
Towards Discovering Neural Architectures from Scratch

Simon Schrodi¹ Danny Stoll¹ Binxin Ru²
Rhea Sukthanker¹ Thomas Brox¹ Frank Hutter^{1,3}

¹University of Freiburg, ²University of Oxford, ³Bosch Center for Artificial Intelligence
{schrodi, stoll, sukthank, brox, fh}@cs.uni-freiburg.de, robin@robots.ox.ac.uk

Abstract

The discovery of neural architectures from scratch is the long-standing goal of Neural Architecture Search (NAS). Searching over a wide spectrum of neural architectures can facilitate the discovery of previously unconsidered but well-performing architectures. In this work, we take a large step towards discovering neural architectures from scratch by expressing architectures algebraically. This algebraic view leads to a more general method for designing search spaces, which allows us to compactly represent search spaces that are 100s of orders of magnitude larger than common spaces from the literature. Further, we propose a Bayesian Optimization strategy to efficiently search over such huge spaces, and demonstrate empirically that both our search space design and our search strategy can be superior to existing baselines. We open source our algebraic NAS approach and provide APIs for PyTorch and TensorFlow at https://github.com/automl/towards_nas_from_scratch.

1 Introduction

Neural Architecture Search (NAS), a field with over 1 000 papers in the last two years [1], is widely touted to automatically discover novel, well-performing architectural patterns. However, while state-of-the-art performance has already been demonstrated in hundreds of NAS papers (prominently, e.g., [2–4]), success in automatically finding truly novel architectural patterns has been very scarce [5, 6]. There is an accumulating amount of evidence that over-engineered, restrictive search spaces (e.g., cell-based ones) are major impediments for NAS to discover truly novel architectures [7–9].

In this work, we introduce a general formalism for the representation of hierarchical search spaces, allowing both for layer diversity and a flexible macro architecture, taking a large step towards discovering neural architectures from scratch. The key observation is that any neural architecture can be represented algebraically; e.g., two residual blocks followed by a fully-connected layer in a linear macro topology can be represented as the algebraic term

$$\text{Linear}(\text{Residual}(\text{conv}, \text{id}, \text{conv}), \text{Residual}(\text{conv}, \text{id}, \text{conv}), \text{fc}) \quad . \quad (1)$$

We build upon this observation and employ Context-Free Grammars (CFGs) to construct large spaces of such *algebraic architecture terms*. Although a particular search space is of course limited in its overall expressiveness, with this approach, we could effectively represent any neural architecture.

Due to the hierarchical structure of algebraic terms, the number of candidate neural architectures scales exponentially with the number of hierarchical levels, leading to search spaces 100s of orders of magnitudes larger than commonly used ones. To search in these huge spaces, we propose an efficient search strategy, Bayesian Optimization for Algebraic Neural Architecture Terms (BANAT), which leverages hierarchical information, capturing the topological patterns across the hierarchical levels, in its tailored kernel design.

Our contributions are as follows:

- We present an algebraic notion of NAS that views neural architectures as algebraic terms and propose to design these spaces with CFGs (Sec. 2).
- We propose BANAT, a Bayesian Optimization (BO) strategy that uses a tailored kernel to efficiently and effectively search over our huge search spaces (Sec. 3).
- We find that search spaces of algebraic architecture terms perform on par or better than common cell-based spaces on different datasets, show the superiority of BANAT over common baselines, and demonstrate the importance of incorporating hierarchical information (Sec. 5).

For discussion as well as limitations of our approach and broader impact statement, please refer to App. A or B, respectively.

2 Algebraic neural architecture search

In this section, we introduce an algebraic notion of Neural Architecture Search (NAS) by representing neural architectures with algebraic terms and propose to use CFGs to construct them.

Neural architectures as algebraic terms We introduce *algebraic architecture terms* as a string representation for neural architectures from a (term) algebra. Formally, an algebra (A, \mathcal{F}) consists of a non-empty set A (universe) and a set of operators $f: A^n \rightarrow A \in \mathcal{F}$ of different arities $n \geq 0$ [10]. In our case, A corresponds to the set of all (sub-)architectures and we distinguish between two types of operators: (i) nullary operators representing primitive computations (e.g., `conv()` or `fc()`) and (ii) k -ary operators with $k > 0$ representing topological operators (e.g., `Linear(·, ·, ·)` or `Residual(·, ·, ·)`). For sake of notational simplicity, we omit parenthesis for nullary operators (i.e., we write `conv`). Term algebras [11] are a special type of algebra which map an algebraic expression to its string representation. E.g., we can represent a neural architecture as the algebraic architecture term ω as shown in Eq. 1. Term algebras also allow for variables x_i that are set to terms themselves that can be re-used across a term. In our case, the *intermediate variables* x_i can therefore share patterns across the architecture, e.g., a shared cell. For example, we could define the intermediate variable x_1 to map to the residual block in ω from Eq. 1 as follows: $\omega' = \text{Linear}(x_1, x_1, \text{fc})$, $x_1 = \text{Residual}(\text{conv}, \text{id}, \text{conv})$.

Algebraic NAS Consequently, we formulate our algebraic view on NAS, where we search over algebraic architecture terms $\omega \in \Omega$ representing their associated architectures $\Phi(\omega)$ as follows: $\arg \min_{\omega \in \Omega} f(\Phi(\omega))$, where $f(\cdot)$ is an error measure that we seek to minimize, e.g., validation error. Thus, our search problem is to discover algebraic architecture terms ω , which can be part of very expressive search spaces Ω .

Constructing neural architecture terms with context-free grammars We propose to use *Context-Free Grammars (CFGs)* [12] since they can naturally generate (hierarchical) algebraic architecture terms. Compared to other search space designs, CFGs give us a formally grounded way to naturally and compactly define very expressive hierarchical search spaces (e.g., see Sec. 5). We can also unify popular search spaces from the literature with our general search space design in one framework (App. G) and CFGs provide a simple mechanism to evolve architectures while staying within the defined search space (Sec. 3).

Formally, a CFG $\mathcal{G} = \langle N, \Sigma, P, S \rangle$ consists of a finite set of nonterminals N , a finite set of terminals Σ (with $N \cap \Sigma = \emptyset$), a finite set of production rules $P = \{A \rightarrow \beta \mid A \in N, \beta \in (N \cup \Sigma)^*\}$, where the asterisk $*$ denotes the Kleene star operation [13], and a start symbol $S \in N$. For example, consider the following CFG in extended Backus-Naur form [14] (see App. D for background): $S ::= \text{Linear}(S, S, S) \mid \text{Residual}(S, S, S) \mid \text{conv} \mid \text{id} \mid \text{fc}$. Fig. 1 depicts the derivation of the algebraic architecture term from Eq. 1 and makes the connection to the associated architecture explicit. The set of all (potentially infinite) algebraic terms generated by a CFG \mathcal{G} is the language $L(\mathcal{G})$, which naturally forms our search space Ω . We can construct very expressive search spaces (e.g., see Sec. 5) as well as unify popular search spaces from the literature with our search space design (App. G). However, there is of course no single search space that can construct any neural architecture. In App. E we further augment the capabilities of CFGs so that we can incorporate constraints, foster regularity, and handle changes in the spatial resolution (and number of channels).

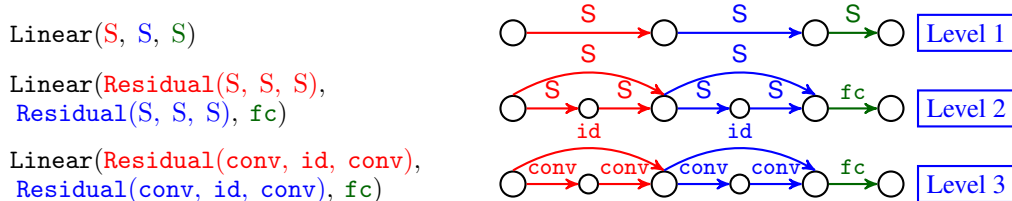


Figure 1: Derivations of algebraic terms (left) correspond to edge replacements [15–17] in the associated neural architecture (right). See App. C for the topological operators and primitive computations.

