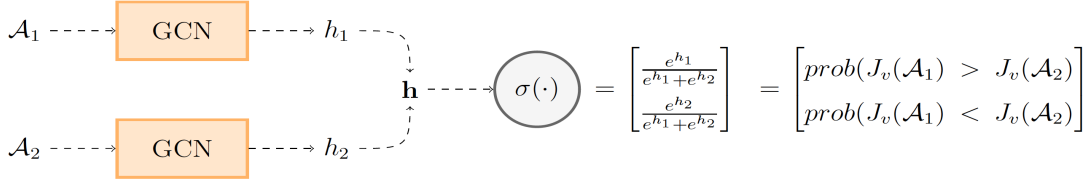


Figure 2. Working principles of a GCN used as a ranking network

Dataset	Learning rate	Weight decay	Drop lr	Optimizer	Batch-size
F-MNIST	0.100	0.0002	of 10 at epoch 40, 80	SGD	128
C10	0.097	0.0006	of 10 at epoch 40, 80	SGD	128
C100	0.065	0.0015	of 5 at epoch 40, 80, 100	SGD	128
Tiny	0.012	0.0011	of 10 at epoch 30, 60, 90	SGD	64

Table 2. List of hyper-parameters derived from Optuna optimization for ResNet-18. The columns show the learning rate (lr), the weight decay (wd), the drop of the learning rate (drop lr), the optimizer, and the size of the batches.

Units per layer	Lr	Weight decay	Optimizer
265	0.019041	0.001126	Adagrad

Table 3. List of hyper-parameters derived from Optuna optimization for GRASP_GCN.

Algorithm 2 Predictor-based neural architecture search

- 1: \mathcal{S} = Architecture search space
- 2: $f(\mathcal{A}, \theta) : \mathcal{A} \rightarrow \mathbb{R}$: Predictor that outputs the predicted performance given the architecture
- 3: $N^{(k)}$: Number of architectures to sample in the k -th iteration
- 4: $k = 1$
- 5: $\tilde{\mathcal{S}} = \emptyset$
- 6: **while** $k \leq \text{MAX_ITER}$ **do**
- 7: Sample a subset of architectures $\mathcal{C}^{(k)} = \{\mathcal{A}_j^{(k)}\}_{j=1, \dots, N^{(k)}}$ from \mathcal{S} utilizing $f(\mathcal{A}, \theta)$
- 8: Evaluate architectures in $\mathcal{S}^{(k)}$, get $\tilde{\mathcal{C}}^{(k)} = \{\mathcal{A}_j^{(k)}, y_j^{(k)}\}_{j=1, \dots, N^{(k)}}$ (y is the performance)
- 9: $\mathcal{S} = \mathcal{S} - \mathcal{C}$
- 10: $\tilde{\mathcal{S}} = \tilde{\mathcal{S}} \cup \tilde{\mathcal{C}}$
- 11: Optimizing $f(\mathcal{A}, \theta)$ using the ground-truth architecture evaluation data $\tilde{\mathcal{S}}$
- 12: **end while**
- 13: Output $\mathcal{A}_{j^*} \in \tilde{\mathcal{S}}$ with best corresponding y^* ; Or, $\mathcal{A}^* = \text{argmax}_{\tilde{\mathcal{S}}} f(\mathcal{A}, \theta)$

NDCG DCG (Discounted Cumulative Gain) is founded on the idea that when assessing search results, highly relevant documents ranked lower should receive a larger penalty than less relevant documents wrongly ranked. This is because the graded relevance value decreases logarithmically as the position of the result worsens. This measure applies very well to predictor-based algorithms since our focus is on allowing the predictor to find the best-forming

architectures in such a way that they are placed on the top of the list, while we don't care about how bad networks are ranked. Eq. 6 exactly shows how highly relevant objects that appear lower in the ranked list are penalized by reducing the graded relevance value logarithmically proportional to the position of the result:

$$DCG_k = \sum_{i=1}^k \frac{2^{rel_i} - 1}{\log_2(i + 1)}. \quad (6)$$

Starting from 6, the NDCG can be easily obtained normalizing w.r.t the Ideal Discounted Cumulative Gain (IDCG), as shown in Eq. 7:

$$NDCG_k = \frac{DCG_k}{IDCG_k} \in [0, 1], \quad (7)$$

$$IDCG_k = \sum_{i=1}^{|REL_k|} \frac{2^{rel_i} - 1}{\log_2(i + 1)}$$


where rel_i is the relevance value assigned to each object and REL_k is the list of relevant objects ordered by their relevance up to position k . In a perfect ranking algorithm, the DCG_p will be the same as the $IDCG_p$ producing an NDCG of 1. We used the NDCG in two different variants: the NDCG@2092 was used to get a big picture of the general behavior all architectures have, while the NDCG@10 was used to focus on the ranking quality of the top-10 performing architectures, which are those we are interested in.

Precision@k It is defined as the proportion of recommended items in the top-k set that are relevant. It is mathematically defined as:

$$Precision@k = \frac{\# \text{of items @k that are relevant}}{k} \in [0, 1], \quad (8)$$

where again, as it can be observed, the concept of relevance is involved.

	F-mnist	C10	C100	Tiny
F-mnist	0	0.0273	0.0461	0.0516
C10	0.1781	0	0.1032	0.1293
C100	0.4559	0.2653	0	0.1743
Tiny	0.5157	0.3456	0.1817	0

Table 4. 1-NDCG@2092. The columns display the training dataset, rows the validation dataset. Color bar:  (lowest to highest).

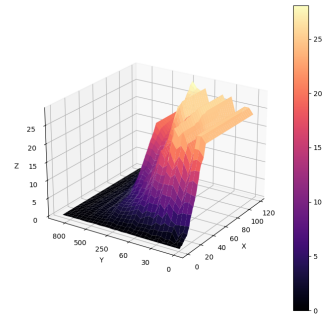
Kendall’s τ The Kendall rank correlation coefficient is used to measure the ordinal association between two measured quantities. It’s range is in $[-1, 1]$, with Kendall’s $\tau = 0$ that indicates absence of correlation. Let (x_1, x_2, \dots, x_n) and (y_1, y_2, \dots, y_n) be a set of observations of the joint variables X and Y , such that all the values of (x_i) and (y_i) are unique (ties are neglected for simplicity). Any pair of observations (x_i, y_i) and (x_j, y_j) , where $i < j$ are said to be concordant if the sort of order (x_i, x_j) and (y_i, y_j) agrees: that is, if either both $x_i > x_j$ and $y_i > y_j$ holds if both $x_i < x_j$ and $y_i < y_j$; otherwise, they are said to be discordant. Eq. 9 defines, based on this, the coefficient:

$$\tau = \frac{(\text{\#of concordant pairs}) - (\text{\#of discordant pairs})}{\binom{n}{2}}, \tag{9}$$

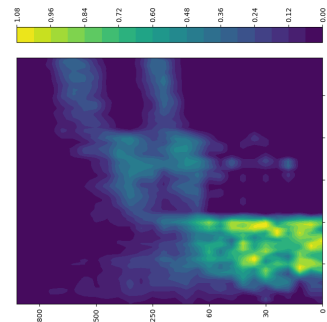
6.1. Dataset

The four datasets were chosen based on the possibility to exploit a variety of conditions to investigate the transferability of the predictor, and by the desire to extend and be partially comparable with previous works that involve three datasets, as in NATS-Bench [4]. Two different scenarios were defined: i) transferability when different latent data is involved *but* the observed data is **the same**, ii) transferability when different latent data is involved *and* the observed data is **not the same**. Hence, we chose (a) Cifar-10 as a baseline dataset, (b) Cifar-100 as it is composed of Cifar-10 images, labeled differently, (c) Fashion-MIST as it is composed of black and white images (different latent data), belonging to completely different categories with respect to Cifar-10, and (d) Tiny-ImageNET to consider a more complex dataset with a larger number of classes. 2000 *unique* architectures were sampled from the search space and trained over the aforementioned datasets.

