

# A Study on Encodings for Neural Architecture Search

**Colin White**

*abacus.ai (prev. RealityEngines.AI)*

COLIN@ABACUS.AI

**Willie Neiswanger**

*Carnegie Mellon University and Petuum, Inc.*

WILLIE@CS.CMU.EDU

**Sam Nolen**

*abacus.ai (prev. RealityEngines.AI)*

SAM@ABACUS.AI

**Yash Savani**

*abacus.ai (prev. RealityEngines.AI)*

YASH@ABACUS.AI

## Abstract

Neural architecture search (NAS) has been extensively studied in the past few years. A popular approach is to represent each neural architecture in the search space as a directed acyclic graph (DAG), and then search over all DAGs by encoding the adjacency matrix and list of operations as a set of hyperparameters. Recent work has demonstrated that even small changes to the way each architecture is encoded can have a significant effect on the performance of NAS algorithms (White et al., 2019; Ying et al., 2019).

In this work, we present the first formal study on the effect of architecture encodings for NAS, including a theoretical grounding and an empirical study. First we formally define architecture encodings and give a theoretical characterization on the scalability of the encodings we study. Then we identify the main encoding-dependent subroutines which NAS algorithms employ, running experiments to show which encodings work best with each subroutine for many popular algorithms. The experiments act as an ablation study for prior work, disentangling the algorithmic and encoding-based contributions, as well as a guideline for future work. Our results demonstrate that NAS encodings are an important design decision which can have a significant impact on overall performance.<sup>1</sup>

## 1. Introduction

In the past few years, the field of neural architecture search (NAS) has seen a steep rise in interest (Elsken et al., 2018), due to the promise of automatically designing specialized neural architectures for any given problem. Techniques for NAS span evolutionary search, Bayesian optimization, reinforcement learning, gradient-based methods, and neural predictor methods. Many NAS instantiations can be described by the optimization problem  $\min_{a \in A} f(a)$ , where  $A$  denotes a large set of neural architectures, and  $f(a)$  denotes the objective function of interest for  $a$ , which is usually a combination of validation accuracy, latency, or number of parameters. A popular approach is to describe each neural architecture  $a$  as a labeled directed acyclic graph (DAG), where each node or edge represents an operation.

Due to the complexity of DAG structures and the large size of the space, neural architecture search is typically a highly non-convex, challenging optimization problem. A natural

---

1. This work was extended to a full-length paper here: <https://arxiv.org/abs/2007.04965>. Our code is available at <https://github.com/naszilla/nas-encodings>.

consideration when designing a NAS algorithm is therefore, *how should we encode the neural architectures to maximize performance?* For example, NAS algorithms may involve manipulating or perturbing architectures, or training a model to predict the accuracy of a given architecture; as a consequence, the representation of the architectures may significantly change the outcome of these subroutines. The majority of prior work has not explicitly considered this question, opting to use a standard encoding consisting of the adjacency matrix of the DAG along with a list of the operations. Two recent papers have shown that even small changes to the architecture encoding can make a substantial difference in the performance of the NAS algorithm (Ying et al., 2019; White et al., 2019). It is not obvious how to formally define an encoding for NAS, as papers define encodings in different ways, inadvertently using encodings which are incompatible with other NAS algorithms.

In this work, we provide the first formal study on NAS encoding schemes, including a theoretical grounding as well as a set of experimental results in different settings. We define an encoding as a multi-function from an architecture to a real-valued tensor. We define a number of the most common encodings from prior work, and we theoretically characterize the scalability of each one. Next, we identify three major encoding-dependent subroutines used in NAS algorithms: *sample random architecture*, *perturb architecture*, and *train predictor model*. We show which encodings perform best for each subroutine by testing each encoding within each subroutine for many popular NAS algorithms. Our experiments retroactively provide an ablation study for prior work by disentangling the algorithmic contributions with the encoding-based contributions. Overall, our results show that NAS encodings are an important design decision which must be taken into account not only at the algorithmic level, but at the subroutine level, and which can have a significant impact on the final performance. Our experimental results follow the guidelines laid out in the NAS research checklist (Lindauer and Hutter, 2019). In particular, we experiment on NASBench-101 (Ying et al., 2019), a popular NAS benchmark, and we release our code.

**Our contributions.** We summarize our main contributions below.

- We demonstrate that the choice of encoding is an important, nontrivial question that should be considered not only at the algorithmic level, but at the subroutine level.
- We give a theoretical grounding for NAS encodings, including a characterization of the scalability of each encoding.
- We give an experimental study of architecture encodings for NAS algorithms, disentangling the algorithmic contributions from the encoding-based contributions of prior work, and laying out recommendations for best encodings to use in different settings as guidance for future work.

## 2. Related Work

**Neural architecture search.** NAS has been studied for at least two decades and has received significant attention in recent years (Kitano, 1990; Stanley and Miikkulainen, 2002; Zoph and Le, 2017). Some of the most popular techniques for NAS include evolutionary algorithms (Maziarz et al., 2018), reinforcement learning (Pham et al., 2018; Tan and Le, 2019), Bayesian optimization (BO) (Kandasamy et al., 2018), gradient descent (Liu et al., 2018), neural predictors (Wen et al., 2019), and local search (White et al., 2020). Recent

papers have highlighted the need for fair and reproducible NAS comparisons (Li and Talwalkar, 2019; Ying et al., 2019; Lindauer and Hutter, 2019). See the recent survey on NAS (Elsken et al., 2018) for more information on NAS research.

**Encoding schemes.** Most prior NAS work has used the adjacency matrix encoding. A continuous-valued variant has been shown to be more effective for some NAS algorithms (Ying et al., 2019). The path encoding works well for neural predictor methods (Wei et al., 2020), and truncating the path encoding leads to a small information loss (White et al., 2019). Some previous work use graph convolutional networks (GCN) as a subroutine in NAS (Shi et al., 2019; Zhang et al., 2019), which requires retraining for each new dataset or search space. Other work has used intermediate encodings to reduce the complexity of the DAG (Stanley et al., 2009; Irwin-Harris et al., 2019), or added summary statistics to the encoding of feedforward networks (Sun et al., 2019). To the best of our knowledge, no paper has conducted a formal study of encodings involving more than two encodings.