3 Bayesian Optimization for algebraic neural architecture search

We propose a BO strategy, Bayesian Optimization for Algebraic Neural Architecture Terms (BANAT), to search efficiently in the huge search spaces spanned by our algebraic architecture terms. Our search strategy, BANAT combines a Gaussian Process (GP) surrogate model with a tailored kernel that leverages the hierarchical structure of algebraic neural architecture terms (see below). For a more detailed explanation of the BO mechanism and our parallelization strategy, please refer to App. H.

Inspired by the state-of-the-art BO approach for NAS [18], we adopt the WL graph kernel [19] in a GP surrogate, modeling the architectures $\Phi(\omega_i)$ instead of its algebraic neural architecture terms ω_i . However, modeling solely based on the final architecture ignores the useful hierarchical information inherent in our algebraic representation. Moreover, the large size of the final architectures at the highest hierarchical level also makes it difficult to use a single WL kernel to capture the more global topological patterns. Note that our hierarchical construction of neural architectures can be viewed as a series of gradually unfolding architectures, with the final architecture containing only primitive computations (see Fig. 1). Thus, we propose a novel hierarchical kernel design which assigns a WL kernel to each hierarchy and combines them. To this end, we introduce fold operators F_l , that removes algebraic terms beyond the l -th hierarchical level. E.g., for the algebraic term ω in Eq. 1.

$$F_3(\omega) = \omega = \text{Linear}(\text{Residual}(\text{conv}, \text{id}, \text{conv}), \text{Residual}(\text{conv}, \text{id}, \text{conv}), \text{fc}), \quad (2)$$

$$F_2(\omega) = \text{Linear}(\text{Residual}, \text{Residual}, \text{fc}) \quad , \quad F_1(\omega) = \text{Linear} \quad .$$

Note the similarity to the derivations in Fig. 1. Furthermore note that, in practice, we also add the corresponding nonterminals to integrate information from our construction process. Consider two architectures $\Phi(\omega_i)$ and $\Phi(\omega_j)$ with algebraic architecture terms ω_i and ω_j , respectively, constructed over a hierarchy of L levels. We then define our hierarchical kernel as follows:

$$k_{hWL}(\omega_i, \omega_j) = \sum_{l=2}^L \lambda_l \cdot k_{WL}(\Phi(F_l(\omega_i)), \Phi(F_l(\omega_j))) \quad , \quad (3)$$

where the weights λ_l govern the importance of the learned graph information at different hierarchical levels and can be tuned (along with other hyperparameters of the GP) by maximizing the marginal likelihood. We omit $l = 1$ in the additive kernel as $F_1(\omega)$ does not contain any edge features which are required to apply our hierarchical WL kernel k_{hWL} . Our proposed kernel can efficiently capture the information in all algebraic term construction levels, which substantially improves its empirical performance on our search space as demonstrated in Sec. 5.

4 Related work

Previous approaches used, e.g., reinforcement learning [20, 21], evolution [22], gradient descent [23], or Bayesian Optimization (BO) [18, 24, 25]. To enable the effective use of BO on graph-like inputs for NAS, previous works have proposed to use a GP with specialized kernels [18, 24], encoding schemes [25, 26], or graph neural networks as surrogate model [27–29]. Different to prior works, we explicitly leverage the hierarchical construction of architectures for modeling. Most previous works focused on finding a shared cell [30] with a fixed macro architecture while only few works considered more expressive hierarchical [4, 9, 31–36] or grammar-based [37–54] approaches. Similar to the former, our formalism allows to design search spaces covering a general set of architecture design choices, but also permits the search for macro architectures with spatial resolution changes and

multiple branches. Different to the latter grammar-based works we construct entire architectures with spatial resolution changes across multiple branches, and propose techniques to incorporate constraints and foster regularity. Please refer to App. I for an extended discussion of related work.

5 Experiments

In this section, we investigate potential benefits of hierarchical search spaces and of our search strategy, BANAT. More specifically, we address the following questions: **Q1** Can hierarchical search spaces yield on par or superior architectures compared to cell-based search spaces with a limited number of evaluations?; **Q2** Can our search strategy BANAT improve performance over common baselines?; and **Q3** Does leveraging the hierarchical information improve performance? To answer these questions, we introduce a hierarchical search space (see below) based on the popular NAS-Bench-201 search space [55] and experimentally answer them in the affirmative. For evaluation, implementation, and training details as well as further results and analyses, please refer to App. J

Hierarchical NAS-Bench-201 search space We propose a hierarchical variant of the popular cell-based NAS-Bench-201 search space [55] by adding a hierarchical macro space (i.e., spatial resolution flow and wiring at the macro-level) and parameterizable convolutional blocks (i.e., choice of convolutions, activations, and normalizations):

$$\begin{aligned}
 D2 &::= \text{Linear3}(D1, D1, D0) \mid \text{Linear3}(D0, D1, D1) \mid \text{Linear4}(D1, D1, D0, D0) \\
 D1 &::= \text{Linear3}(C, C, D) \mid \text{Linear4}(C, C, C, D) \mid \text{Residual3}(C, C, D, D) \\
 D0 &::= \text{Linear3}(C, C, CL) \mid \text{Linear4}(C, C, C, CL) \mid \text{Residual3}(C, C, CL, CL) \\
 D &::= \text{Linear2}(CL, \text{down}) \mid \text{Linear3}(CL, CL, \text{down}) \mid \text{Residual2}(C, \text{down}, \text{down}) \\
 C &::= \text{Linear2}(CL, CL) \mid \text{Linear3}(CL, CL) \mid \text{Residual2}(CL, CL, CL) \\
 CL &::= \text{Cell}(OP, OP, OP, OP, OP, OP) \\
 OP &::= \text{zero} \mid \text{id} \mid \text{BLOCK} \mid \text{avg_pool} \\
 \text{BLOCK} &::= \text{Linear3}(\text{ACT}, \text{CONV}, \text{NORM}) \\
 \text{ACT} &::= \text{relu} \mid \text{hardswish} \mid \text{mish} \\
 \text{CONV} &::= \text{conv1x1} \mid \text{conv3x3} \mid \text{dconv3x3} \\
 \text{NORM} &::= \text{batch} \mid \text{instance} \mid \text{layer}
 \end{aligned} \tag{4}$$

See App. C for the terminal vocabulary of topological operators and primitive computations. The productions with the nonterminals $\{D2, D1, D0, D\}$ define the spatial resolution flow and together with $\{C\}$ define the macro architecture containing possibly multiple branches. The productions for $\{CL, OP\}$ construct the NAS-Bench-201 cell and $\{\text{BLOCK}, \text{ACT}, \text{CONV}, \text{NORM}\}$ parameterize the convolutional block. To ensure that we use the same distribution over the primitive computations as in NAS-Bench-201, we reweigh the sampling probabilities of the productions generated by the nonterminal OP, i.e., all production choices have sampling probability of 20 %, but BLOCK has 40 %. Note that we omit the stem (i.e., 3x3 convolution followed by batch normalization) and classifier (i.e., batch normalization followed by ReLU, global average pooling, and fully-connected layer) for simplicity. We implemented the merge operation as element-wise summation. Different to the cell-based NAS-Bench-201 search space, we exclude degenerated architectures by introducing a constraint that ensures that each subterm maps the input to the output (i.e., in the associated computational graph there is at least one path from source to sink).

