

# Chapter 3

## Neural Architecture Search

Thomas Elsken and Jan Hendrik Metzen and Frank Hutter

### Abstract

Deep Learning has enabled remarkable progress over the last years on a variety of tasks, such as image recognition, speech recognition, and machine translation. One crucial aspect for this progress are novel neural architectures. Currently employed architectures have mostly been developed manually by human experts, which is a time-consuming and error-prone process. Because of this, there is growing interest in automated *neural architecture search* methods. We provide an overview of existing work in this field of research and categorize them according to three dimensions: search space, search strategy, and performance estimation strategy.

### 3.1 Introduction

The success of deep learning in perceptual tasks is largely due to its automation of the feature engineering process: hierarchical feature extractors are learned in an end-to-end fashion from data rather than manually designed. This success has been accompanied, however, by a rising demand for *architecture engineering*, where increasingly more complex neural architectures are designed manually. *Neural Architecture Search* (NAS), the process of automating architecture engineering, is thus a logical next step in automating machine learning. NAS can be seen as subfield of AutoML and has significant overlap with hyperparameter optimization and meta-learning (which are described in Chapters 1 and 2 of this book, respectively). We categorize methods for NAS according to three dimensions: search space, search strategy, and performance estimation strategy:

- **Search Space.** The search space defines which architectures can be represented in principle. Incorporating prior knowledge about properties well-suited for a task can reduce the size of the search space and simplify the

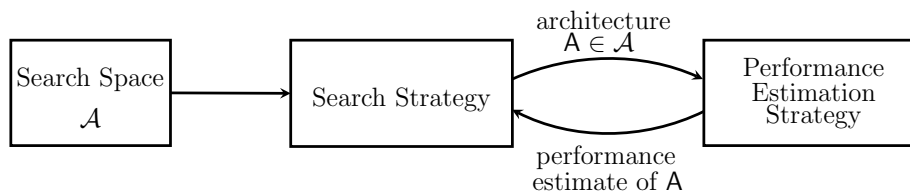


Figure 3.1: Abstract illustration of Neural Architecture Search methods. A search strategy selects an architecture  $A$  from a predefined search space  $\mathcal{A}$ . The architecture is passed to a performance estimation strategy, which returns the estimated performance of  $A$  to the search strategy.

search. However, this also introduces a human bias, which may prevent finding novel architectural building blocks that go beyond the current human knowledge.

- **Search Strategy.** The search strategy details how to explore the search space. It encompasses the classical exploration-exploitation trade-off since, on the one hand, it is desirable to find well-performing architectures quickly, while on the other hand, premature convergence to a region of suboptimal architectures should be avoided.
- **Performance Estimation Strategy.** The objective of NAS is typically to find architectures that achieve high predictive performance on unseen data. *Performance Estimation* refers to the process of estimating this performance: the simplest option is to perform a standard training and validation of the architecture on data, but this is unfortunately computationally expensive and limits the number of architectures that can be explored. Much recent research therefore focuses on developing methods that reduce the cost of these performance estimations.

We refer to Figure 3.1 for an illustration. The article is also structured according to these three dimensions: we start with discussing search spaces in Section 3.2, cover search strategies in Section 3.3, and outline approaches to performance estimation in Section 3.4. We conclude with an outlook on future directions in Section 3.5.

## 3.2 Search Space

The search space defines which neural architectures a NAS approach might discover in principle. We now discuss common search spaces from recent works.

A relatively simple search space is the space of *chain-structured neural networks*, as illustrated in Figure 3.2 (left). A chain-structured neural network architecture  $A$  can be written as a sequence of  $n$  layers, where the  $i$ 'th layer  $L_i$  receives its input from layer  $i - 1$  and its output serves as the input for layer  $i + 1$ , i.e.,  $A = L_n \circ \dots \circ L_1 \circ L_0$ . The search space is then parametrized by:

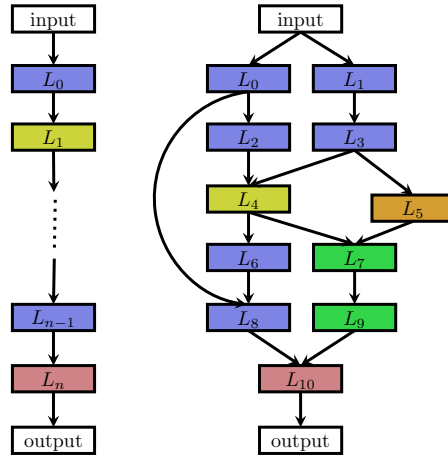


Figure 3.2: An illustration of different architecture spaces. Each node in the graphs corresponds to a layer in a neural network, e.g., a convolutional or pooling layer. Different layer types are visualized by different colors. An edge from layer  $L_i$  to layer  $L_j$  denotes that  $L_i$  receives the output of  $L_j$  as input. Left: an element of a chain-structured space. Right: an element of a more complex search space with additional layer types and multiple branches and skip connections.

(i) the (maximum) number of layers  $n$  (possibly unbounded); (ii) the type of operation every layer can execute, e.g., pooling, convolution, or more advanced layer types like depthwise separable convolutions [13] or dilated convolutions [67]; and (iii) hyperparameters associated with the operation, e.g., number of filters, kernel size and strides for a convolutional layer [4, 58, 10], or simply number of units for fully-connected networks [40]. Note that the parameters from (iii) are conditioned on (ii), hence the parametrization of the search space is not fixed-length but rather a conditional space.