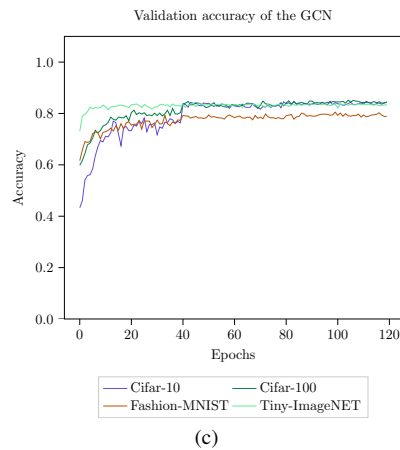
Models trained over different input distributions are not guaranteed to perform in the same way. Do they “overfit” the dataset specializing during training? Can early stopping be applied to reduce the training time it takes to get the true validation accuracy of the architectures?



(a)



(b)



(c)

Figure 3. Cifar-10 results. Ranking evolution during training with a cumulative (a) and derivative (b) plot. The x-axis shows the architectures in a descent order with respect to their accuracy. The y-axis carries the training epochs. The heatmap displays large number of rank changes (yellow) to no changes (blue). (c) Validation accuracy the GCN can obtain when trained with the previous rankings.

Distribution Shift Tab. 4, to be read column-wise, highlights that the order induced by complex datasets generalizes better for simpler datasets than vice-versa. This result is complemented by Fig. 4, that displays the ranking stability of each architecture on one of the four datasets. The plot

	F-mnist	C10	C100	Tiny
F-mnist	85.1 ± 0.7	76.4 ± 0.5	68.2 ± 0.6	68.4 ± 0.8
C10	74.2 ± 0.6	87.9 ± 0.4	84.9 ± 0.4	83.4 ± 0.7
C100	70.8 ± 0.6	85.0 ± 0.1	87.0 ± 0.4	85.6 ± 0.3
Tiny	72.3 ± 0.3	83.6 ± 0.9	85.3 ± 0.2	87.9 ± 0.2

(a) Average validation accuracy (%) of the predictor.

	F-mnist	C10	C100	Tiny
F-mnist	100 ± 0	93 ± 2	94 ± 2	96 ± 1
C10	80 ± 4	99 ± 1	97 ± 2	95 ± 3
C100	36 ± 3	60 ± 2	65 ± 1	63 ± 1
Tiny	41 ± 2	61 ± 3	57 ± 3	69 ± 2

(b) Average Precision@10 (×100) of the predictor.

Table 5. Performance of our GRASP-GCN provided with input matrix \mathbf{A} , features \mathbf{X} , vertex shapes \mathbf{V} and trained with the validation accuracy at epoch 40. For each element of the table, the mean and the std of four runs are given. The columns are the training datasets, the rows are the validation datasets.

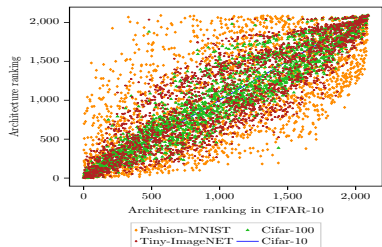


Figure 4. The ranking of each architecture on one of the four datasets, sorted by the ranking in Cifar-10. Correlation measured through Kendall’s τ .

was obtained by sorting the architectures trained on Cifar-10 in descending order with respect to their accuracy on the dataset, and inducing the same ID-order on the architectures trained on different datasets. In this way, it is possible to observe what rank has been assigned to the architectures trained over different datasets, having as baseline Cifar-10 rank. We can observe that (i) Fashion-MNIST architectures cause a higher variability in terms of architectures rank ii) the plot start tighter, increase variance as we move towards bigger ranks, and gets tighter as we approach the “worst” architectures. We can deduce that, not only the worst architectures are such since the very beginning of training, (Fig. 3a, 3b), but as it could be naturally expected, some architectures simply do not have enough capacity to solve any classification task and perform badly independently on the input distribution they are provided, further validating the early stopping method we propose in Sec. 6.2.