Our search space consists of ca. 10^{446} algebraic architecture terms (please refer to App. F on how to compute the search space size), which is significantly larger than other popular search spaces from the literature. For comparison, the cell-based NAS-Bench-201 search space is just a minuscule subspace of size $10^{4.18}$, where we apply only the blue-colored production rules and replace the CL nonterminals with a placeholder terminal x_1 that will be substituted by the searched, shared cell.

Results Fig. 2 (top) compares the results of the cell-based and hierarchical search space design using our search strategy BANAT. Results with BANAT are on par on CIFAR-10/100, superior on ImageNet-16-120, and clearly superior on CIFARtile and AddNIST (answering **Q1**). We emphasize that the NAS community has engineered the cell-based search space to achieve strong performance on those popular image classification datasets for over a decade, making it unsurprising that our improvements are much larger for the novel datasets. Yet, our best found architecture on ImageNet-16-120 from the hierarchical search space also achieves an excellent test error of 52.78 % with only 0.626 MB parameters (App. J.4); this is superior to the state-of-the-art method Shapley-NAS (i.e.,

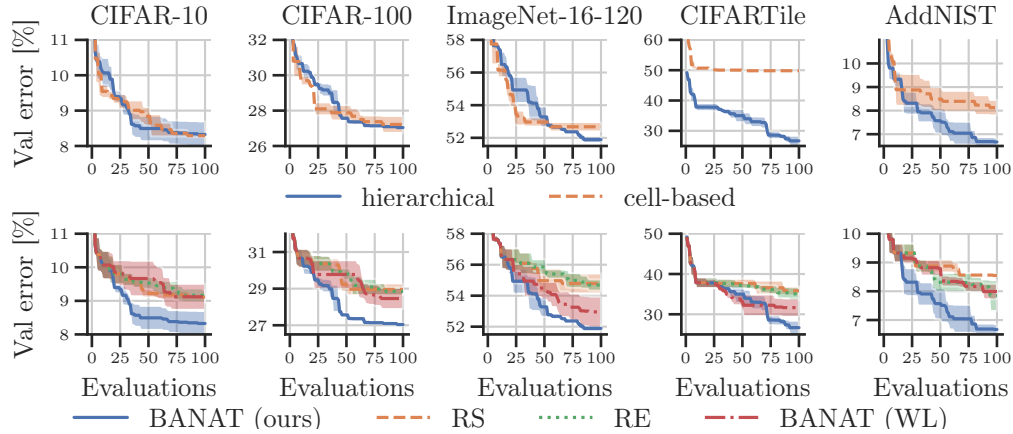


Figure 2: We compare cell-based vs. hierarchical search space using BANAT (top) and compare it to common baselines (bottom). We plot mean and ± 1 standard error.

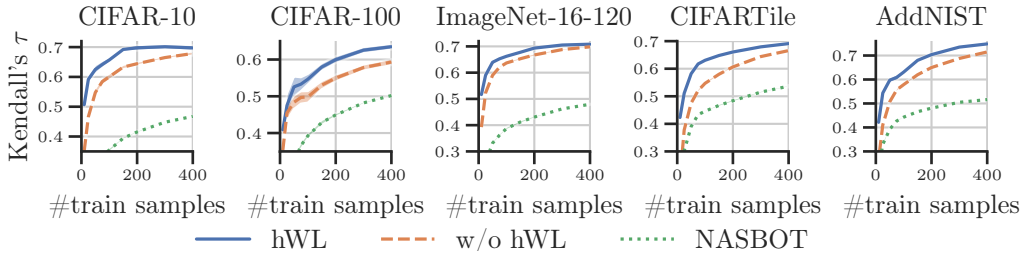


Figure 3: Mean Kendall's τ rank correlation with ± 1 standard error achieved by a GP with our hierarchical WL kernel (hWL), (standard) WL kernel (WL), and NASBOT [24].

53.15%) [56] and on par with the *optimal* architecture of the cell-based NAS-Bench-201 search space (i.e., 52.69% with 0.866 MB). Fig. 2 (bottom) shows that BANAT is superior to common baselines (answering **Q2**) and that leveraging hierarchical information clearly improves performance (answering **Q3**). The analysis in Fig. 3 shows that incorporating hierarchical information improves modeling, especially on smaller amounts of training data; this provides (further answering **Q3**).

6 Conclusion

We introduced very expressive search spaces of algebraic architecture terms constructed with CFGs. To efficiently search over the huge search spaces, we proposed BANAT. Our experiments indicate that both our search space design and algorithm can yield strong performance over existing baselines.

Acknowledgments and Disclosure of Funding

This research was funded by the Deutsche Forschungsgemeinschaft (DFG, German Research Foundation) under grant number 417962828, and the Bundesministerium für Umwelt, Naturschutz, nukleare Sicherheit und Verbraucherschutz (BMUV, German Federal Ministry for the Environment, Nature Conservation, Nuclear Safety and Consumer Protection) based on a resolution of the German Bundestag (67KI2029A). Robert Bosch GmbH is acknowledged for financial support. We gratefully acknowledge support by the European Research Council (ERC) Consolidator Grant “Deep Learning 2.0” (grant no. 101045765). Funded by the European Union. Views and opinions expressed are however those of the author(s) only and do not necessarily reflect those of the European Union or the ERC. Neither the European Union nor the ERC can be held responsible for them.



References

- [1] Difan Deng and Marius Lindauer. Literature on Neural Architecture Search. <https://www.automl.org/automl/literature-on-neural-architecture-search/>, 2022. [Online; accessed 25-September-2022].
- [2] Mingxing Tan and Quoc V. Le. EfficientNet: Rethinking Model Scaling for Convolutional Neural Networks. In *International Conference on Machine Learning*, 2019.
- [3] Mingxing Tan and Quoc V. Le. EfficientNetV2: Smaller Models and Faster Training. In *International Conference on Machine Learning*, 2021.
- [4] Chenxi Liu, Liang-Chieh Chen, Florian Schroff, Hartwig Adam, Wei Hua, Alan L. Yuille, and Li Fei-Fei. Auto-DeepLab: Hierarchical Neural Architecture Search for Semantic Image Segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.
- [5] Prajit Ramachandran, Barret Zoph, and Quoc V. Le. Searching for Activation Functions. *arXiv*, 2017.
- [6] Hanxiao Liu, Andy Brock, Karen Simonyan, and Quoc V. Le. Evolving Normalization-Activation Layers. *Advances in Neural Information Processing Systems*, 2020.
- [7] Antoine Yang, Pedro M. Esperança, and Fabio M. Carlucci. NAS evaluation is frustratingly hard. In *International Conference on Learning Representations*, 2020.
- [8] Saining Xie, Alexander Kirillov, Ross Girshick, and Kaiming He. Exploring Randomly Wired Neural Networks for Image Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.
- [9] Robin Ru, Pedro Esperança, and Fabio Maria Carlucci. Neural architecture generator optimization. In *Advances in Neural Information Processing Systems*, 2020.
- [10] Garrett Birkhoff. On the structure of abstract algebras. *Mathematical Proceedings of the Cambridge Philosophical Society*, 1935.
- [11] Franz Baader and Tobias Nipkow. *Term rewriting and all that*. Cambridge university press, 1999.
- [12] Noam Chomsky. Three models for the description of language. *IRE Transactions on information theory*, 1956.
- [13] Stephen C. Kleene et al. Representation of events in nerve nets and finite automata. *Automata studies*, 1956.
- [14] John Warner Backus. The syntax and semantics of the proposed international algebraic language of the zurich acm-gamm conference. In *International Conference on Information Processing*, 1959.
- [15] Annegret Habel and Hans-Jörg Kreowski. On context-free graph languages generated by edge replacement. In *Graph-Grammars and Their Application to Computer Science*, 1983.
- [16] Annegret Habel and Hans-Jörg Kreowski. Characteristics of graph languages generated by edge replacement. *Theoretical Computer Science*, 1987.
- [17] Frank Drewes, Hans-Jörg Kreowski, and Annegret Habel. Hyperedge replacement graph grammars. In *Handbook Of Graph Grammars And Computing By Graph Transformation: Volume 1: Foundations*. World Scientific, 1997.
- [18] Binxin Ru, Xingchen Wan, Xiaowen Dong, and Michael Osborne. Interpretable Neural Architecture Search via Bayesian Optimisation with Weisfeiler-Lehman Kernels. In *International Conference on Learning Representations*, 2021.
- [19] Nino Shervashidze, Pascal Schweitzer, Erik Jan Van Leeuwen, Kurt Mehlhorn, and Karsten M. Borgwardt. Weisfeiler-Lehman Graph Kernels. *Journal of Machine Learning Research*, 2011.
- [20] Barret Zoph and Quoc V. Le. Neural Architecture Search with Reinforcement Learning. In *International Conference on Learning Representations*, 2017.
- [21] Hieu Pham, Melody Guan, Barret Zoph, Quoc V. Le, and Jeff Dean. Efficient Neural Architecture Search via Parameters Sharing. In *International Conference on Machine Learning*, 2018.
- [22] Esteban Real, Sherry Moore, Andrew Selle, Saurabh Saxena, Yutaka Leon Suematsu, Jie Tan, Quoc V. Le, and Alexey Kurakin. Large-Scale Evolution of Image Classifiers. In *International Conference on Machine Learning*, 2017.