Recent work on NAS [9, 21, 74, 22, 48, 11] incorporate modern design elements known from hand-crafted architectures such as skip connections, which allow to build complex, *multi-branch networks*, as illustrated in Figure 3.2 (right). In this case the input of layer  $i$  can be formally described as a function  $g_i(L_{i-1}^{out}, \dots, L_0^{out})$  combining previous layer outputs. Employing such a function results in significantly more degrees of freedom. Special cases of these multi-branch architectures are (i) the chain-structured networks (by setting  $g_i(L_{i-1}^{out}, \dots, L_0^{out}) = L_{i-1}^{out}$ ), (ii) Residual Networks [27], where previous layer outputs are summed ( $g_i(L_{i-1}^{out}, \dots, L_0^{out}) = L_{i-1}^{out} + L_j^{out}, j < i$ ) and (iii) DenseNets [28], where previous layer outputs are concatenated ( $g_i(L_{i-1}^{out}, \dots, L_0^{out}) = \text{concat}(L_{i-1}^{out}, \dots, L_0^{out})$ ).

Motivated by hand-crafted architectures consisting of repeated motifs [61, 27, 28], Zoph et al. [74] and Zhong et al. [70] propose to search for such motifs, dubbed *cells* or *blocks*, respectively, rather than for whole architectures. Zoph et al. [74] optimize two different kind of cells: a *normal cell* that preserves the dimensionality of the input and a *reduction cell* which reduces the spatial

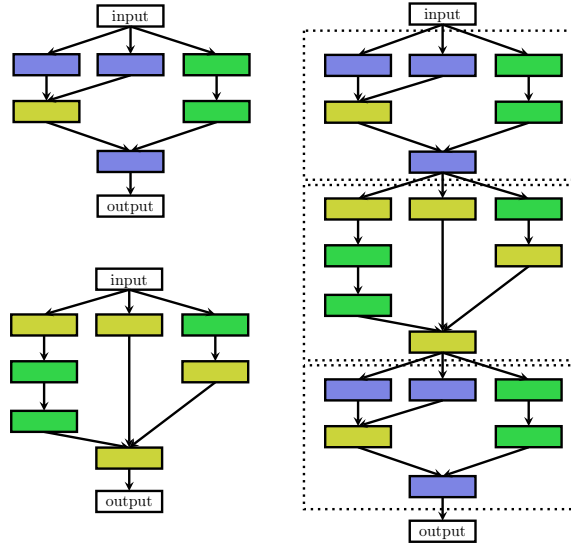


Figure 3.3: Illustration of the cell search space. Left: Two different cells, e.g., a normal cell (top) and a reduction cell (bottom) [74]. Right: an architecture built by stacking the cells sequentially. Note that cells can also be combined in a more complex manner, such as in multi-branch spaces, by simply replacing layers with cells.

dimension. The final architecture is then built by stacking these cells in a predefined manner, as illustrated in Figure 3.3. This search space has two major advantages compared to the ones discussed above:

1. The size of the search space is drastically reduced since cells can be comparably small. For example, Zoph et al. [74] estimate a seven-times speed-up compared to their previous work [73] while achieving better performance.
2. Cells can more easily be transferred to other datasets by adapting the number of cells used within a model. Indeed, Zoph et al. [74] transfer cells optimized on CIFAR-10 to ImageNet and achieve state-of-the-art performance.

Consequently, this cell-based search space was also successfully employed by many later works [48, 36, 45, 22, 11, 38, 71]. However, a new design-choice arises when using a cell-based search space, namely how to choose the *meta-architecture*: how many cells shall be used and how should they be connected to build the actual model? For example, Zoph et al. [74] build a sequential model from cells, in which each cell receives the outputs of the two preceding cells as input, while Cai et al. [11] employ the high-level structure of well-known manually designed architectures, such as DenseNet [28], and use their cells within these models. In principle, cells can be combined arbitrarily, e.g.,

within the multi-branch space described above by simply replacing layers with cells. Ideally, the meta-architecture should be optimized automatically as part of NAS; otherwise one easily ends up doing meta-architecture engineering and the search for the cell becomes overly simple if most of the complexity is already accounted for by the meta-architecture.

One step in the direction of optimizing meta-architectures is the hierarchical search space introduced by Liu et al. [37], which consists of several levels of motifs. The first level consists of the set of primitive operations, the second level of different motifs that connect primitive operations via a direct acyclic graphs, the third level of motifs that encode how to connect second-level motifs, and so on. The cell-based search space can be seen as a special case of this hierarchical search space where the number of levels is three, the second level motifs corresponds to the cells, and the third level is the hard-coded meta-architecture.

The choice of the search space largely determines the difficulty of the optimization problem: even for the case of the search space based on a single cell with fixed meta-architecture, the optimization problem remains (i) non-continuous and (ii) relatively high-dimensional (since more complex models tend to perform better, resulting in more design choices). We note that the architectures in many search spaces can be written as fixed-length vectors; e.g., the search space for each of the two cells by Zoph et al. [74] can be written as a 40-dimensional search space with categorical dimensions, each of which chooses between a small number of different building blocks and inputs. Similarly, unbounded search spaces can be constrained to have a maximal depth, giving rise to fixed-size search spaces with (potentially many) conditional dimensions.

In the next section, we discuss Search Strategies that are well-suited for these kinds of search spaces.