Network Specialization We compared through NDCG the ranking induced by the validation accuracy v_acc_i with $i = 1, \dots, 120$ with v_acc_{120} at the end of training. Fig. 5 highlights a possible correlation between the change of rank and the epoch where the learning rate is dropped (we refer to Appendix 4 for details on hyper-parameters). Moreover, focusing on the NDCG@10, that considers only the top-10, the plot reaches zero way before the end of the training. Therefore, we formulated the hypothesis that the difference between NDCG2092 and the NDCG10 could be caused by the average performing networks, which are more strongly influenced by the hyper-parameters. We checked the evolution of ranking during training (3a, 3b) by counting how many times the architectures change the relevance value (Eq. 7). We observe that top-performing DNNs stop chang-

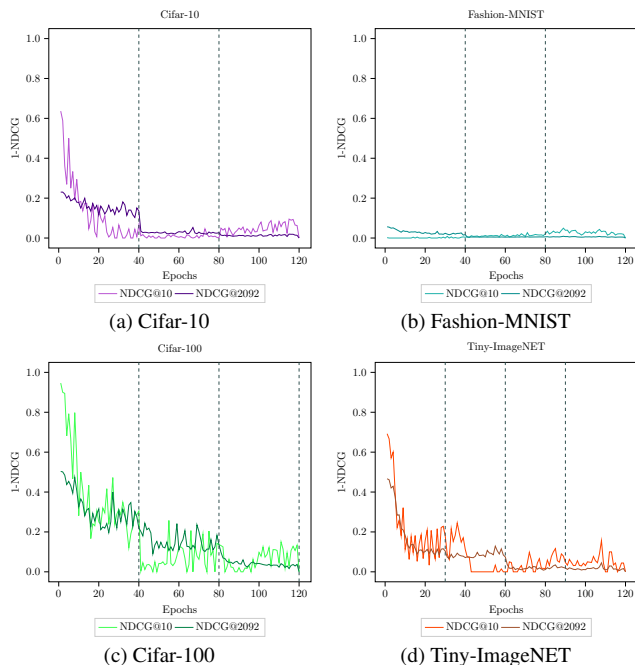


Figure 5. 1-NDCG plot showing the ranking correlation among the sorting the architectures have at epoch i w.r.t the sorting of the architectures at epoch 120. The lower the better. Every figure displays the results for considering one dataset at a time. The dashed lines over each plot highlight the epoch where the learning rate is dropped. Every plot displays both the NDCG@10 (light color) and the NDCG@2092 (dark color).

ing rank at earlier epochs, average performing ones have the peak shifted towards the end of training, and DNNs not able to solve the task are such since the early epochs.

6.2. GRASP-GCN

We use the validation accuracy, and the Precision@10 measures to evaluate the performance of GRASP-GCN. Tab 5a, Tab. 5b show the results obtained for each training dataset (columns) on each of the four validation datasets (rows), while Fig. 3c gives a compact representation of the GCN performance trained with the accuracy of every training epoch. The values in the table are above 80 % when Fashion-MNIST is not involved, which suggests that the predictor trained over the datasets involved in the lower part of the table, e.g. Cifar-10, Cifar-100, Tiny- ImageNET, can

	F-mnist			C10			C100			Tiny		
F-mnist	85.1	79.1	83.8	76.4	67.7	74.7	68.2	65.8	69.2	68.4	61.7	65.9
C10	74.2	71.9	73.4	87.9	84.2	86.6	84.9	82.3	84.0	83.4	80.1	83.4
C100	70.8	68.9	71.2	85.0	82.5	84.3	87.0	84.4	86.5	85.6	82.7	84.8
Tiny	72.3	62.5	70.9	83.6	81.1	83.6	85.3	83.0	85.0	87.9	84.0	87.3

(a) Accuracy (%).

	F-mnist			C10			C100			Tiny		
F-mnist	1.00	1.00	1.00	0.93	0.90	0.93	0.94	0.88	0.95	0.96	0.90	0.95
C10	0.80	0.78	0.82	0.99	0.98	0.98	0.97	0.80	0.95	0.95	0.90	0.91
C100	0.36	0.35	0.41	0.60	0.45	0.41	0.65	0.45	0.50	0.63	0.40	0.60
Tiny	0.41	0.10	0.50	0.61	0.40	0.58	0.57	0.43	0.55	0.69	0.50	0.64

(b) Average Precision@10.