- [23] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: Differentiable Architecture Search. In *International Conference on Learning Representations*, 2019.
- [24] Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabás Póczos, and Eric P. Xing. Neural Architecture Search with Bayesian Optimisation and Optimal Transport. In *Advances in Neural Information Processing Systems*, 2018.
- [25] Colin White, Willie Neiswanger, and Yash Savani. BANANAS: Bayesian Optimization with Neural Architectures for Neural Architecture Search. In *Proceedings of the National Conference on Artificial Intelligence*, 2021.
- [26] Chris Ying, Aaron Klein, Esteban Real, Eric Christiansen, Kevin Murphy, and Frank Hutter. NAS-Bench-101: Towards Reproducible Neural Architecture Search. In *International Conference on Machine Learning*, 2019.
- [27] Lizheng Ma, Jiaxu Cui, and Bo Yang. Deep Neural Architecture Search with Deep Graph Bayesian Optimization. In *IEEE/WIC/ACM International Conference on Web Intelligence*. IEEE Computer Society Press, 2019.
- [28] Han Shi, Renjie Pi, Hang Xu, Zhenguo Li, James Kwok, and Tong Zhang. Bridging the Gap between Sample-based and One-shot Neural Architecture Search with BONAS. In *Advances in Neural Information Processing Systems*, 2020.
- [29] Chris Zhang, Mengye Ren, and Raquel Urtasun. Graph HyperNetworks for Neural Architecture Search. In *International Conference on Learning Representations*, 2019.
- [30] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V. Le. Learning Transferable Architectures for Scalable Image Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018.
- [31] Hanxiao Liu, Karen Simonyan, Oriol Vinyals, Chrisantha Fernando, and Koray Kavukcuoglu. Hierarchical Representations for Efficient Architecture Search. In *International Conference on Learning Representations*, 2018.
- [32] Haokui Zhang, Ying Li, Hao Chen, and Chunhua Shen. Memory-Efficient Hierarchical Neural Architecture Search for Image Denoising. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020.
- [33] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V. Le. MnasNet: Platform-Aware Neural Architecture Search for Mobile. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.
- [34] Risto Miikkulainen, Jason Liang, Elliot Meyerson, Aditya Rawal, Daniel Fink, Olivier Francon, Bala Raju, Hormoz Shahrzad, Arshak Navruzyan, Nigel Duffy, et al. *Evolving Deep Neural Networks*. In *Artificial Intelligence in the Age of Neural Networks and Brain Computing*. Elsevier, 2019.
- [35] Renato Negrinho, Matthew Gormley, Geoffrey J. Gordon, Darshan Patil, Nghia Le, and Daniel Ferreira. Towards modular and programmable architecture search. *Advances in Neural Information Processing Systems*, 2019.
- [36] David So, Wojciech Mańke, Hanxiao Liu, Zihang Dai, Noam Shazeer, and Quoc V. Le. Searching for Efficient Transformers for Language Modeling. In *Advances in Neural Information Processing Systems*, 2021.
- [37] Hiroaki Kitano. Designing neural networks using genetic algorithms with graph generation system. *Complex systems*, 1990.
- [38] Egbert J. W. Boers, Herman Kuiper, Bart L. M. Happel, and Ida G. Sprinkhuizen-Kuyper. Biological Metaphors In Designing Modular Artificial Neural Networks. In *International Conference on Artificial Neural Networks*, 1993.
- [39] Frederic Gruau. *Neural Network Synthesis Using Cellular Encoding And The Genetic Algorithm*. PhD thesis, Laboratoire de l’Informatique du Parallélisme, Ecole Normale Supérieure de Lyon, 1994.
- [40] Sean Luke and Lee Spector. Evolving Graphs and Networks with Edge Encoding: Preliminary Report. In *Late-breaking Papers of the Genetic Programming 96 conference*, 1996.
- [41] Edwin D. De Jong and Jordan B. Pollack. Utilizing Bias to Evolve Recurrent Neural Networks. In *International Joint Conference on Neural Networks*, 2001.

- [42] Han Cai, Jiacheng Yang, Weinan Zhang, Song Han, and Yong Yu. Path-Level Network Transformation for Efficient Architecture Search. In *International Conference on Machine Learning*, 2018.
- [43] Martin H. Luerksen and David M. W. Powers. *On the Artificial Evolution of Neural Graph Grammars*. University of New South Wales, 2003.
- [44] Martin H. Luerksen. Graph Grammar Encoding and Evolution of Automata Networks. In *Australasian Conference on Computer Science*, 2005.
- [45] Jean-Baptiste Mouret and Stéphane Doncieux. MENNAG: a modular, regular and hierarchical encoding for neural-networks based on attribute grammars. *Evolutionary Intelligence*, 2008.
- [46] Christian Jacob and Jan Rehder. Evolution of neural net architectures by a hierarchical grammar-based genetic system. In *Artificial Neural Nets and Genetic Algorithms*, 1993.
- [47] Jorge Couchet, Daniel Manrique, and Luis Porras. Grammar-Guided Neural Architecture Evolution. In *International Work-Conference on the Interplay Between Natural and Artificial Computation*, 2007.
- [48] Fardin Ahmadizar, Khabat Soltanian, Fardin AkhlaghianTab, and Ioannis Tsoulos. Artificial neural network development by means of a novel combination of grammatical evolution and genetic algorithm. *Engineering Applications of Artificial Intelligence*, 2015.
- [49] Qadeer Ahmad, Atif Rafiq, Muhammad Adil Raja, and Noman Javed. Evolving MIMO multi-layered artificial neural networks using grammatical evolution. In *Proceedings of the ACM Symposium on Applied Computing*, 2019.
- [50] Filipe Assunção, Nuno Lourenço, Penousal Machado, and Bernardete Ribeiro. Automatic generation of neural networks with structured grammatical evolution. In *IEEE Congress on Evolutionary Computation*, 2017.
- [51] Filipe Assunção, Nuno Lourenço, Penousal Machado, and Bernardete Ribeiro. DENSER: deep evolutionary network structured representation. *Genetic Programming and Evolvable Machines*, 2019.
- [52] Ricardo H. R. Lima, Aurora T. R. Pozo, and Roberto Santana. Automatic Design of Convolutional Neural Networks using Grammatical Evolution. In *Brazilian Conference on Intelligent Systems*, 2019.
- [53] Víctor de la Fuente Castillo, Alberto Díaz-Álvarez, Miguel-Ángel Manso-Callejo, and Francisco Seradilla Garcia. Grammar Guided Genetic Programming for Network Architecture Search and Road Detection on Aerial Orthophotography. *Applied Sciences*, 2020.
- [54] Xilai Li, Xi Song, and Tianfu Wu. AOGNets: Compositional grammatical architectures for deep learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019.
- [55] Xuanyi Dong and Yi Yang. NAS-Bench-201: Extending the Scope of Reproducible Neural Architecture Search. In *International Conference on Learning Representations*, 2020.
- [56] Han Xiao, Ziwei Wang, Zheng Zhu, Jie Zhou, and Jiwen Lu. Shapley-NAS: Discovering Operation Contribution for Neural Architecture Search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022.
- [57] William Ogden. A helpful result for proving inherent ambiguity. *Mathematical systems theory*, 1968.
- [58] Kenneth O. Stanley and Risto Miikkulainen. Evolving Neural Networks through Augmenting Topologies. *Evolutionary Computation*, 2002.
- [59] Jie Hu, Li Shen, and Gang Sun. Squeeze-and-Excitation Networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018.
- [60] Mark Sandler, Andrew Howard, Menglong Zhu, Andrey Zhmoginov, and Liang-Chieh Chen. MobileNetV2: Inverted Residuals and Linear Bottlenecks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2018.
- [61] Duncan J. Watts and Steven H. Strogatz. Collective dynamics of ‘small-world’ networks. *Nature*, 1998.
- [62] Paul Erdős, Alfréd Rényi, et al. On the evolution of random graphs. *Publ. Math. Inst. Hung. Acad. Sci.*, 1960.
- [63] Eric Brochu, Vlad M. Cora, and Nando De Freitas. A Tutorial on Bayesian Optimization of Expensive Cost Functions, with Application to Active User Modeling and Hierarchical Reinforcement Learning. *arXiv*, 2010.