### 3.3 Search Strategy

Many different search strategies can be used to explore the space of neural architectures, including random search, Bayesian optimization, evolutionary methods, reinforcement learning (RL), and gradient-based methods. Historically, evolutionary algorithms were already used by many researchers to evolve neural architectures (and often also their weights) decades ago [see, e.g., 2, 55, 24, 54]. Yao [66] provides a literature review of work earlier than 2000.

Bayesian optimization celebrated several early successes in NAS since 2013, leading to state-of-the-art vision architectures [7], state-of-the-art performance for CIFAR-10 without data augmentation [19], and the first automatically-tuned neural networks to win competition datasets against human experts [40]. NAS became a mainstream research topic in the machine learning community after Zoph and Le [73] obtained competitive performance on the CIFAR-10 and Penn Treebank benchmarks with a search strategy based on reinforcement learning. While Zoph and Le [73] use vast computational resources to achieve this result (800 GPUs for three to four weeks), after their work, a wide variety of methods

have been published in quick succession to reduce the computational costs and achieve further improvements in performance.

To frame NAS as a *reinforcement learning* (RL) problem [4, 73, 70, 74], the generation of a neural architecture can be considered to be the agent’s action, with the action space identical to the search space. The agent’s reward is based on an estimate of the performance of the trained architecture on unseen data (see Section 3.4). Different RL approaches differ in how they represent the agent’s policy and how they optimize it: Zoph and Le [73] use a recurrent neural network (RNN) policy to sequentially sample a string that in turn encodes the neural architecture. They initially trained this network with the REINFORCE policy gradient algorithm, but in follow-up work use Proximal Policy Optimization (PPO) instead [74]. Baker et al. [4] use Q-learning to train a policy which sequentially chooses a layer’s type and corresponding hyperparameters. An alternative view of these approaches is as sequential decision processes in which the policy samples actions to generate the architecture sequentially, the environment’s “state” contains a summary of the actions sampled so far, and the (undiscounted) reward is obtained only after the final action. However, since no interaction with an environment occurs during this sequential process (no external state is observed, and there are no intermediate rewards), we find it more intuitive to interpret the architecture sampling process as the sequential generation of a single action; this simplifies the RL problem to a stateless multi-armed bandit problem.

A related approach was proposed by Cai et al. [10], who frame NAS as a sequential decision process: in their approach the state is the current (partially trained) architecture, the reward is an estimate of the architecture’s performance, and the action corresponds to an application of function-preserving mutations, dubbed network morphisms [12, 62], see also Section 3.4, followed by a phase of training the network. In order to deal with variable-length network architectures, they use a bi-directional LSTM to encode architectures into a fixed-length representation. Based on this encoded representation, actor networks decide on the sampled action. The combination of these two components constitute the policy, which is trained end-to-end with the REINFORCE policy gradient algorithm. We note that this approach will not visit the same state (architecture) twice so that strong generalization over the architecture space is required from the policy.

An alternative to using RL are *neuro-evolutionary* approaches that use evolutionary algorithms for optimizing the neural architecture. The first such approach for designing neural networks we are aware of dates back almost three decades: Miller et al. [43] use genetic algorithms to propose architectures and use backpropagation to optimize their weights. Many neuro-evolutionary approaches since then [2, 55, 54] use genetic algorithms to optimize both the neural architecture and its weights; however, when scaling to contemporary neural architectures with millions of weights for supervised learning tasks, SGD-based weight optimization methods currently outperform evolutionary ones<sup>1</sup>. More

---

<sup>1</sup>Some recent work shows that evolving even millions of weights is competitive to gradient-

recent neuro-evolutionary approaches [49, 58, 37, 48, 42, 65, 22] therefore again use gradient-based methods for optimizing weights and solely use evolutionary algorithms for optimizing the neural architecture itself. Evolutionary algorithms evolve a population of models, i.e., a set of (possibly trained) networks; in every evolution step, at least one model from the population is sampled and serves as a parent to generate offsprings by applying mutations to it. In the context of NAS, mutations are local operations, such as adding or removing a layer, altering the hyperparameters of a layer, adding skip connections, as well as altering training hyperparameters. After training the offsprings, their fitness (e.g., performance on a validation set) is evaluated and they are added to the population.

Neuro-evolutionary methods differ in how they sample parents, update populations, and generate offsprings. For example, Real et al. [49], Real et al. [48], and Liu et al. [37] use tournament selection [26] to sample parents, whereas Elsken et al. [22] sample parents from a multi-objective Pareto front using an inverse density. Real et al. [49] remove the worst individual from a population, while Real et al. [48] found it beneficial to remove the oldest individual (which decreases greediness), and Liu et al. [37] do not remove individuals at all. To generate offspring, most approaches initialize child networks randomly, while Elsken et al. [22] employ Lamarckian inheritance, i.e. knowledge (in the form of learned weights) is passed on from a parent network to its children by using network morphisms. Real et al. [49] also let an offspring inherit all parameters of its parent that are not affected by the applied mutation; while this inheritance is not strictly function-preserving it might also speed up learning compared to a random initialization. Moreover, they also allow mutating the learning rate which can be seen as a way for optimizing the learning rate schedule during NAS.

Real et al. [48] conduct a case study comparing RL, evolution, and random search (RS), concluding that RL and evolution perform equally well in terms of final test accuracy, with evolution having better anytime performance and finding smaller models. Both approaches consistently perform better than RS in their experiments, but with a rather small margin: RS achieved test errors of approximately 4% on CIFAR-10, while RL and evolution reached approximately 3.5% (after “model augmentation” where depth and number of filters was increased; the difference on the actual, non-augmented search space was approx. 2%). The difference was even smaller for Liu et al. [37], who reported a test error of 3.9% on CIFAR-10 and a top-1 validation error of 21.0% on ImageNet for RS, compared to 3.75% and 20.3% for their evolution-based method, respectively.