Table 6. Performance measures of the predictor with vertex shapes and early stopping (black values), without vertex shapes (purple), without early stopping (blue). The datasets used for training are displayed in the columns, while the rows show the validation ones. Best results are obtained when both vertex shapes and early stopping techniques are employed.

transfer knowledge over one of these same datasets. On the other hand, when Fashion-MNIST is involved, either as the training or the validation dataset, the performance drops drastically compared to the average performance the predictor has with Cifar-10/Cifar-100/Tiny-ImageNet. It is worth noting however, that despite the drop in performance, when Fashion-MNIST is the validation dataset (first row) the drop does not represent a problem. Indeed as the Precision@10 highlights, the top-performing architecture are still correctly ranking. More precisely, as almost all architectures well-solve Fashion-mnist, a wrong ranking will probably affect architecture with a small difference in validation accuracy, thus resulting in a good ranking order. Finally, Tab. 7 shows that GRASP-GCN surpasses all other methods when trained and evaluated over Cifar-10 and when directly applied to new datasets.

	F-mnist	C10	C100	Tiny
F-mnist	84.1 82.2 79.8	74.4 71.1 68.2	68.2 64.8 68.2	68.4 63.5 65.5
C10	74.2 73.8 73.2	87.9 84.1 86.2	84.9 81.1 80.8	83.4 81.1 80.7
C100	70.8 68.5 70.7	85.0 80.1 79.1	87.0 82.1 84.3	85.6 82.4 85.0
Tiny	72.3 66.4 70.1	83.6 79.2 83.2	85.3 80.1 83.4	87.9 84.0 86.6

Table 7. Comparison between (black) ours, (blue) BRP-NAS, (red), MetaD2A.

6.3. Ablation

Tab. 6a highlights how the predictor significantly improves the performance when trained and validated over the same dataset (diagonal values) improving the baseline of more than 3%. If we focus on transferability, looking out of the diagonal, we have even larger improvements, with gaps exceeding 9%. Consistent results are obtained in 6b with Precision@10 measure. Finally, Fig. 3c ablates on the early-stopping mechanism we propose to employ in the NAS procedure. We can observe that if we use as training set the accuracies the architectures have after the first drop of the learning rate (epoch 40) the performance of the predictor

is not affected. This is further validated by Tab. 6a, which compares our best results with and without early stopping.

7. Conclusions

In our work, we face the problem of analyzing the transferability of a predictor under data distribution shift. For this reason, we created our small dataset composed of architectures trained on four datasets. Our ranking analysis on the trained networks showed an association between the drop of the learning rate and the epoch where architectures stop changing their ranking during training, highlighting that top-performing architectures keep their ranking since early training epochs, while the worst ones are such since the early epochs. We moreover spotted that the ranking induced by complex datasets generalizes better than those induced by simple datasets. Given this, we improve the naive predictor training by including the vertex shapes as input, and employing an early stopping procedure. Our method surpasses state-of-the-art performance on predictor-based algorithms addressed by distribution shifts. We believe that such a result can help during the search phase of NAS, as to get the true validation accuracy of the top-k ranked architectures to pick the best-performing ones the training can be stopped early reducing significantly the searching time.

Limitations and future works Possible future works could include enlarging the datasets by considering new tasks such as object detection and new modalities such as videos. Another interesting direction could be studying how the method improves with the inclusion of gradients information as input to the predictor.

Acknowledgement

This work has been partially supported by the Spanish project PID2022-136436NB-I00, by ICREA under the ICREA Academia programme, and by unibz startup fund IN 2814 and IN 2902.