### 3. Encodings for NAS

We denote a set of neural architectures  $a$  by  $A$  (called a search space), and we define an objective function  $\ell : A \rightarrow \mathbb{R}$ , where  $\ell(a)$  is typically a combination of the neural network accuracy, model parameters, or FLOPS. We define a neural network encoding as an integer  $d$  and a multifunction  $e : A \rightarrow \mathbb{R}^d$  from a set of neural architectures  $A$  to a  $d$ -dimensional Euclidean space  $\mathbb{R}^d$ , and we define a NAS algorithm  $\mathcal{A}$  which takes as input a triple  $(A, \ell, e)$ , and outputs an architecture  $a$ , with the goal that  $\ell(a)$  is as close to  $\max_{a \in A} \ell(a)$  as possible. Based on this definition, we consider an encoding  $e$  to be a fixed transformation, independent of  $\ell$ . In particular, NAS components that use  $\ell$  to learn a transformation of an input architecture (such as graph convolutional networks or autoencoders), are considered part of the NAS algorithm rather than the encoding. This is consistent with prior definitions of encodings (Ying et al., 2019; Talbi, 2020).

We define eight different encodings split into two paradigms: adjacency matrix-based and path-based encodings. We assume that each architecture is represented by a DAG with at most  $n$  nodes, at most  $k$  edges, and  $q$  choices of operations on each node. For brevity, we focus on the case where nodes represent operations, though our analysis extends similarly to formulations where edges represent operations. Most of the following encodings have been defined in prior work (Ying et al., 2019; White et al., 2019; Talbi, 2020), and we will see in the next section that each encoding is useful for some part of the NAS pipeline.

**Adjacency matrix encodings.** We first consider a class of encodings that are based on representations of the adjacency matrix. These are the most common types of encodings used in current NAS research. For visualizations of these encodings, see Figure A.1 (b).

The *one-hot adjacency matrix encoding* is created by row-major vectorizing (i.e. flattening) the architecture adjacency matrix and concatenating it with a list of node operation labels. Each position in the operation list is a single integer-valued feature, where each operation is denoted by a different integer. The total dimension is  $n(n-1)/2 + n$ . In the *categorical adjacency matrix encoding*, the adjacency matrix is first flattened (similar to the one-hot encoding described previously), and is then defined as a list of the indices each of which specifies one of the  $n(n-1)/2$  possible edges in the adjacency matrix. To

ensure a fixed length encoding, each architecture is represented by  $k$  features, where  $k$  is the maximum number of possible edges. We again concatenate this representation with a list of operations, yielding a total dimensionality of  $k + n$ . Finally, the *continuous adjacency matrix encoding* is similar to the one-hot encoding, but each of the features for each edge can take on any real value in  $[0, 1]$ , rather than just  $\{0, 1\}$ . We also add a feature representing the number of edges,  $1 \leq K \leq k$ . The list of operations is encoded the same way as before. The architecture is created by choosing the  $K$  edges with the largest continuous features. The dimension is  $n(n-1)/2 + n + 1$ . The disadvantage of adjacency matrix-based encodings is that nodes are arbitrarily assigned indices in the matrix, which means one architecture can have many different representations (in other words,  $e^{-1}$  is not onto).

**Path-based encodings.** Path-based encodings are representations of a neural architecture that are based on the set of paths from input to output that are present within the architecture DAG. For visualizations of these encodings, see Figure A.1 (c).

The *one-hot path encoding* is created by giving a binary feature to each possible path from the input node to the output node in the DAG (for example: `input-conv1x1-maxpool3x3-output`). The total dimension is  $\sum_{i=0}^n q^i = (q^{n+1} - 1)/(q - 1)$ . The *truncated one-hot path encoding*, simply truncates this encoding to only include paths of length  $x$ . The new dimension is  $\sum_{i=0}^x q^i$ . The *categorical path encoding*, is defined as a list of indices each of which specifies one of the  $\sum_{i=0}^n q^i$  possible paths. The *continuous path encoding* consists of a real-valued feature  $[0, 1]$  for each potential path, as well as a feature representing the number of paths. Just like the one-hot path encoding, the continuous path encoding can be truncated. Path-based encodings have the advantage that nodes are not arbitrarily assigned indices, and also that isomorphisms are automatically mapped to the same encoding. Path-based encodings have the disadvantage that different architectures can map to the same encoding ( $e$  is not onto).

### 3.1 The scalability of encodings

In this section, we discuss the scalability of the NAS encodings with respect to architecture size. We focus on the one-hot variants of the encodings, but our analysis extends to all encodings. We show that the path encoding can be truncated significantly while maintaining its performance, while the adjacency matrix cannot be truncated at all without sacrificing performance, and we back up our theoretical results with experimental observations in Appendix B. The one-hot path encoding has shown to be very effective on the smaller NASBench-101 and NASBench-201 datasets (Wei et al., 2020). However, other work has questioned whether its exponential scaling allows it to perform well on very large search spaces (Talbi, 2020). The vast majority of features in the one-hot encoding correspond to single line paths using the full set of nodes, which is not common during NAS algorithms, nor is it likely to be effective. Prior work has made an attempt to prove that truncating the path encoding does not harm the performance of NAS algorithms (White et al., 2019).

Consider the popular *sample random architecture* method: given  $n$ ,  $r$ , and  $k < \frac{n(n-1)}{2}$ , (1) denote  $n$  nodes by 1 to  $n$ , and label each node with one of  $r$  operations. (2) for all  $i < j$ , add edge  $(i, j)$  with probability  $\frac{2k}{n(n-1)}$ . (3) if there is no path from node 1 to node  $n$ , goto(1). In our analysis, we consider a random graph  $G'_{n,k,r}$  returned after the first two

steps (taking into account the third step later on). Let  $a_{n,k,\ell}$  denote the expected number of paths from node 1 to node  $n$  of length  $\ell$  in  $G'_{n,k,r}$ , and define  $b(k, x) = \frac{\sum_{\ell=1}^x a_{n,k,\ell}}{\sum_{\ell=1}^n a_{n,k,\ell}}$ .

Prior work was only able to show that  $b(k, x) > 1 - 1/n^2$  when  $k < n + O(1)$  and  $x > \log n$  (White et al., 2019). Now we present our main result, which gives a full characterization of  $b(k, x)$  up to constant factors. We show that, for the purposes of NAS, truncating the path encoding to length  $r^{k/n}$  contains almost exactly the same information as the full path encoding, and it cannot be any smaller. If  $k \leq n \log n$ , this is shorter than the one-hot adjacency matrix encoding. In Appendix A, we give a proof of the following theorem, and we also show that it is not possible to truncate the adjacency matrix encoding by a single bit without a significant loss of information.

**Theorem 3.1.** Given  $10 \leq n \leq k \leq \frac{n(n-1)}{2}$ , and  $c > 3$ , for  $x > 2ec \cdot \frac{k}{n}$ ,  $b(k, x) > 1 - c^{-x+1}$ , and for  $x < \frac{1}{2ec} \cdot \frac{k}{n}$ ,  $b(k, x) < -2^{\frac{k}{2n}}$ .