*Bayesian Optimization* (BO, see, e.g., [52]) is one of the most popular methods for hyperparameter optimization (see also Chapter 1 of this book), but it has not been applied to NAS by many groups since typical BO toolboxes are based on Gaussian processes and focus on low-dimensional continuous optimization

---

based optimization when only high-variance estimates of the gradient are available, e.g., for reinforcement learning tasks [50, 56, 15]. Nonetheless, for supervised learning tasks gradient-based optimization is by far the most common approach.

problems. Swersky et al. [59] and Kandasamy et al. [30] derive kernel functions for architecture search spaces in order to use classic GP-based BO methods, but so far without achieving new state-of-the-art performance. In contrast, several works use tree-based models (in particular, treed Parzen estimators [8], or random forests [29]) to effectively search very high-dimensional conditional spaces and achieve state-of-the-art performance on a wide range of problems, optimizing both neural architectures and their hyperparameters jointly [7, 19, 40, 68]. While a full comparison is lacking, there is preliminary evidence that these approaches can also outperform evolutionary algorithms [32].

Architectural search spaces have also been explored in a *hierarchical* manner, e.g., in combination with evolution [37] or by sequential model-based optimization [36]. Negrinho and Gordon [44] and Wistuba [64] exploit the tree-structure of their search space and use *Monte Carlo Tree Search*. Elsken et al. [21] propose a simple yet well performing *hill climbing* algorithm that discovers high-quality architectures by greedily moving in the direction of better performing architectures without requiring more sophisticated exploration mechanisms.

In contrast to the gradient-free optimization methods above, Liu et al. [38] propose a continuous relaxation of the search space to enable *gradient-based optimization*: instead of fixing a single operation  $o_i$  (e.g., convolution or pooling) to be executed at a specific layer, the authors compute a convex combination from a set of operations  $\{o_1, \dots, o_m\}$ . More specifically, given a layer input  $x$ , the layer output  $y$  is computed as  $y = \sum_{i=1}^m \lambda_i o_i(x)$ ,  $\lambda_i \geq 0$ ,  $\sum_{i=1}^m \lambda_i = 1$ , where the convex coefficients  $\lambda_i$  effectively parameterize the network architecture. Liu et al. [38] then optimize both the network weights and the network architecture by alternating gradient descent steps on training data for weights and on validation data for architectural parameters such as  $\lambda$ . Eventually, a discrete architecture is obtained by choosing the operation  $i$  with  $i = \arg \max_i \lambda_i$  for every layer. Shin et al. [53] and Ahmed and Torresani [1] also employ gradient-based optimization of neural architectures, however they only consider optimizing layer hyperparameters or connectivity patterns, respectively.

### 3.4 Performance Estimation Strategy

The search strategies discussed in Section 3.3 aim at finding a neural architecture  $A$  that maximizes some performance measure, such as accuracy on unseen data. To guide their search process, these strategies need to estimate the performance of a given architecture  $A$  they consider. The simplest way of doing this is to train  $A$  on training data and evaluate its performance on validation data. However, training each architecture to be evaluated from scratch frequently yields computational demands in the order of thousands of GPU days for NAS [73, 49, 74, 48].

To reduce this computational burden, performance can be estimated based on *lower fidelities* of the actual performance after full training (also denoted as proxy metrics). Such lower fidelities include shorter training times [74, 68], training on a subset of the data [33], on lower-resolution images [14], or with



less filters per layer [74, 48]. While these low-fidelity approximations reduce the computational cost, they also introduce bias in the estimate as performance will typically be underestimated. This may not be problematic as long as the search strategy only relies on ranking different architectures and the relative ranking remains stable. However, recent results indicate that this relative ranking can change dramatically when the difference between the cheap approximations and the “full” evaluation is too big [68], arguing for a gradual increase in fidelities [34, 23].

Another possible way of estimating an architecture’s performance builds upon learning curve extrapolation [60, 19, 31, 5, 47]. Domhan et al. [19] propose to extrapolate initial learning curves and terminate those predicted to perform poorly to speed up the architecture search process. Swersky et al. [60], Klein et al. [31], Baker et al. [5], Rawal and Miikkulainen [47] also consider architectural hyperparameters for predicting which partial learning curves are most promising. Training a surrogate model for predicting the performance of novel architectures is also proposed by Liu et al. [36], who do not employ learning curve extrapolation but support predicting performance based on architectural/cell properties and extrapolate to architectures/cells with larger size than seen during training. The main challenge for predicting the performances of neural architectures is that, in order to speed up the search process, good predictions in a relatively large search space need to be made based on relatively few evaluations.

Another approach to speed up performance estimation is to initialize the weights of novel architectures based on weights of other architectures that have been trained before. One way of achieving this, dubbed *network morphisms* [63], allows modifying an architecture while leaving the function represented by the network unchanged [10, 11, 21, 22]. This allows increasing capacity of networks successively and retaining high performance without requiring training from scratch. Continuing training for a few epochs can also make use of the additional capacity introduced by network morphisms. An advantage of these approaches is that they allow search spaces without an inherent upper bound on the architecture’s size [21]; on the other hand, strict network morphisms can only make architectures larger and may thus lead to overly complex architectures. This can be attenuated by employing approximate network morphisms that allow shrinking architectures [22].