- [64] Bobak Shahriari, Kevin Swersky, Ziyu Wang, Ryan P. Adams, and Nando De Freitas. Taking the Human Out of the Loop: A Review of Bayesian Optimization. *Proceedings of the IEEE*, 2015.
- [65] Jonas Mockus, Vytautas Tiesis, and Antanas Zilinskas. The Application of Bayesian Methods for Seeking the Extremum. *Towards Global Optimization*, 1978.
- [66] David Ginsbourger, Rodolphe Le Riche, and Laurent Carraro. Kriging is well-suited to parallelize optimization. In *Computational Intelligence in Expensive Optimization Problems*. Springer, 2010.
- [67] Robert I. McKay, Nguyen Xuan Hoai, Peter A. Whigham, Yin Shan, and Michael O’neill. Grammar-based genetic programming: a survey. *Genetic Programming and Evolvable Machines*, 2010.
- [68] Henry Moss, David Leslie, Daniel Beck, Javier González, and Paul Rayson. BOSS: Bayesian Optimization over String Spaces. In *Advances in Neural Information Processing Systems*, 2020.
- [69] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural Architecture Search: A Survey. *Journal of Machine Learning Research*, 2019.
- [70] Kenneth O Stanley, Jeff Clune, Joel Lehman, and Risto Miikkulainen. Designing neural networks through neuroevolution. *Nature Machine Intelligence*, 2019.
- [71] Marcin Andrychowicz, Misha Denil, Sergio Gomez, Matthew W Hoffman, David Pfau, Tom Schaul, Brendan Shillingford, and Nando De Freitas. Learning to learn by gradient descent by gradient descent. *Advances in Neural Information Processing Systems*, 2016.
- [72] Ke Li and Jitendra Malik. Learning to Optimize. *International Conference on Learning Representations*, 2017.
- [73] Yutian Chen, Matthew W Hoffman, Sergio Gómez Colmenarejo, Misha Denil, Timothy P. Lillicrap, Matt Botvinick, and Nando Freitas. Learning to Learn without Gradient Descent by Gradient Descent. In *International Conference on Machine Learning*, 2017.
- [74] Tianlong Chen, Xiaohan Chen, Wuyang Chen, Howard Heaton, Jialin Liu, Zhangyang Wang, and Wotao Yin. Learning to Optimize: A Primer and A Benchmark. *Journal of Machine Learning Research*, 2022.
- [75] Irwan Bello, Barret Zoph, Vijay Vasudevan, and Quoc V. Le. Neural Optimizer Search with Reinforcement Learning. In *International Conference on Machine Learning*, 2017.
- [76] Ruo Chen Wang, Yuanhao Xiong, Minhao Cheng, and Cho-Jui Hsieh. Efficient Non-Parametric Optimizer Search for Diverse Tasks. *Advances in Neural Information Processing Systems*, 2022.
- [77] Esteban Real, Chen Liang, David So, and Quoc V. Le. AutoML-Zero: Evolving Machine Learning Algorithms From Scratch. In *International Conference on Machine Learning*, 2020.
- [78] Xiangning Chen, Chen Liang, Da Huang, Esteban Real, Yao Liu, Kaiyuan Wang, Cho-Jui Hsieh, Yifeng Lu, and Quoc V. Le. Evolved Optimizer for Vision. In *First Conference on Automated Machine Learning (Late-Breaking Workshop)*, 2022.
- [79] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An Image is Worth 16x16 Words: Transformers for Image Recognition at Scale. In *International Conference on Learning Representations*, 2021.
- [80] Alexander L Gaunt, Marc Brockschmidt, Nate Kushman, and Daniel Tarlow. Differentiable programs with neural libraries. In *International Conference on Machine Learning*, 2017.
- [81] Lazar Valkov, Dipak Chaudhari, Akash Srivastava, Charles Sutton, and Swarat Chaudhuri. Houdini: Lifelong Learning as Program Synthesis. *Advances in Neural Information Processing Systems*, 2018.
- [82] Ameesh Shah, Eric Zhan, Jennifer Sun, Abhinav Verma, Yisong Yue, and Swarat Chaudhuri. Learning Differentiable Programs with Admissible Neural Heuristics. *Advances in Neural Information Processing Systems*, 2020.
- [83] Guofeng Cui and He Zhu. Differentiable Synthesis of Program Architectures. In *Advances in Neural Information Processing Systems*, 2021.
- [84] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V. Le. Regularized Evolution for Image Classifier Architecture Search. In *Proceedings of the National Conference on Artificial Intelligence*, 2019.

- [85] Ilya Loshchilov and Frank Hutter. Decoupled Weight Decay Regularization. In *International Conference on Learning Representations*, 2019.
- [86] Rob Geada, Stephen McGough, Amir A. Abarghouei, Isabelle Guyon, and Sébastien Treguer. CVPR-NAS-Datasets (codebase for the “CVPR-NAS 2021 Unseen Data Track”). <https://github.com/RobGeada/cvpr-nas-datasets>, 2021.
- [87] Alex Krizhevsky, Geoffrey Hinton, et al. Learning Multiple Layers of Features from Tiny Images, 2009.
- [88] Patryk Chrabaszcz, Ilya Loshchilov, and Frank Hutter. A Downsampled Variant of ImageNet as an Alternative to the CIFAR Datasets. *arXiv*, 2017.
- [89] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep Residual Learning for Image Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2016.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? **[Yes]** Our abstract and introduction accurately reflect our paper’s contribution.
 - (b) Did you describe the limitations of your work? **[Yes]** See Sec. A.
 - (c) Did you discuss any potential negative societal impacts of your work? **[Yes]** See App. B.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? **[Yes]** We discuss the ethics guidelines in App. B.
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? **[N/A]** We did not include theoretical results.
 - (b) Did you include complete proofs of all theoretical results? **[N/A]** We did not include theoretical results.
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? **[Yes]** https://github.com/automl/towards_nas_from_scratch.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? **[Yes]** We report all training details in App. J.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? **[Yes]** We plotted ± 1 standard error for our experiments in Sec. 5 over 3 or 20 trials for the search or surrogate experiments, respectively.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? **[Yes]** We report compute time and type of resource in App. J.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? **[Yes]** We cite the creators of all datasets and search strategies used.
 - (b) Did you mention the license of the assets? **[Yes]** See App. J.3.
 - (c) Did you include any new assets either in the supplemental material or as a URL? **[N/A]** We do not include new assets.
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? **[N/A]** We did not use or release any datasets with personal data.
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? **[N/A]** The data does not contain personal information.
5. If you used crowdsourcing or conducted research with human subjects...