## 4. Experiments

We run experiments on the NASBench-101 dataset (Ying et al., 2019), which consists of over 423,000 pretrained neural architectures on the CIFAR-10 dataset. We split up our experiments based on the three encoding-dependent subroutines: *sample random architecture*, *perturb architecture*, and *train predictor model*. These three subroutines are the only encoding-dependent building blocks necessary for most NAS algorithms.

**Sample random architecture.** Most NAS algorithms use a subroutine to draw an architecture randomly from the search space. Although this operation is more generally parameterized by a distribution over the search space, it is often instantiated with the choice of architecture encoding. Given an encoding, we define a subroutine by randomly sampling each feature uniformly at random. We also compared to sampling each *architecture* uniformly at random from the search space (which does not correspond to any encoding).

**Perturb architecture.** Another common subroutine in NAS algorithms is to make a small change to a given architecture. The type of modification depends on the encoding. For example, a perturbation might be to change an operation, add or remove an edge, or add or remove a path. Given an encoding, we define a perturbation subroutine by resampling each feature of the encoding uniformly at random with a fixed probability.

**Train predictor model.** Many families of NAS algorithms use a subroutine which learns a model based on already queried architectures. This can take the form of a Gaussian process within Bayesian optimization (BO), or, more recently, a neural network. In the case of a Gaussian process model, the algorithm needs a distance metric between neural architectures, which is typically chosen as the edit distance between architecture encodings. In the case of a neural network model, the input are the encodings of the architectures, and the goal is typically to predict the accuracy of unseen architectures.

Many NAS algorithms use more than one subroutines, so in each experiment, we fix the encodings for all subroutines except the one we are testing. For each NAS subroutine, we choose experiments with algorithms which are highly dependent on the subroutine: Random Search for random sampling, regularized evolution (Real et al., 2019) and local

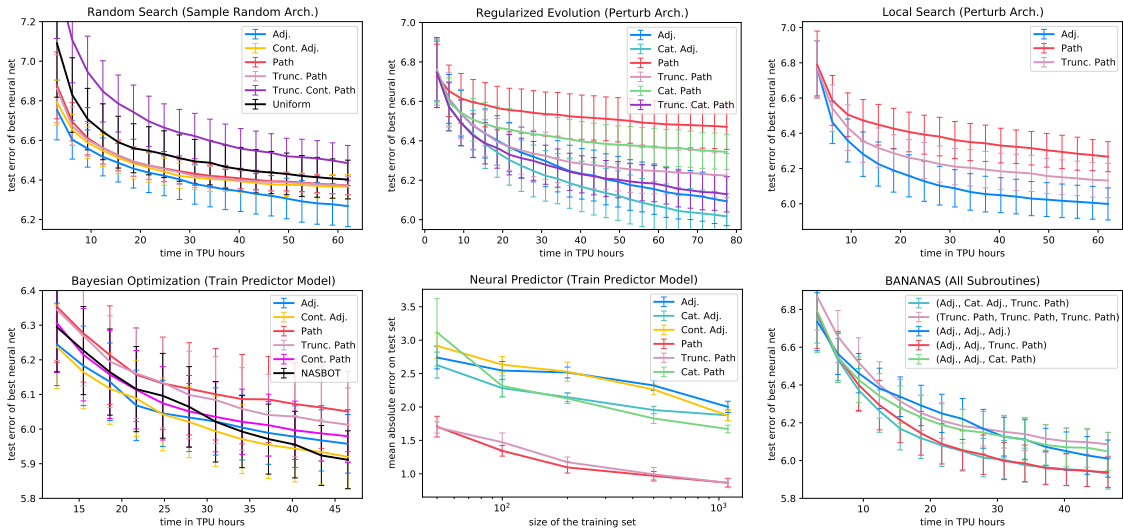


Figure 4.1: NAS experiments with different encodings, keeping all but one subroutine fixed: random sampling (top left), perturb architecture (top middle, top right), train predictor model (bottom left, bottom middle), or varying all three subroutines (bottom right).

search (White et al., 2020) for perturb architecture, and BO and BANANAS (White et al., 2020) for train predictor model. We run a final experiment on BANANAS varying all three subroutines. See Figure 4.1. We used the open source code for each algorithm, and we ran 300 trials for each experiment. In Appendix B, we give more details on our implementation of the baseline algorithms and our adherence to best practices (Lindauer and Hutter, 2019). We also present more experiments for NASBench-201, and we test each encoding’s ability to generalize beyond the training search space.

Depending on the subroutine, two encodings might behave identically, which is why not all encodings appear in each experiment. There is no overall best encoding; instead, different encodings perform best within different subroutines and algorithms. Categorical, one-hot, adjacency-based, path-based, and continuous encodings are all best in certain settings. Truncating the path encoding improves performance, consistent with our theoretical results in Section 3. Some of our findings explain the success of prior algorithms, e.g., regularized evolution uses the categorical adjacency encoding, BANANAS uses the path encoding, and local search uses the adjacency encoding.

## 5. Conclusion

In this paper, we give the first formal study of encoding schemes for neural architecture search. We define eight different encodings and characterize the scalability of each one. We then identify three encoding-dependent subroutines used by NAS algorithms, and we run experiments to find the best encoding for each subroutine in many popular algorithms. Our experimental results allow us to disentangle the algorithmic and encoding-based contributions of prior work, and act as a guideline for the encodings to use in future work. Overall, we show that encodings are an important, nontrivial design decision in the field of NAS.

## References

- Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey. *arXiv preprint arXiv:1808.05377*, 2018.
- William Irwin-Harris, Yanan Sun, Bing Xue, and Mengjie Zhang. A graph-based encoding for evolutionary convolutional neural network architecture design. In *2019 IEEE Congress on Evolutionary Computation (CEC)*, pages 546–553. IEEE, 2019.
- Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabas Poczos, and Eric P Xing. Neural architecture search with bayesian optimisation and optimal transport. In *Advances in Neural Information Processing Systems*, pages 2016–2025, 2018.
- Hiroaki Kitano. Designing neural networks using genetic algorithms with graph generation system. *Complex systems*, 4(4):461–476, 1990.
- Liam Li and Ameet Talwalkar. Random search and reproducibility for neural architecture search. *arXiv preprint arXiv:1902.07638*, 2019.
- Marius Lindauer and Frank Hutter. Best practices for scientific research on neural architecture search. *arXiv preprint arXiv:1909.02453*, 2019.
- Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. *arXiv preprint arXiv:1806.09055*, 2018.
- Krzysztof Maziarczyk, Andrey Khorlin, Quentin de Laroussilhe, and Andrea Gesmundo. Evolutionary-neural hybrid agents for architecture search. *arXiv preprint arXiv:1811.09828*, 2018.
- Willie Neiswanger, Kirthevasan Kandasamy, Barnabas Poczos, Jeff Schneider, and Eric Xing. Probo: a framework for using probabilistic programming in bayesian optimization. *arXiv preprint arXiv:1901.11515*, 2019.
- Hieu Pham, Melody Y Guan, Barret Zoph, Quoc V Le, and Jeff Dean. Efficient neural architecture search via parameter sharing. *arXiv preprint arXiv:1802.03268*, 2018.
- 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.
- Han Shi, Renjie Pi, Hang Xu, Zhenguo Li, James T Kwok, and Tong Zhang. Multi-objective neural architecture search via predictive network performance optimization. *arXiv preprint arXiv:1911.09336*, 2019.
- Pantelimon Stanica. Good lower and upper bounds on binomial coefficients. *Journal of Inequalities in Pure and Applied Mathematics*, 2001.
- Kenneth O Stanley and Risto Miikkulainen. Evolving neural networks through augmenting topologies. *Evolutionary computation*, 10(2):99–127, 2002.