*One-Shot Architecture Search* is another promising approach for speeding up performance estimation, which treats all architectures as different subgraphs of a supergraph (the one-shot model) and shares weights between architectures that have edges of this supergraph in common [51, 9, 45, 38, 6]. Only the weights of a single one-shot model need to be trained (in one of various ways), and architectures (which are just subgraphs of the one-shot model) can then be evaluated without any separate training by inheriting trained weights from the one-shot model. This greatly speeds up performance estimation of architectures, since no training is required (only evaluating performance on validation data). This approach typically incurs a large bias as it underestimates the actual performance of architectures severely; nevertheless, it allows ranking ar-

architectures reliably, since the estimated performance correlates strongly with the actual performance [6]. Different one-shot NAS methods differ in how the one-shot model is trained: ENAS [45] learns an RNN controller that samples architectures from the search space and trains the one-shot model based on approximate gradients obtained through REINFORCE. DARTS [38] optimizes all weights of the one-shot model jointly with a continuous relaxation of the search space obtained by placing a mixture of candidate operations on each edge of the one-shot model. Bender et al. [6] only train the one-shot model once and show that this is sufficient when deactivating parts of this model stochastically during training using path dropout. While ENAS and DARTS optimize a distribution over architectures during training, the approach of Bender et al. [6] can be seen as using a fixed distribution. The high performance obtainable by the approach of Bender et al. [6] indicates that the combination of weight sharing and a fixed (carefully chosen) distribution might (perhaps surprisingly) be the only required ingredients for one-shot NAS. Related to these approaches is meta-learning of hypernetworks that generate weights for novel architectures and thus requires only training the hypernetwork but not the architectures themselves [9]. The main difference here is that weights are not strictly shared but generated by the shared hypernetwork (conditional on the sampled architecture).

A general limitation of one-shot NAS is that the supergraph defined a-priori restricts the search space to its subgraphs. Moreover, approaches which require that the entire supergraph resides in GPU memory during architecture search will be restricted to relatively small supergraphs and search spaces accordingly and are thus typically used in combination with cell-based search spaces. While approaches based on weight-sharing have substantially reduced the computational resources required for NAS (from thousands to a few GPU days), it is currently not well understood which biases they introduce into the search if the sampling distribution of architectures is optimized along with the one-shot model. For instance, an initial bias in exploring certain parts of the search space more than others might lead to the weights of the one-shot model being better adapted for these architectures, which in turn would reinforce the bias of the search to these parts of the search space. This might result in premature convergence of NAS and might be one advantage of a fixed sampling distribution as used by Bender et al. [6]. In general, a more systematic analysis of biases introduced by different performance estimators would be a desirable direction for future work.

### 3.5 Future Directions

In this section, we discuss several current and future directions for research on NAS. Most existing work has focused on NAS for image classification. On the one hand, this provides a challenging benchmark since a lot of manual engineering has been devoted to finding architectures that perform well in this domain and are not easily outperformed by NAS. On the other hand, it is relatively easy to define a well-suited search space by utilizing knowledge from manual en-

gineering. This in turn makes it unlikely that NAS will find architectures that substantially outperform existing ones considerably since the found architectures cannot differ fundamentally. We thus consider it important to go beyond image classification problems by applying NAS to less explored domains. Notable first steps in this direction are applying NAS to language modeling [73], music modeling [47], image restoration [57] and network compression [3]; applications to reinforcement learning, generative adversarial networks, semantic segmentation, or sensor fusion could be further promising future directions.

An alternative direction is developing NAS methods for multi-task problems [35, 41] and for multi-objective problems [22, 20, 72], in which measures of resource efficiency are used as objectives along with the predictive performance on unseen data. Likewise, it would be interesting to extend RL/bandit approaches, such as those discussed in Section 3.3, to learn policies that are conditioned on a state that encodes task properties/resource requirements (i.e., turning the setting into a contextual bandit). A similar direction was followed by Ramachandran and Le [46] in extending one-shot NAS to generate different architectures depending on the task or instance on-the-fly. Moreover, applying NAS to searching for architectures that are more robust to adversarial examples [17] is an intriguing recent direction.

Related to this is research on defining more general and flexible search spaces. For instance, while the cell-based search space provides high transferability between different image classification tasks, it is largely based on human experience on image classification and does not generalize easily to other domains where the hard-coded hierarchical structure (repeating the same cells several times in a chain-like structure) does not apply (e.g., semantic segmentation or object detection). A search space which allows representing and identifying more general hierarchical structure would thus make NAS more broadly applicable, see Liu et al. [37] for first work in this direction. Moreover, common search spaces are also based on predefined building blocks, such as different kinds of convolutions and pooling, but do not allow identifying novel building blocks on this level; going beyond this limitation might substantially increase the power of NAS.