- (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A] We did not run experiments with human subjects.
- (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A] We did not run experiments with human subjects.
- (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A] We did not run experiments with human subjects.

A Discussion and limitations

While our grammar-based construction mechanism is a powerful mechanism to construct huge hierarchical search space, we can not construct any architecture with our grammar-based construction approach (Sec. 2 and App. E) since we are limited to context-free languages; e.g., architectures of the type $\{a^n b^n c^n | n \in \mathbb{N}_{>0}\}$ cannot be generated by CFGs (this can be proven using Odgen’s lemma [57]). Further, due to the discrete nature of CFGs we can not integrate continuous design choices, e.g., dropout probability. Furthermore, our grammar-based mechanism does not (generally) support simple scalability of discovered neural architectures (e.g., repetition of building blocks) without special consideration in the search space design. Nevertheless, our search spaces still significantly increase the expressiveness, including the ability to represent common search spaces from the literature (see App. G for how we can represent the search spaces of DARTS, Auto-Deeplab, the hierarchical cell search space of Liu et al. [31], the Mobile-net search space, and the hierarchical random graph generator search space), as well as allowing search for entire neural architectures, e.g., based around the popular NAS-Bench-201 search space (Sec. 5). Thus, our search space design can facilitate the discovery of novel well-performing neural architectures in those huge search spaces of algebraic architecture terms.

However, there is an inherent trade-off between the expressiveness and the difficulty of search. The much greater expressiveness facilitates search in a richer set of architectures that may include better architectures than in more restrictive search spaces, which however need not exist. Besides that, the (potential) existence of such a well-performing architecture does not lead a search strategy inevitably discovering it, even with large amounts of compute available. Note that the trade-off manifests itself also in the acquisition function optimization of our search strategy BANAT.

In addition, a well-performing neural architecture may not work with current training protocols and hyperparameters due to interaction effects, i.e., training protocols and hyperparameters may be over-optimized for specific types of neural architectures. To overcome this limitation, one could consider a joint optimization of neural architectures, training protocols, and hyperparameters. However, this further fuels the trade-off between expressiveness and the difficulty of search.

B Broader impact

NAS has immense potential to facilitate systematic, automated discovery of high-performing (novel) architecture designs. However, the restrictive cell-based search spaces most commonly used in NAS render it impossible to discover truly novel neural architectures. With our general formalism based on algebraic terms, we hope to provide fertile foundation towards discovering high-performing and efficient architectures; potentially from scratch. However, search in such huge search spaces is expensive, particularly in the context of the ongoing detrimental climate crisis. While on the one hand, the discovered neural architectures, like other AI technologies, could potentially be exploited to have a negative societal impact; on the other hand, our work could also lead to advances across scientific disciplines like healthcare and chemistry.

C From terminals to primitive computations and topological operators

Tab. 1 and Fig. 4 describe the primitive computations and topological operators used throughout our experiments in Sec. 5 and App. J, respectively. Note that by adding more primitive computations and/or topological operators we could construct even more expressive search spaces.

Table 1: Primitive computations. "Name" corresponds to the string terminals in our CFGs and "Function" is the associated implementation of the primitive computation in pseudocode. The subscripts g , k , s , and p are abbreviations for groups, kernel size, strides, and padding, respectively. During assembly of neural architectures $\Phi(\omega)$, we replace string terminals with the associated primitive computation.

Name	Function
avg_pool	$\text{AvgPool}_{k=3,s=1,p=1}(x)$
batch	$\text{BN}(x)$
conv1x1	$\text{Conv}_{k=1,s=1,p=0}(x)$
conv3x3	$\text{Conv}_{k=3,s=1,p=1}(x)$
dconv3x3	$\text{Conv}_{g=C,k=3,s=1,p=1}(x)$
down	$\text{conv3x3}_{s=1}(\text{conv3x3}_{s=2}(x)) + \text{Conv}_{k=1,s=1}(\text{AvgPool}_{k=2,s=2}(x))$
hardswish	$\text{Hardswish}(x)$
id	$\text{Identity}(x)$
instance	$\text{IN}(x)$
layer	$\text{LN}(x)$
mish	$\text{Mish}(x)$
relu	$\text{ReLU}(x)$
zero	$\text{Zeros}(x)$

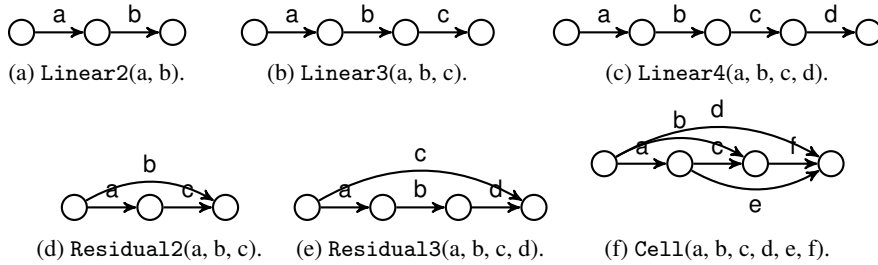


Figure 4: Topological operators. Each subfigure makes the connection between the topological operator and associated computational graph explicit, i.e., the arguments of the graph operators (a, b, ...) are mapped to the respective edges in the computational graph.

D Extended Backus-Naur form

The (extended) Backus-Naur form [14] is a meta-language to describe the syntax of CFGs. We use meta-rules of the form $S ::= \alpha$ where $S \in N$ is a nonterminal and $\alpha \in (N \cup \Sigma)^*$ is a string of nonterminals and/or terminals. We denote nonterminals in UPPER CASE, terminals corresponding to topological operators in Initial upper case/teletype, and terminals corresponding to primitive computations in lower case/teletype, e.g., $S ::= \text{Residual}(S, S, \text{id})$. To compactly express production rules with the same left-hand side nonterminal, we use the vertical bar $|$ to indicate a choice of production rules with the same left-hand side, e.g., $S ::= \text{Linear}(S, S, S) | \text{Residual}(S, S, \text{id}) | \text{conv}$.

E Augmenting the capabilities of context-free grammars

Below we augment the capabilities of CFGs and leverage properties of context-free languages so that we can incorporate constraints, foster regularity, and handle changes in the spatial resolution (and number of channels).

Constraints In many search space designs, we want to adhere to some constraints, e.g., to limit the number of nodes or to ensure that for all architectures in the search space there exists at least one path from the input to the output. We can simply do so by allowing only the application of production rules which guarantee compliance to such constraints. For example, to ensure that there is at least

one path from the input to the output, it is sufficient to ensure that each derivation connects its input to the output due to the recursive nature of CFGs. Note that this makes CFGs context-sensitive w.r.t. those constraints.