- KO Stanley, DB D'Ambrosio, and J Gauci. A hypercube-based indirect encoding for evolving large-scale neural networks. In *Artificial Life*, volume 15(2), pages 185–212, 2009.
- Yanan Sun, Bing Xue, Mengjie Zhang, and Gary G Yen. Evolving deep convolutional neural networks for image classification. *IEEE Transactions on Evolutionary Computation*, 2019.
- El-Ghazali Talbi. Optimization of deep neural networks: a survey and unified taxonomy. 2020.
- Mingxing Tan and Quoc V Le. Efficientnet: Rethinking model scaling for convolutional neural networks. *arXiv preprint arXiv:1905.11946*, 2019.
- Chen Wei, Chuang Niu, Yiping Tang, and Jimin Liang. Npenas: Neural predictor guided evolution for neural architecture search. *arXiv preprint arXiv:2003.12857*, 2020.
- Wei Wen, Hanxiao Liu, Hai Li, Yiran Chen, Gabriel Bender, and Pieter-Jan Kindermans. Neural predictor for neural architecture search. *arXiv preprint arXiv:1912.00848*, 2019.
- Colin White, Willie Neiswanger, and Yash Savani. Bananas: Bayesian optimization with neural architectures for neural architecture search. *arXiv preprint arXiv:1910.11858*, 2019.
- Colin White, Sam Nolen, and Yash Savani. Local search is state of the art for nas benchmarks. *arXiv preprint arXiv:2005.02960*, 2020.
- Antoine Yang, Pedro M Esperança, and Fabio M Carlucci. Nas evaluation is frustratingly hard. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2020.
- Chris Ying, Aaron Klein, Esteban Real, Eric Christiansen, Kevin Murphy, and Frank Hutter. Nas-bench-101: Towards reproducible neural architecture search. *arXiv preprint arXiv:1902.09635*, 2019.
- Muhan Zhang, Shali Jiang, Zhicheng Cui, Roman Garnett, and Yixin Chen. D-vae: A variational autoencoder for directed acyclic graphs. In *Proceedings of the Annual Conference on Neural Information Processing Systems (NIPS)*, 2019.
- Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2017.



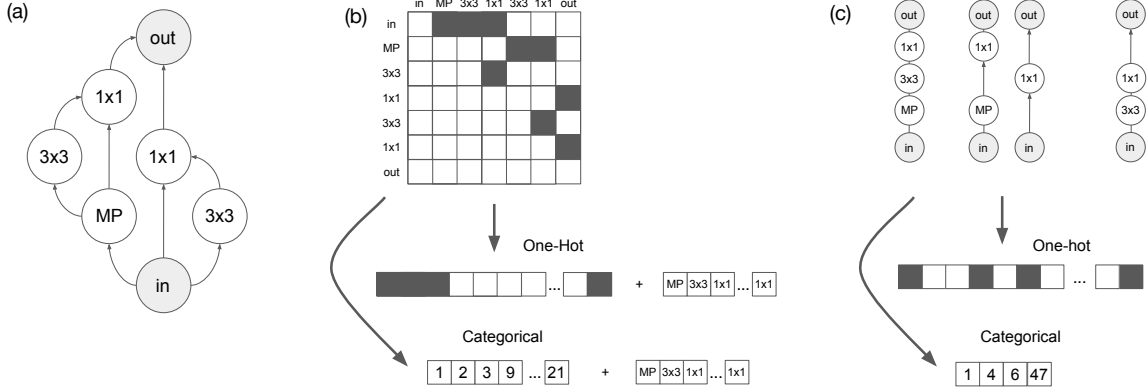


Figure A.1: (a) An example neural architecture  $a$ . (b) An adjacency matrix representation of  $a$ , showing two encodings. (c) A path-based representation of  $a$ , showing two encodings.

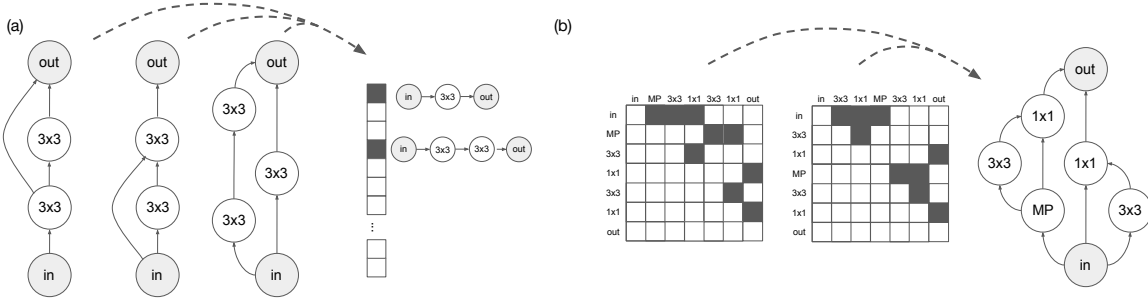


Figure A.2: (a) An example of three architectures that map to the same path encoding. (b) An example of two adjacency matrices that map to the same architecture.

## Appendix A. Details from Section 3 (Encodings for NAS)

We give the details from Section 3. We restate the random graph model here more formally.

**Definition A.1.** Given nonzero integers  $n, r$ , and  $k < n(n-1)/2$ , a random graph  $G_{n,k,r}$  is generated as follows:

- (1) Denote  $n$  nodes by 1 to  $n$  and label each node randomly with one of  $r$  operations.
- (2) For all  $i < j$ , add edge  $(i, j)$  with probability  $\frac{2k}{n(n-1)}$ .
- (3) If there is no path from node 1 to node  $n$ , **goto** (1).