The comparison of different methods for NAS is complicated by the fact that measurements of an architecture’s performance depend on many factors other than the architecture itself. While most authors report results on the CIFAR-10 dataset, experiments often differ with regard to search space, computational budget, data augmentation, training procedures, regularization, and other factors. For example, for CIFAR-10, performance substantially improves when using a cosine annealing learning rate schedule [39], data augmentation by CutOut [18], by MixUp [69] or by a combination of factors [16], and regularization by Shake-Shake regularization [25] or scheduled drop-path [74]. It is therefore conceivable that improvements in these ingredients have a larger impact on reported performance numbers than the better architectures found by NAS. We thus consider the definition of common benchmarks to be crucial for a fair comparison of different NAS methods. A first step in this direction is the definition of a benchmark for joint architecture and hyperparameter search for a fully connected neural network with two hidden layers [32]. In this bench-

mark, nine discrete hyperparameters need to be optimized that control both architecture and optimization/regularization. All 62.208 possible hyperparameter combinations have been pre-evaluated such that different methods can be compared with low computational resources. However, the search space is still very simple compared to the spaces employed by most NAS methods. It would also be interesting to evaluate NAS methods not in isolation but as part of a full open-source AutoML system, where also hyperparameters [40, 49, 68], and data augmentation pipeline [16] are optimized along with NAS.

While NAS has achieved impressive performance, so far it provides little insights into why specific architectures work well and how similar the architectures derived in independent runs would be. Identifying common motifs, providing an understanding why those motifs are important for high performance, and investigating if these motifs generalize over different problems would be desirable.

### Acknowledgements

We would like to thank Esteban Real, Arber Zela, Gabriel Bender, Kenneth Stanley and Thomas Pfeil for feedback on earlier versions of this survey. This work has partly been supported by the European Research Council (ERC) under the European Union’s Horizon 2020 research and innovation programme under grant no. 716721.

### Bibliography

- [1] Ahmed, K., Torresani, L.: Maskconnect: Connectivity learning by gradient descent. In: European Conference on Computer Vision (ECCV) (2018)
- [2] Angeline, P.J., Saunders, G.M., Pollack, J.B.: An evolutionary algorithm that constructs recurrent neural networks. *IEEE transactions on neural networks* 5 1, 54–65 (1994)
- [3] Ashok, A., Rhinehart, N., Beainy, F., Kitani, K.M.: N2n learning: Network to network compression via policy gradient reinforcement learning. In: International Conference on Learning Representations (2018)
- [4] Baker, B., Gupta, O., Naik, N., Raskar, R.: Designing neural network architectures using reinforcement learning. In: International Conference on Learning Representations (2017)
- [5] Baker, B., Gupta, O., Raskar, R., Naik, N.: Accelerating Neural Architecture Search using Performance Prediction. In: NIPS Workshop on Meta-Learning (2017)
- [6] Bender, G., Kindermans, P.J., Zoph, B., Vasudevan, V., Le, Q.: Understanding and simplifying one-shot architecture search. In: International Conference on Machine Learning (2018)

- [7] Bergstra, J., Yamins, D., Cox, D.D.: Making a science of model search: Hyperparameter optimization in hundreds of dimensions for vision architectures. In: ICML (2013)
- [8] Bergstra, J.S., Bardenet, R., Bengio, Y., Kégl, B.: Algorithms for hyperparameter optimization. In: Shawe-Taylor, J., Zemel, R.S., Bartlett, P.L., Pereira, F., Weinberger, K.Q. (eds.) *Advances in Neural Information Processing Systems* 24. pp. 2546–2554 (2011)
- [9] Brock, A., Lim, T., Ritchie, J.M., Weston, N.: SMASH: one-shot model architecture search through hypernetworks. In: *NIPS Workshop on Meta-Learning* (2017)
- [10] Cai, H., Chen, T., Zhang, W., Yu, Y., Wang, J.: Efficient architecture search by network transformation. In: *Association for the Advancement of Artificial Intelligence* (2018)
- [11] Cai, H., Yang, J., Zhang, W., Han, S., Yu, Y.: Path-Level Network Transformation for Efficient Architecture Search. In: *International Conference on Machine Learning* (Jun 2018)
- [12] Chen, T., Goodfellow, I.J., Shlens, J.: Net2net: Accelerating learning via knowledge transfer. In: *International Conference on Learning Representations* (2016)
- [13] Chollet, F.: Xception: Deep learning with depthwise separable convolutions. *arXiv:1610.02357* (2016)
- [14] Chrabaszcz, P., Loshchilov, I., Hutter, F.: A downsampled variant of imagenet as an alternative to the CIFAR datasets. *CoRR abs/1707.08819* (2017)
- [15] Chrabaszcz, P., Loshchilov, I., Hutter, F.: Back to basics: Benchmarking canonical evolution strategies for playing atari. In: *Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18*. pp. 1419–1426. *International Joint Conferences on Artificial Intelligence Organization* (Jul 2018)
- [16] Cubuk, E.D., Zoph, B., Mane, D., Vasudevan, V., Le, Q.V.: AutoAugment: Learning Augmentation Policies from Data. In: *arXiv:1805.09501* (May 2018)
- [17] Cubuk, E.D., Zoph, B., Schoenholz, S.S., Le, Q.V.: Intriguing Properties of Adversarial Examples. In: *arXiv:1711.02846* (Nov 2017)
- [18] Devries, T., Taylor, G.W.: Improved regularization of convolutional neural networks with cutout. *arXiv preprint abs/1708.04552* (2017)