To implement the above constraint "only consider valid neural architectures", we note that our search space design only creates neural architectures where neither the spatial resolution nor the channels can be mismatched; please refer to Sec. 2 for details. Thus, the only way a neural architecture can become invalid is through zero operations, which could remove edges from the computational graph and possibly disassociate the input from the output. Since we recursively assemble neural architectures, it is sufficient to ensure that the derived algebraic architecture term (i.e., the associated computational graph) is compliant with the constraint, i.e., there is at least one path from input to output. Thus, during sampling (and similarly during evolution), we modify the current production rule choices when an application of the zero operation would disassociate the input from the output.

Fostering regularity through substitution To implement intermediate variables x_i (Sec. 2) we leverage that context-free languages are closed under substitution: we map terminals, representing the intermediate variables x_i , from one language to algebraic terms of other languages, e.g., a shared cell. For example, we can split a CFG \mathcal{G} , constructing entire algebraic architecture terms, into the CFGs \mathcal{G}_{macro} and \mathcal{G}_{cell} for the macro- or cell-level, respectively. Further, we add a single (or multiple) intermediate terminal(s) x_1 to \mathcal{G}_{macro} which maps to an algebraic term $\omega_1 \in L(\mathcal{G}_{cell})$, e.g., the searchable cell. Thus, we effectively search over the macro-level as well as a single, shared cell. Note that by using a fixed macro architecture (i.e., $|L(\mathcal{G}_{macro})| = 1$), we can represent cell-based search spaces, e.g., NAS-Bench-201 [55], while also being able to represent more expressive search spaces (e.g., see Sec. 5). More generally, we could extend this by adding further intermediate terminals which map to other languages $L(\mathcal{G}_j)$, or by adding intermediate terminals to \mathcal{G}_2 which map to languages $L(\mathcal{G}_{j \neq 1})$. In this way, we can effectively foster regularity.

Representing common architecture patterns for object recognition Neural architectures for object recognition commonly build a hierarchy of features that are gradually downsampled, e.g., by pooling operations. However, previous works in NAS were either limited to a fixed macro architecture [30], only allowed for linear macro architectures [4], or required post-sampling testing for resolution mismatches [9, 58]. While this produced impressive performance on popular benchmarks [2–4], it is still an open research question whether a different type of macro architecture (e.g., one with multiple branches) could yield even better performance.

To accommodate flexible macro architectures, we propose to overload the nonterminals. In particular, the nonterminals indicate how often we apply downsampling operations in the subsequent derivations of the nonterminal. Consider the production rule $D2 \rightarrow \text{Residual}(D1, D2, D1)$, where D_i with $i \in \{1, 2\}$ are a nonterminals which indicate that i downsampling operations have to be applied in their subsequent derivations. That is, in both paths of the residual the input features will be downsampled twice and, consequently, the merging paths will have the same spatial resolution. Thereby, this mechanism distributes the downsampling operations recursively across the architecture. For the channels, we adopted the common design to double the number of channels whenever we halve the spatial resolution in our experiments. Note that we could also handle a varying number of channels by using, e.g., depthwise concatenation as merge operation.

F Search space size

In this section, we show how to efficiently compute the size of our search spaces constructed by CFGs. There are two cases to consider: (i) a CFG contains cycles (i.e., part of the derivation can be repeated infinitely many times), yielding an open-ended, infinite search space; and (ii) a CFG contains no cycles, yielding in a finite search space whose size we can compute.

Consider a production $A \rightarrow \text{Residual}(B, B, B)$ where Residual is a terminal, and A and B are nonterminals with $B \rightarrow \text{conv} \mid \text{id}$. Consequently, there are $2^3 = 8$ possible instances of the residual block. If we add another production choice for the nonterminal A , e.g., $A \rightarrow \text{Linear}(B, B, B)$, we would have $2^3 + 2^3 = 16$ possible instances. Further, adding a production $C \rightarrow \text{Linear}(A, A, A)$ would yield a search space size of $(2^3 + 2^3)^3 = 4096$.

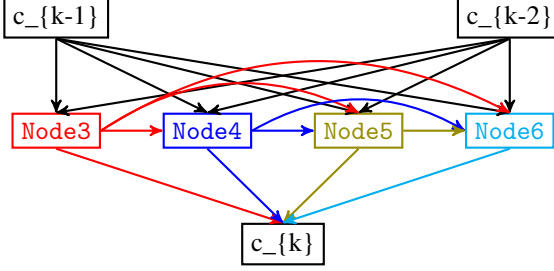


Figure 5: Visualization of the Darts topological operator.

More generally, we introduce the function P_A that returns the set of productions for nonterminal $A \in N$, and the function $\mu : P \rightarrow N$ that returns all the nonterminals for a production $p \in P$. We can then recursively compute the size of the search space as follows:

$$f(A) = \sum_{p \in P_A} \begin{cases} 1 & , \mu(p) = \emptyset, \\ \prod_{A' \in \mu(p)} f(A') & , otherwise . \end{cases} \quad (5)$$

When a CFG contains some constraint, we ensure to only account for valid architectures (i.e., compliant with the constraints) by ignoring productions which would lead to invalid architectures.

G Common search spaces from the literature

In Sec. 5, we demonstrated how to construct the popular NAS-Bench-201 search space within our algebraic search space design, and below we show how to reconstruct the following popular search spaces: DARTS search space [23], Auto-DeepLab search space [4], hierarchical cell search space [31], Mobile-net search space [33], and hierarchical random graph generator search space [9]. For implementation details we refer to the respective works.

DARTS search space

The DARTS search space [23] consists of a fixed macro architecture and a cell, i.e., a seven node directed acyclic graph (Darts; see Fig. 5 for the topological operator). We omit the fixed macro architecture from our search space design for simplicity. Each cell receives the feature maps from the two preceding cells as input and outputs a single feature map. All intermediate nodes (i.e., Node3, Node4, Node5, and Node6) is computed based on all of its predecessors. Thus, we can define the DARTS search space as follows:

$$\begin{aligned} \text{DARTS} &::= \text{Darts}(\text{NODE3}, \text{NODE4}, \text{NODE5}, \text{NODE6}) \\ \text{NODE3} &::= \text{Node3}(\text{OP}, \text{OP}) \\ \text{NODE4} &::= \text{Node4}(\text{OP}, \text{OP}, \text{OP}) \\ \text{NODE5} &::= \text{Node5}(\text{OP}, \text{OP}, \text{OP}, \text{OP}) \\ \text{NODE6} &::= \text{Node6}(\text{OP}, \text{OP}, \text{OP}, \text{OP}, \text{OP}) \\ \text{OP} &::= \text{sep_conv_3x3} \mid \text{sep_conv_5x5} \mid \text{dil_conv_3x3} \mid \text{dil_conv_5x5} \\ &\quad \mid \text{max_pool} \mid \text{avg_pool} \mid \text{id} \mid \text{zero} , \end{aligned} \quad (6)$$

where the topological operator Node3 receives two inputs, applies the operations separately on them, and sums them up. Similarly, Node4, Node5, and Node6 apply their operations separately to the given inputs and sum them up. The topological operator Darts feeds the corresponding feature maps into each of those topological operators and finally concatenates all intermediate feature maps.