References

- [1] Gabriel Bender, Pieter-Jan Kindermans, Barret Zoph, Vijay Vasudevan, and Quoc Le. Understanding and simplifying one-shot architecture search, 2018. [1](#)
- [2] Han Cai, Chuang Gan, Tianzhe Wang, Zhekai Zhang, and Song Han. Once-for-all: Train one network and specialize it for efficient deployment, 2020. [1](#)
- [3] Patryk Chrabaszcz, Ilya Loshchilov, and Frank Hutter. A downsampled variant of imagenet as an alternative to the cifar datasets, 2017. [1](#)
- [4] Xuanyi Dong, Lu Liu, Katarzyna Musial, and Bogdan Gabrys. Nats-bench: Benchmarking nas algorithms for architecture topology and size, 2020. [1](#), [6](#)
- [5] Yawen Duan, Xin Chen, Hang Xu, Zewei Chen, Xiaodan Liang, Tong Zhang, and Zhenguo Li. Transnas-bench-101: Improving transferability and generalizability of cross-task neural architecture search. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages 5251–5260, 2021. [1](#)
- [6] Łukasz Dudziak, Thomas Chau, Mohamed Abdelfattah, Royson Lee, Hyeji Kim, and Nicholas Lane. Brp-nas: Prediction-based nas using gcns, 2020. [1](#), [2](#), [4](#)
- [7] Minbin Huang, Zhijian Huang, Changlin Li, Xin Chen, Hang Xu, Zhenguo Li, and Xiaodan Liang. Arch-graph: Acyclic architecture relation predictor for task-transferable neural architecture search, 2022. [1](#), [2](#)
- [8] Thomas N. Kipf and Max Welling. Semi-supervised classification with graph convolutional networks, 2017. [3](#)
- [9] Aaron Klein, Stefan Falkner, Simon Bartels, Philipp Hennig, and Frank Hutter. Fast bayesian optimization of machine learning hyperparameters on large datasets. ArXiv, abs/1605.07079, 2016. [1](#)
- [10] Hayeon Lee, Eunyoung Hyung, and Sung Ju Hwang. Rapid neural architecture search by learning to generate graphs from datasets, 2021. [1](#), [2](#)
- [11] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: Differentiable architecture search. In International Conference on Learning Representations, 2019. [1](#)
- [12] Renqian Luo, Fei Tian, Tao Qin, Enhong Chen, and Tie-Yan Liu. Neural architecture optimization, 2019. [2](#)
- [13] Xuefei Ning, Yin Zheng, Tianchen Zhao, Yu Wang, and Huazhong Yang. A generic graph-based neural architecture encoding scheme for predictor-based nas, 2020. [1](#)
- [14] Hieu Pham, Melody Guan, Barret Zoph, Quoc Le, and Jeff Dean. Efficient neural architecture search via parameters sharing. In Proceedings of the 35th International Conference on Machine Learning, pages 4095–4104. PMLR, 2018. [1](#)
- [15] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc Le. Regularized evolution for image classifier architecture search. Proceedings of the AAAI Conference on Artificial Intelligence, 33, 2019. [1](#)
- [16] Chen Wei, Chuang Niu, Yiping Tang, Yue Wang, Haihong Hu, and Jimin Liang. Npenas: Neural predictor guided evolution for neural architecture search. IEEE Transactions on Neural Networks and Learning Systems, 34(11):8441–8455, 2023. [2](#)
- [17] Wei Wen, Hanxiao Liu, Hai Li, Yiran Chen, Gabriel Bender, and Pieter-Jan Kindermans. Neural predictor for neural architecture search, 2019. [1](#), [2](#)
- [18] Saining Xie, Alexander Kirillov, Ross Girshick, and Kaiming He. Exploring randomly wired neural networks for image recognition, 2019. [3](#)
- [19] Sirui Xie, Hehui Zheng, Chunxiao Liu, and Liang Lin. SNAS: stochastic neural architecture search. In International Conference on Learning Representations, 2019. [1](#)
- [20] Chris Ying, Aaron Klein, Esteban Real, Eric Christiansen, Kevin P. Murphy, and Frank Hutter. Nas-bench-101: Towards reproducible neural architecture search. In International Conference on Machine Learning, 2019. [1](#)
- [21] Arber Zela, Aaron Klein, Stefan Falkner, and Frank Hutter. Towards automated deep learning: Efficient joint neural architecture and hyperparameter search. CoRR, abs/1807.06906, 2018. [1](#)
- [22] B. Zoph, V. Vasudevan, J. Shlens, and Q. V. Le. Learning transferable architectures for scalable image recognition. In 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR), pages 8697–8710, Los Alamitos, CA, USA, 2018. IEEE Computer Society. [1](#), [3](#)