Let  $G'_{n,k,r}$  denote the random graph outputted by the above procedure without step (3). Since the number of pairs  $(i, j)$  such that  $i < j$  is  $\frac{n(n-1)}{2}$ , the expected number of edges of  $G'_{n,k,r}$  is  $k$ . Define  $a_{n,k,\ell}$  as the expected number of paths from node 1 to node  $n$  of length  $\ell$  in  $G'_{n,k,r}$ . Formally, we set  $\mathcal{P} = \{\text{paths from node 1 to } n \text{ in } G'_{n,k,r}\}$ , and define

$$a_{n,k,\ell} = \mathbb{E} [|p \in \mathcal{P}| | |p| = \ell].$$

Recall that

$$b(k, x) = \frac{\sum_{\ell=1}^x a_{n,k,\ell}}{\sum_{\ell=1}^n a_{n,k,\ell}}.$$

In the next theorem, we give a full characterization of  $b(k, x)$  in terms of  $k$  and  $n$ , up to constant factors. We prove there exists a phase transition for  $b(k, x)$  at  $x = \frac{k}{n}$ . As noted by prior work (White et al., 2019), there are two caveats when applying this type of theorem to NAS performance. The theorem considers the distribution from Definition A.1, not the distribution of architectures encountered in a real search, and the most common paths in the distribution are not necessarily the ones with the most entropy in predicting whether an architecture has a high accuracy. However, two prior works have experimentally showed that truncating the path encoding does not decrease performance (Wei et al., 2020; White et al., 2019), and we gave even more experimental evidence in Section 4.

**Theorem 3.1.** Given  $10 \leq n \leq k \leq \frac{n(n-1)}{2}$ , and  $c > 3$ , for  $x > 2ec \cdot \frac{k}{n}$ ,  $b(k, x) > 1 - c^{-x+1}$ , and for  $x < \frac{1}{2ec} \cdot \frac{k}{n}$ ,  $b(k, x) < -2\frac{k}{2^n}$ .

To prove Theorem 3.1, we use the well-known bounds on binomial coefficients, e.g. (Stanica, 2001).

**Theorem A.2.** Given  $0 \leq \ell \leq n$ ,

$$\left(\frac{n}{\ell}\right)^\ell \leq \binom{n}{\ell} \leq \left(\frac{en}{\ell}\right)^\ell.$$

Now we give upper and lower bounds on  $a_{n,k,\ell}$  which will be used for the rest of the proofs. The next fact is similar to Lemma C.3 from BANANAS (White et al., 2019).

**Fact A.3.** Given  $n \leq k \leq \frac{n(n-1)}{2}$ , and  $0 < x < n$ , we have

$$\frac{2k}{n(n-1)} \left(\frac{2k(n-2)}{(\ell-1)n(n-1)}\right)^{\ell-1} \leq a_{n,k,\ell} \leq \frac{2k}{n(n-1)} \left(\frac{2ek(n-2)}{(\ell-1)n(n-1)}\right)^{\ell-1}$$

*Proof.* First, we have

$$a_{n,k,\ell} = \binom{n-2}{\ell-1} \left(\frac{2k}{n(n-1)}\right)^\ell$$

because on a path from node 1 to node  $n$  with length  $\ell$ , there are  $\binom{n-2}{\ell-1}$  choices of intermediate nodes from 1 to  $n$ . Once the nodes are chosen, we need all  $\ell$  edges between the nodes to exist, and each edge exists independently with probability  $\frac{2}{n(n-1)} \cdot k$ . Then we achieve the desired result by applying Theorem A.2.  $\square$

Now we prove the lower bound of Theorem 3.1.

**Lemma A.4.** Given  $n \leq k \leq \frac{n(n-1)}{2}$  and  $c > 2$ , for  $x > \frac{2eck}{n}$ ,  $b(k, x) > 1 - c^{-x+1}$ .

*Proof.* Given  $n \leq k \leq \frac{n(n-1)}{2}$  and  $x > \frac{2eck}{n}$ , we give a lower bound for  $\sum_{\ell=1}^x a_{n,k,\ell}$  and an upper bound for  $\sum_{\ell=x+1}^n a_{n,k,\ell}$ .

When  $\ell = 1$ , we have  $\binom{n-2}{\ell-1} = 1$ . Therefore,

$$\sum_{\ell=1}^x a_{n,k,\ell} \geq a_{n,k,1} = \frac{2k}{n(n-1)}.$$

Now we upper bound  $\sum_{\ell=x+1}^n a_{n,k,\ell}$ .

$$\begin{aligned}
 \sum_{\ell=x+1}^n a_{n,k,\ell} &\leq \sum_{\ell=x+1}^n \frac{2k}{n(n-1)} \left( \frac{2ek(n-2)}{(\ell-1)n(n-1)} \right)^{\ell-1} \\
 &= \frac{2k}{n(n-1)} \sum_{\ell=x+1}^n \left( \frac{2ek(n-2)}{(\ell-1)n(n-1)} \right)^{\ell-1} \\
 &\leq \frac{2k}{n(n-1)} \sum_{\ell=x+1}^n \left( \frac{1}{c} \right)^{\ell-1} \tag{A.1}
 \end{aligned}$$

$$\begin{aligned}
 &\leq \left( \frac{2k}{n(n-1)} \right) \left( \frac{1}{c} \right)^x \sum_{\ell=0}^{\infty} \left( \frac{1}{c} \right)^{\ell} \\
 &= \left( \frac{2k}{n(n-1)} \right) \left( \frac{1}{c} \right)^x \left( \frac{1}{1-\frac{1}{c}} \right) \\
 &= \left( \frac{2k}{n(n-1)} \right) \left( \frac{1}{c} \right)^x \left( \frac{c}{c-1} \right) \tag{A.2} \\
 &= \left( \frac{2k}{n(n-1)} \right) \left( \frac{1}{c} \right)^{x-1}
 \end{aligned}$$

In inequality A.1, we use the fact that for all  $\ell \geq x+1$ ,

$$\ell \geq x+1 > \frac{2eck}{n} + 1 \implies \frac{2ek(n-2)}{(\ell-1)n(n-1)} \leq \frac{2ek}{(\ell-1)n} \leq \frac{1}{c}$$

and in inequality A.1, we use the fact that  $c > 2$ .