- [19] Domhan, T., Springenberg, J.T., Hutter, F.: Speeding up automatic hyperparameter optimization of deep neural networks by extrapolation of learning curves. In: Proceedings of the 24th International Joint Conference on Artificial Intelligence (IJCAI) (2015)
- [20] Dong, J.D., Cheng, A.C., Juan, D.C., Wei, W., Sun, M.: Dpp-net: Device-aware progressive search for pareto-optimal neural architectures. In: European Conference on Computer Vision (2018)
- [21] Elsken, T., Metzen, J.H., Hutter, F.: Simple And Efficient Architecture Search for Convolutional Neural Networks. In: NIPS Workshop on Meta-Learning (2017)
- [22] Elsken, T., Metzen, J.H., Hutter, F.: Efficient Multi-objective Neural Architecture Search via Lamarckian Evolution. ArXiv e-prints (Apr 2018)
- [23] Falkner, S., Klein, A., Hutter, F.: BOHB: Robust and efficient hyperparameter optimization at scale. In: Dy, J., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 80, pp. 1436–1445. PMLR, Stockholmsmässan, Stockholm Sweden (10–15 Jul 2018)
- [24] Floreano, D., Dürr, P., Mattiussi, C.: Neuroevolution: from architectures to learning. *Evolutionary Intelligence* 1(1), 47–62 (2008)
- [25] Gastaldi, X.: Shake-shake regularization. In: International Conference on Learning Representations Workshop (2017)
- [26] Goldberg, D.E., Deb, K.: A comparative analysis of selection schemes used in genetic algorithms. In: Foundations of Genetic Algorithms. pp. 69–93. Morgan Kaufmann (1991)
- [27] He, K., Zhang, X., Ren, S., Sun, J.: Deep Residual Learning for Image Recognition. In: Conference on Computer Vision and Pattern Recognition (2016)
- [28] Huang, G., Liu, Z., Weinberger, K.Q.: Densely Connected Convolutional Networks. In: Conference on Computer Vision and Pattern Recognition (2017)
- [29] Hutter, F., Hoos, H., Leyton-Brown, K.: Sequential model-based optimization for general algorithm configuration. In: LION. pp. 507–523 (2011)
- [30] Kandasamy, K., Neiswanger, W., Schneider, J., Póczos, B., Xing, E.: Neural Architecture Search with Bayesian Optimisation and Optimal Transport. arXiv:1802.07191 (Feb 2018)
- [31] Klein, A., Falkner, S., Springenberg, J.T., Hutter, F.: Learning curve prediction with Bayesian neural networks. In: International Conference on Learning Representations (2017)

- [32] Klein, A., Christiansen, E., Murphy, K., Hutter, F.: Towards reproducible neural architecture and hyperparameter search. In: ICML 2018 Workshop on Reproducibility in ML (RML 2018) (2018)
- [33] Klein, A., Falkner, S., Bartels, S., Hennig, P., Hutter, F.: Fast Bayesian Optimization of Machine Learning Hyperparameters on Large Datasets. In: Singh, A., Zhu, J. (eds.) Proceedings of the 20th International Conference on Artificial Intelligence and Statistics. Proceedings of Machine Learning Research, vol. 54, pp. 528–536. PMLR, Fort Lauderdale, FL, USA (20–22 Apr 2017)
- [34] Li, L., Jamieson, K., DeSalvo, G., Rostamizadeh, A., Talwalkar, A.: Hyperband: bandit-based configuration evaluation for hyperparameter optimization. In: International Conference on Learning Representations (2017)
- [35] Liang, J., Meyerson, E., Miikkulainen, R.: Evolutionary Architecture Search For Deep Multitask Networks. In: arXiv:1803.03745 (Mar 2018)
- [36] Liu, C., Zoph, B., Neumann, M., Shlens, J., Hua, W., Li, L.J., Fei-Fei, L., Yuille, A., Huang, J., Murphy, K.: Progressive Neural Architecture Search. In: European Conference on Computer Vision (2018)
- [37] Liu, H., Simonyan, K., Vinyals, O., Fernando, C., Kavukcuoglu, K.: Hierarchical Representations for Efficient Architecture Search. In: International Conference on Learning Representations (2018)
- [38] Liu, H., Simonyan, K., Yang, Y.: Darts: Differentiable architecture search. In: arXiv:1806.09055 (2018)
- [39] Loshchilov, I., Hutter, F.: Sgdr: Stochastic gradient descent with warm restarts. In: International Conference on Learning Representations (2017)
- [40] Mendoza, H., Klein, A., Feurer, M., Springenberg, J., Hutter, F.: Towards Automatically-Tuned Neural Networks. In: International Conference on Machine Learning, AutoML Workshop (Jun 2016)
- [41] Meyerson, E., Miikkulainen, R.: Pseudo-task Augmentation: From Deep Multitask Learning to Intratask Sharing and Back. In: arXiv:1803.03745 (Mar 2018)
- [42] Miikkulainen, R., Liang, J., Meyerson, E., Rawal, A., Fink, D., Francon, O., Raju, B., Shahrzad, H., Navruzyan, A., Duffy, N., Hodjat, B.: Evolving Deep Neural Networks. In: arXiv:1703.00548 (Mar 2017)
- [43] Miller, G., Todd, P., Hedge, S.: Designing neural networks using genetic algorithms. In: 3rd International Conference on Genetic Algorithms (ICGA’89) (1989)
- [44] Negrinho, R., Gordon, G.: DeepArchitect: Automatically Designing and Training Deep Architectures. arXiv:1704.08792 (2017)