Therefore, we have

$$\begin{aligned}
 b(k, x) &= \frac{\sum_{\ell=1}^x a_{n,k,\ell}}{\sum_{\ell=1}^n a_{n,k,\ell}} \\
 &= \frac{\sum_{\ell=1}^x a_{n,k,\ell}}{\sum_{\ell=1}^x a_{n,k,\ell} + \sum_{\ell=x+1}^n a_{n,k,\ell}} \\
 &\geq \frac{\frac{2k}{n(n-1)}}{\frac{2k}{n(n-1)} + \left( \frac{2k}{n(n-1)} \right) \left( \frac{1}{c} \right)^{x-1}} \tag{A.3} \\
 &= \frac{1}{1 + \left( \frac{1}{c} \right)^{x-1}} \\
 &\geq 1 - c^{-x+1}.
 \end{aligned}$$

In inequality A.3, we use the fact that for all  $0 \leq A, B, C$ , we know that  $A \geq B$  implies  $\frac{A}{A+C} \geq \frac{B}{B+C}$ . □

Now we prove the upper bound for Theorem 3.1.

**Lemma A.5.** Given  $10 \leq n \leq k \leq \frac{n(n-1)}{2}$  and  $c > 3$ , for  $x < \frac{k}{2ecn}$ ,  $b(k, x) < 2^{-\frac{k}{2n}}$ .

*Proof.* Given  $n \leq k \leq \frac{n(n-1)}{2}$  and  $x < \frac{k}{2ecn}$ , now we give an upper bound for  $\sum_{\ell=1}^x a_{n,k,\ell}$  and a lower bound for  $\sum_{\ell=1}^n a_{n,k,\ell}$ .

First we make the following claim. For all  $1 \leq \ell \leq x < \frac{k}{2ecn}$ , we have

$$\left( \frac{2ek(n-2)}{(\ell-1)n(n-1)} \right)^{\ell-1} < \left( \frac{4e^2c(n-2)}{n-1} \right)^{\frac{k}{2ecn}}. \quad (\text{A.4})$$

Now we prove the claim. In the following inequalities, we take log to have base 2.

$$\begin{aligned} \left( \frac{2ek(n-2)}{(\ell-1)n(n-1)} \right)^{\ell-1} &= 2^{(\ell-1) \log\left(\frac{2ek(n-2)}{(\ell-1)n(n-1)}\right)} \\ &= 2^{(\ell-1) \log \frac{1}{\ell-1} + (\ell-1) \log\left(\frac{2ek(n-2)}{n(n-1)}\right)} \\ &\leq 2^{\frac{k}{2ecn} \log\left(\frac{2ecn}{k}\right) + \frac{k}{2ecn} \log\left(\frac{k}{2ecn} \cdot \frac{4e^2c(n-2)}{n-1}\right)} \\ &= 2^{\frac{k}{2ecn} \left( \log\left(\frac{2ecn}{k}\right) + \log\left(\frac{k}{2ecn}\right) + \log\left(\frac{4e^2c(n-2)}{n-1}\right) \right)} \\ &= 2^{\frac{k}{2ecn} \log\left(\frac{4e^2c(n-2)}{n-1}\right)} \\ &= \left( \frac{4e^2c(n-2)}{n-1} \right)^{\frac{k}{2ecn}} \end{aligned} \quad (\text{A.5})$$

In inequality A.5, we use the fact that for any  $1 \leq A \leq B$ ,  $A \log\left(\frac{1}{A}\right) \leq B \log\left(\frac{1}{B}\right)$  (specifically, we used  $A = \ell - 1$  and  $B = \frac{k}{2ecn}$ ).

Now we have

$$\begin{aligned} \sum_{\ell=1}^x a_{n,k,\ell} &\leq \sum_{\ell=1}^x \frac{2k}{n(n-1)} \left( \frac{2ek(n-2)}{(\ell-1)n(n-1)} \right)^{\ell-1} \\ &= \frac{2k}{n(n-1)} \sum_{\ell=1}^x \left( \frac{2ek(n-2)}{(\ell-1)n(n-1)} \right)^{\ell-1} \\ &\leq \frac{2k}{n(n-1)} \sum_{\ell=1}^x \left( \frac{4e^2c(n-2)}{n-1} \right)^{\frac{k}{2ecn}} \\ &\leq \left( \frac{2k}{n(n-1)} \right) \cdot x \cdot \left( \frac{4e^2c(n-2)}{n-1} \right)^{\frac{k}{2ecn}}. \end{aligned}$$

Now we give the lower bound for the other summation which goes from  $\ell = 1$  to  $n$ . We lower bound the whole summation by a single term of the summation,  $\ell = \frac{k(n-2)}{n(n-1)} + 1$ . Recall that  $k \leq \frac{n(n-1)}{2}$ , which implies  $\frac{k}{n} \leq \frac{n-1}{2} < n$ , so  $\ell = \frac{k(n-2)}{n(n-1)} + 1$  is indeed between 1 and  $n$ .

$$\begin{aligned}
\sum_{\ell=1}^n a_{n,k,\ell} &= \sum_{\ell=1}^n \frac{2k}{n(n-1)} \left( \frac{2k(n-2)}{(\ell-1)n(n-1)} \right)^{\ell-1} \\
&\geq \frac{2k}{n(n-1)} \sum_{\ell=\frac{k(n-2)}{n(n-1)}+1}^{\frac{k(n-2)}{n(n-1)}+1} \left( \frac{2k(n-2)}{(\ell-1)n(n-1)} \right)^{\ell-1} \\
&= \frac{2k}{n(n-1)} (2)^{\frac{k(n-2)}{n(n-1)}}
\end{aligned}$$

Therefore,

$$\begin{aligned}
b(k, x) &= \frac{\sum_{\ell=1}^x a_{n,k,\ell}}{\sum_{\ell=1}^n a_{n,k,\ell}} \\
&\leq \frac{\left( \frac{2k}{n(n-1)} \right) \cdot x \cdot \left( \frac{4e^2 c(n-2)}{n-1} \right)^{\frac{k}{2ecn}}}{\frac{2k}{n(n-1)} (2)^{\frac{k(n-2)}{n(n-1)}}} \\
&\leq x \cdot (4e^2 c)^{\frac{k}{2ecn}} (2)^{-\frac{k(n-2)}{n(n-1)}} \\
&\leq 2^{\log\left(\frac{k}{2ecn}\right) + \frac{k}{2ecn} \log(4e^2 c) - \frac{k(n-2)}{n(n-1)}} \\
&= 2^{-\frac{k}{n} \left( \frac{n-2}{n-1} - \frac{1}{2ec} \log(4e^2 c) \right) - \log\left(\frac{k}{2ecn}\right)} \\
&\leq 2^{-\frac{k}{2n}}
\end{aligned}$$

The final inequality holds because  $n \geq 10$ ,  $c > 3$ , and  $\frac{k}{n} \geq 1$ . □

The proof of Theorem 3.1 follows immediately from combining Lemmas A.4 and A.5.