- [45] Pham, H., Guan, M.Y., Zoph, B., Le, Q.V., Dean, J.: Efficient neural architecture search via parameter sharing. In: International Conference on Machine Learning (2018)
- [46] Ramachandran, P., Le, Q.V.: Dynamic Network Architectures. In: AutoML 2018 (ICML workshop) (2018)
- [47] Rawal, A., Miikkulainen, R.: From Nodes to Networks: Evolving Recurrent Neural Networks. In: arXiv:1803.04439 (Mar 2018)
- [48] Real, E., Aggarwal, A., Huang, Y., Le, Q.V.: Regularized Evolution for Image Classifier Architecture Search. In: arXiv:1802.01548 (Feb 2018)
- [49] Real, E., Moore, S., Selle, A., Saxena, S., Suematsu, Y.L., Le, Q.V., Kurakin, A.: Large-scale evolution of image classifiers. International Conference on Machine Learning (2017)
- [50] Salimans, T., Ho, J., Chen, X., Sutskever, I.: Evolution strategies as a scalable alternative to reinforcement learning. arXiv preprint (2017)
- [51] Saxena, S., Verbeek, J.: Convolutional neural fabrics. In: Lee, D.D., Sugiyama, M., Luxburg, U.V., Guyon, I., Garnett, R. (eds.) Advances in Neural Information Processing Systems 29, pp. 4053–4061. Curran Associates, Inc. (2016)
- [52] Shahriari, B., Swersky, K., Wang, Z., Adams, R.P., de Freitas, N.: Taking the human out of the loop: A review of bayesian optimization. Proceedings of the IEEE 104(1), 148–175 (Jan 2016)
- [53] Shin, R., Packer, C., Song, D.: Differentiable neural network architecture search. In: International Conference on Learning Representations Workshop (2018)
- [54] Stanley, K.O., D’Ambrosio, D.B., Gauci, J.: A hypercube-based encoding for evolving large-scale neural networks. *Artif. Life* 15(2), 185–212 (Apr 2009), <http://dx.doi.org/10.1162/artl.2009.15.2.15202>
- [55] Stanley, K.O., Miikkulainen, R.: Evolving neural networks through augmenting topologies. *Evolutionary Computation* 10, 99–127 (2002)
- [56] Such, F.P., Madhavan, V., Conti, E., Lehman, J., Stanley, K.O., Clune, J.: Deep neuroevolution: Genetic algorithms are a competitive alternative for training deep neural networks for reinforcement learning. arXiv preprint (2017)
- [57] Suganuma, M., Ozay, M., Okatani, T.: Exploiting the potential of standard convolutional autoencoders for image restoration by evolutionary search. In: Dy, J., Krause, A. (eds.) Proceedings of the 35th International Conference on Machine Learning. Proceedings of Machine Learning Research, vol. 80, pp. 4771–4780. PMLR, Stockholmsmässan, Stockholm Sweden (10–15 Jul 2018)



- [58] Suganuma, M., Shirakawa, S., Nagao, T.: A genetic programming approach to designing convolutional neural network architectures. In: Genetic and Evolutionary Computation Conference (2017)
- [59] Swersky, K., Duvenaud, D., Snoek, J., Hutter, F., Osborne, M.: Raiders of the lost architecture: Kernels for bayesian optimization in conditional parameter spaces. In: NIPS Workshop on Bayesian Optimization in Theory and Practice (2013)
- [60] Swersky, K., Snoek, J., Adams, R.P.: Freeze-thaw bayesian optimization (2014)
- [61] Szegedy, C., Vanhoucke, V., Ioffe, S., Shlens, J., Wojna, Z.: Rethinking the Inception Architecture for Computer Vision. In: Conference on Computer Vision and Pattern Recognition (2016)
- [62] Wei, T., Wang, C., Chen, C.W.: Modularized morphing of neural networks. arXiv:1701.03281 (2017)
- [63] Wei, T., Wang, C., Rui, Y., Chen, C.W.: Network morphism. In: International Conference on Machine Learning (2016)
- [64] Wistuba, M.: Finding Competitive Network Architectures Within a Day Using UCT. In: arXiv:1712.07420 (Dec 2017)
- [65] Xie, L., Yuille, A.: Genetic CNN. In: International Conference on Computer Vision (2017)
- [66] Yao, X.: Evolving artificial neural networks. Proceedings of the IEEE 87(9), 1423–1447 (Sept 1999)
- [67] Yu, F., Koltun, V.: Multi-scale context aggregation by dilated convolutions (2016)
- [68] Zela, A., Klein, A., Falkner, S., Hutter, F.: Towards automated deep learning: Efficient joint neural architecture and hyperparameter search. In: ICML 2018 Workshop on AutoML (AutoML 2018) (2018)
- [69] Zhang, H., Cissé, M., Dauphin, Y.N., Lopez-Paz, D.: mixup: Beyond empirical risk minimization. arXiv preprint abs/1710.09412 (2017)
- [70] Zhong, Z., Yan, J., Wu, W., Shao, J., Liu, C.L.: Practical block-wise neural network architecture generation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2423–2432 (2018)
- [71] Zhong, Z., Yang, Z., Deng, B., Yan, J., Wu, W., Shao, J., Liu, C.L.: Block-qnn: Efficient block-wise neural network architecture generation. arXiv preprint (2018)
- [72] Zhou, Y., Ebrahimi, S., Arık, S., Yu, H., Liu, H., Diamos, G.: Resource-efficient neural architect. In: arXiv:1806.07912 (2018)

- [73] Zoph, B., Le, Q.V.: Neural architecture search with reinforcement learning. In: International Conference on Learning Representations (2017)
- [74] Zoph, B., Vasudevan, V., Shlens, J., Le, Q.V.: Learning transferable architectures for scalable image recognition. In: Conference on Computer Vision and Pattern Recognition (2018)