In Figure B.2, we plot  $b(k, x)$  for NASBench-101, which supports Theorem 3.1. Next, we may ask whether the one-hot adjacency matrix encoding can be truncated. However, even removing one bit from the adjacency matrix encoding can be very costly, because each single edge makes the difference between a path from the input node to the output node vs. no path from the input node to the output node. In the next theorem, we show that the probability of a random graph containing any individual edge is at least  $2k/(n(n-1))$ . Therefore, truncating the adjacency matrix encoding even by a single bit results in significant information loss. In the following theorem, let  $E_{n,k,r}$  denote the edge set of  $G_{n,k,r}$ . Given  $1 \leq z \leq n(n-1)/2$ , we slightly abuse notation by writing  $z \in E_{n,k,r}$  if the edge with index  $z$  in the adjacency matrix is in  $E_{n,k,r}$ .

**Theorem A.6.** Given  $n \leq k \leq \frac{n(n-1)}{2}$  and  $1 \leq z \leq n(n-1)/2$ , we have  $P(z \in E_{n,k,r}) > \frac{2k}{n(n-1)}$ .

**Proof.** Recall that *sample random architecture* adds each edge with probability  $2k/(n(n-1))$  and rejects in step (3) if there is no path from the input to the output. Define  $G'_{n,k,r}$

and *valid* as in the proof of Theorem 3.1 and let  $E'_{n,k,r}$  denote the edge set of  $G'_{n,k,r}$ . Then

$$\frac{P(G'_{n,k,r} \text{ is valid} \mid z \in E'_{n,k,r})}{P(G'_{n,k,r} \text{ is valid})} = \frac{P(z \in E'_{n,k,r} \mid G'_{n,k,r} \text{ is valid})}{P(z \in E'_{n,k,r})} > 1,$$

where the first equality comes from Bayes' theorem, and the inequality follows because there is a natural bijection  $\phi$  from graphs with  $z$  to graphs without  $z$  given by removing  $z$ , where  $G$  is valid if  $\phi(G)$  is valid but the reverse does not hold. Therefore,

$$\begin{aligned} P(z \in E_{n,k,r}) &= P(z \in E'_{n,k,r} \mid G'_{n,k,r} \text{ is valid}) = \frac{P(G'_{n,k,r} \text{ is valid} \mid z \in E'_{n,k,r})P(z \in E'_{n,k,r})}{P(G'_{n,k,r} \text{ is valid})} \\ &> P(z \in E'_{n,k,r}) = \frac{2k}{n(n-1)}. \end{aligned} \quad \square$$

Our theoretical results show that the path encoding can be heavily truncated, while the adjacency matrix cannot be truncated. We also verify this experimentally (Figure B.2).

## Appendix B. Details from Section 4 (Experiments)

In this section, we give more details from Section 4, and we give more experiments on NASBench-201. First we describe the algorithms used in the experiments in Section 4.

- Random Search consists of randomly choosing architectures and then training them, until the runtime budget is exceeded.
- Regularized evolution (Real et al., 2019) consists of maintaining a population of neural architectures. In each iteration, a subset is selected and the best architecture from the subset is mutated. The mutation replaces the oldest architecture from the population. We used a population size of 30. We also found that replacing the worst architecture (not the oldest) performed better, so we used this version.
- Local search (White et al., 2020) is a simple greedy algorithm that has only recently been applied to NAS. We use the simplest instantiation (often called the hill-climbing algorithm).
- Bayesian optimization (BO) is a strong method for zeroth order optimization. We use the ProBO (Neiswanger et al., 2019) implementation, which uses a Gaussian process kernel and expected improvement as the acquisition function.
- NASBOT (Kandasamy et al., 2018) is a BO-based NAS algorithm. It was not originally defined for cell-based search spaces, so we use a variant that works for cell-based spaces (White et al., 2019).
- BANANAS (White et al., 2019) is a BO-based method which uses a neural predictor model.

Table 1: Ability of neural predictor with different encodings to generalize beyond the search space.

Encoding	Validation error		Test error	
	Top 10 avg.	Top 1 avg.	Top 10 avg.	Top 1 avg.
Adjacency	<b>5.888</b>	<b>5.505</b>	<b>6.454</b>	<b>6.056</b>
Categorical Adjacency	7.589	6.191	8.155	7.086
Path	5.967	5.606	6.616	6.335
Truncated Path	6.082	5.644	6.712	6.452
Categorical Path	6.357	5.703	6.939	6.489
Truncated Categorical Path	6.339	5.895	6.918	6.766

### B.1 Experiments on NASBench-201

In this section, we give similar experiments to Figure 4.1, but with NASBench-201 instead of NASBench-101. Note that NASBench-201 is not as good for encoding experiments because every single architecture has the same graph structure - a clique of size 4. The only differences are the operations. Therefore, many encodings are functionally equivalent. For example, the one-hot, categorical, and continuous adjacency matrix encodings are all identical because the only difference is the way they encode the adjacency matrix. I.e., these encodings will all look like a set of operations, plus some adjacency matrix encoding that is the same for every architecture in the search space. The one-hot adjacency matrix encoding, path encoding, and truncated path encoding are all distinct from one another, so we run experiments with these encodings. See Figure B.1. We see largely the same trends as in NASBench-101 (Figure 4.1). Note that on the ImageNet-16-120 dataset, some algorithms such as NASBOT overfit to the training set, causing performance to decline over time.

**Outside search space experiment.** In the set of experiments above, we tested the effect of encodings on a neural predictor model by computing the mean absolute error between the predicted vs. actual errors on the test set, and also by evaluating the performance of BANANAS when changing the encoding of its neural predictor model. The latter experiment tests the predictor model’s ability to predict the *best* architectures, not just all architectures on average. We take this one step further and test the ability of the neural predictor to generalize beyond the search space on which it was trained. We set up the experiment as follows. We define the training search space as a subset of NASBench-101: architectures with at most 6 nodes and 7 edges. We define the disjoint test search space as architectures with 6 nodes and 7 to 9 edges. The neural predictor is trained on 1000 architectures and predicts the validation loss of the 5000 architectures from the test search space. We evaluate the losses of the ten architectures with the highest predicted validation loss. We run 200 trials for each encoding and average the results. See Table 1. The adjacency encoding performed the best. An explanation is that for the path encoding, there are features (paths) in architectures from the test set that do not exist in the training set. This is not the case for the adjacency encoding: all features (edges) from architectures in the test set have shown up in the training set.

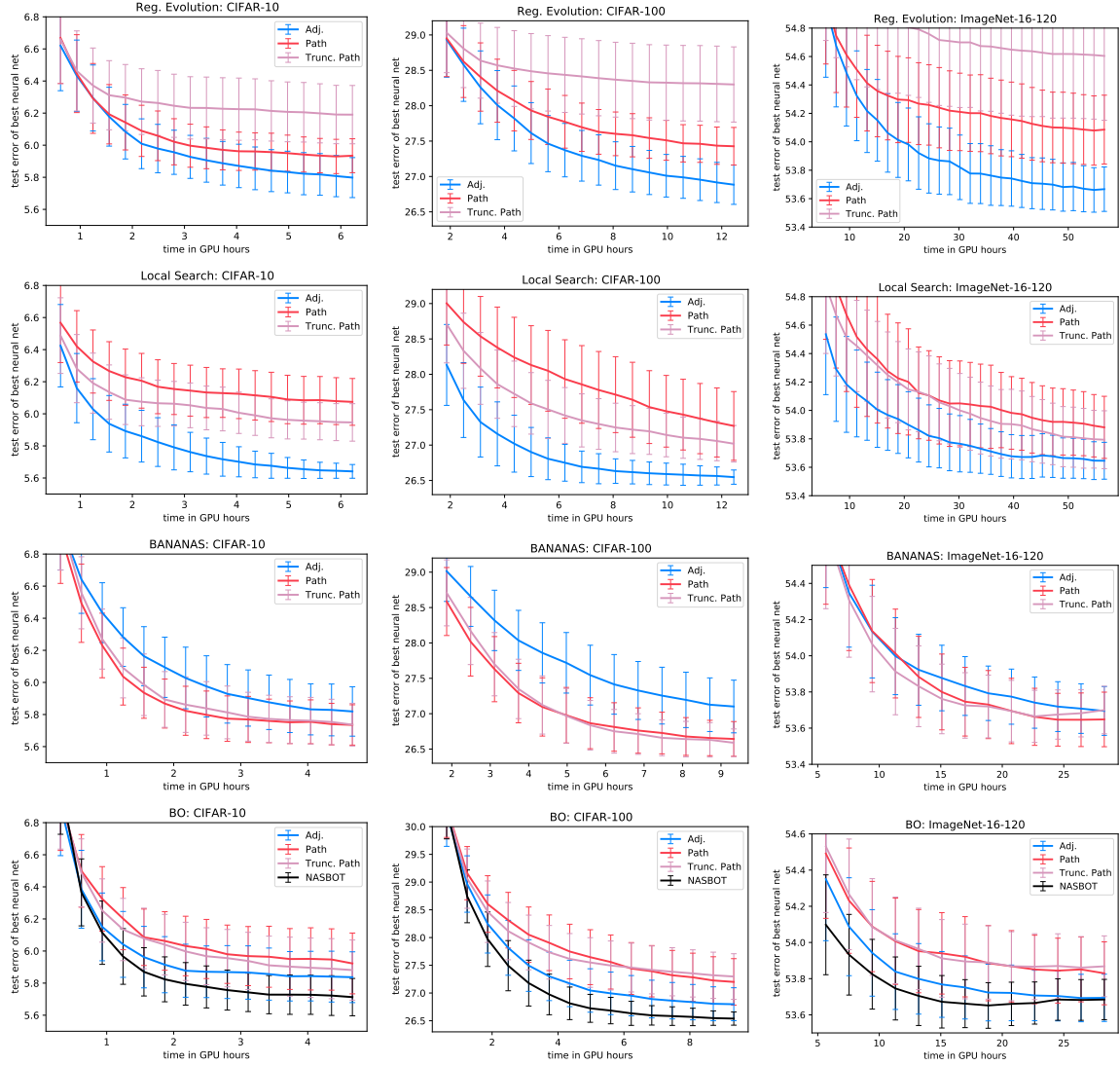


Figure B.1: Experiments on NASBench-201 with different encodings, keeping all but one subroutine fixed: *perturb architecture* (Reg. evolution (top row), local search (second row)), *train predictor model* (BANANAS (third row), Bayesian optimization (bottom row)).



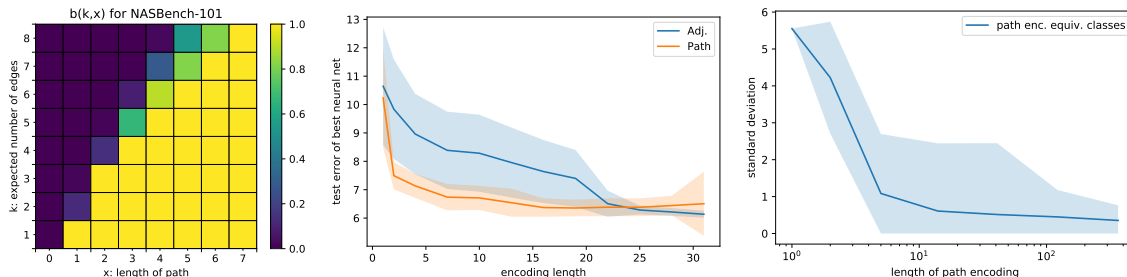


Figure B.2: Plot of  $b(k, x)$  on NASBench-101 (left), which is consistent with Theorem 3.1. Truncation of encodings for regularized evolution on NASBench-101 (middle). Average standard deviation of accuracies within each equivalence class defined by the path encoding at different levels of truncation on NASBench-101 (right).

**Equivalence class experiments.** Recall that the path encoding function  $e$  is not onto (see Figure A.2). In general, this is not desirable because information is lost when two architectures map to the same encoding. However, if the encoding function only maps architectures with similar accuracies to the same encoding, then the behavior is beneficial. On the NASBench-101 dataset, we compute the path encoding of all 423k architectures, and then we compute the average standard deviation of accuracies among architectures with the same encoding (i.e., we look at the standard deviations within the equivalence classes defined by the encoding). See Figure B.2. The result is an average standard deviation of 0.353%, compared to the 5.55% standard deviation over the entire set of architectures.

## B.2 Best practices for NAS

Many authors have called for improving the reproducibility and fairness in experimental comparisons in NAS research (Li and Talwalkar, 2019; Ying et al., 2019; Yang et al., 2020), which has led to the release of a NAS best practices checklist (Lindauer and Hutter, 2019). We address each section and we encourage future work to do the same.

- **Best practices for releasing code.** We released our code publicly. We used the NASBench-101 and NASBench-201 datasets, so questions about training pipeline, evaluation, and hyperparameters for the final evaluation do not apply.
- **Best practices for comparing NAS methods.** We made fair comparisons due to our use of NASBench-101 and NASBench-201. We did run ablation studies and ran random search. We performed 300 trials of each experiment on NASBench-101 and NASBench-201.
- **Best practices for reporting important details.** We used the hyperparameters straight from the open source repositories, with a few small exceptions listed earlier in this section.