Auto-DeepLab search space

Auto-DeepLab [4] combines a cell-level with a network-level search space to search for segmentation networks, where the cell is shared across the searched macro architecture, i.e., a twelve step (linear) path across different spatial resolutions. The cell-level design is adopted from Liu et al. [23] and,

thus, we can re-use the CFG from Eq. 6. For the network-level, we introduce a constraint that ensures that the path is of length twelve, i.e., we ensure exactly twelve derivations in our CFG. Further, we overload the nonterminals so that they correspond to the respective spatial resolution level, e.g., D4 indicates that the original input is downsampled by a factor of four; please refer to Sec. 2 for details on overloading nonterminals. For the sake of simplicity, we omit the first two layers and atrous spatial pyramid poolings as they are fixed, and hence define the network-level search space as follows:

$$\begin{aligned}
D4 & ::= \text{Same}(\text{CELL}, D4) \mid \text{Down}(\text{CELL}, D8) \\
D8 & ::= \text{Up}(\text{CELL}, D4) \mid \text{Same}(\text{CELL}, D8) \mid \text{Down}(\text{CELL}, D16) \\
D16 & ::= \text{Up}(\text{CELL}, D8) \mid \text{Same}(\text{CELL}, D16) \mid \text{Down}(\text{CELL}, D32) \\
D32 & ::= \text{Up}(\text{CELL}, D16) \mid \text{Same}(\text{CELL}, D32) \quad ,
\end{aligned} \tag{7}$$

where the topological operators Up, Same, and Down upsample/halve, do not change/do not change, or downsample/double the spatial resolution/channels, respectively. The placeholder variable CELL maps to the shared DARTS cell from the language generated by the CFG from Eq. 6.

Hierarchical cell search space

The hierarchical cell search space [31] consists of a fixed (linear) macro architecture and a hierarchically assembled cell with three levels which is shared across the macro architecture. Thus, we can omit the fixed macro architecture from our search space design for simplicity. Their first, second, and third hierarchical levels correspond to the primitive computations (i.e., `id`, `max_pool`, `avg_pool`, `sep_conv`, `depth_conv`, `conv`, `zero`), six densely connected four node directed acyclic graphs (DAG4), and a densely connected five node directed acyclic graph (DAG5), respectively. The `zero` operation could lead to directed acyclic graphs which have fewer nodes. Therefore, we introduce a constraint enforcing that there are always four (level 2) or five (level 3) nodes for every directed acyclic graph. Further, since a densely connected five node directed acyclic graph has ten edges, we need to introduce placeholder variables (i.e., M1, ..., M6) to enforce that only six (possibly) different four node directed acyclic graphs are used, and consequently define a CFG for the third level

$$\text{LEVEL3} ::= \text{DAG5}(\underbrace{\text{LEVEL2}, \dots, \text{LEVEL2}}_{\times 10}) \tag{8}$$

$$\text{LEVEL2} ::= M1 \mid M2 \mid M3 \mid M4 \mid M5 \mid M6 \mid \text{zero} \quad ,$$

mapping the placeholder variables M1, ..., M6 to the six lower-level motifs constructed by the first and second hierarchical level

$$\text{LEVEL2} ::= \text{DAG4}(\underbrace{\text{LEVEL1}, \dots, \text{LEVEL1}}_{\times 6}) \tag{9}$$

$$\text{LEVEL1} ::= \text{id} \mid \text{max_pool} \mid \text{avg_pool} \mid \text{sep_conv} \mid \text{depth_conv} \mid \text{conv} \mid \text{zero} \quad .$$

Mobile-net search space

Factorized hierarchical search spaces, e.g., the Mobile-net search space [33], allow for layer diversity. They factorize a (fixed) macro architecture – often based on an already well-performing reference architecture – into separate blocks (e.g., cells). For the sake of simplicity, we assume here a three sequential blocks (Block) architecture (Linear). In each of those blocks, we search for the convolution operations (CONV), kernel sizes (KSIZE), squeeze-and-excitation ratio (SERATIO) [59], skip connections (SKIP), number of output channels (FSIZE), and number of layers per block (#LAYERS), where the latter two are discretized using a reference architecture, e.g., MobileNetV2 [60]. Consequently, we can express this search space as follows:

$$\begin{aligned}
\text{MACRO} & ::= \text{Linear}(\text{BLOCK}, \text{BLOCK}, \text{BLOCK}) \\
\text{BLOCK} & ::= \text{Block}(\text{CONV}, \text{KSIZE}, \text{SERATIO}, \text{SKIP}, \text{FSIZE}, \text{\#LAYERS}) \\
\text{CONV} & ::= \text{conv} \mid \text{dconv} \mid \text{mbconv} \\
\text{KSIZE} & ::= 3 \mid 5 \\
\text{SERATIO} & ::= 0 \mid 0.25 \\
\text{SKIP} & ::= \text{pooling} \mid \text{id_residual} \mid \text{no_skip} \\
\text{FSIZE} & ::= 0.75 \mid 1.0 \mid 1.25 \\
\text{\#LAYERS} & ::= -1 \mid 0 \mid 1 \quad ,
\end{aligned} \tag{10}$$

where `conv`, `dconv` and `mbconv` correspond to convolution, depthwise convolution, and mobile inverted bottleneck convolution [60], respectively.

Algorithm 1 Bayesian Optimization algorithm [63].

Input: Initial observed data \mathcal{D}_t , a black-box objective function f , total number of BO iterations T
Output: The best recommendation about the global optimizer \mathbf{x}^*
for $t = 1, \dots, T$ **do**
 Select the next \mathbf{x}_{t+1} by maximizing acquisition function $\alpha(\mathbf{x}|\mathcal{D}_t)$
 Evaluate the objective function at $f_{t+1} = f(\mathbf{x}_{t+1})$
 $\mathcal{D}_{t+1} \leftarrow \mathcal{D}_t \cup (\mathbf{x}_{t+1}, f_{t+1})$
 Update the surrogate model with \mathcal{D}_{t+1}
end for

Hierarchical random graph generator search space

The hierarchical random graph generator search space [9] consists of three hierarchical levels of random graph generators (i.e., Watts-Strogatz [61] and Erdős-Rényi [62]). We denote with Watts-Strogatz_i the random graph generated by the Watts-Strogatz model with i nodes. Thus, we can represent the search space as follows:

$$\begin{aligned}
\text{TOP} & ::= \text{Watts-Strogatz}_3(K, \text{Pt})(\text{MID}, \text{MID}, \text{MID}) \mid \dots \\
& \quad \mid \text{Watts-Strogatz}_{10}(K, \text{Pt})(\underbrace{\text{MID}, \dots, \text{MID}}_{\times 10}) \\
\text{MID} & ::= \text{Erdős-Rényi}_1(\text{Pm})(\text{BOT}) \mid \dots \\
& \quad \mid \text{Erdős-Rényi}_{10}(\text{Pm})(\underbrace{\text{BOT}, \dots, \text{BOT}}_{\times 10}) \\
\text{BOT} & ::= \text{Watts-Strogatz}_3(K, \text{Pb})(\text{NODE}, \text{NODE}, \text{NODE}) \mid \dots \\
& \quad \mid \text{Watts-Strogatz}_{10}(K, \text{Pb})(\underbrace{\text{NODE} \dots, \text{NODE}}_{\times 10}) \\
K & ::= 2 \mid 3 \mid 4 \mid 5 \quad ,
\end{aligned} \tag{11}$$

where each terminal Pt, Pm, and Pb maps to a continuous number in $[0.1, 0.9]^1$ and the placeholder variable NODE maps to a primitive computation, e.g., separable convolution. Note that we omit other hyperparameters, such as stage ratio, channel ratio etc., for simplicity.

H More details on the search strategy

In this section, we provide more details and examples for our search strategy Bayesian Optimization for Algebraic Neural Architecture Terms (BANAT) presented in Sec. 3.

H.1 Bayesian Optimization

Bayesian Optimization (BO) is a powerful family of search techniques for finding the global optimum of a black-box objective problem. It is particularly useful when the objective is expensive to evaluate and thus sample efficiency is highly important [63].

To minimize a black-box objective problem with BO, we first need to build a probabilistic surrogate to model the objective based on the observed data so far. Based on the surrogate model, we design an acquisition function to evaluate the utility of potential candidate points by trading off exploitation (where the posterior mean of the surrogate model is low) and exploration (where the posterior variance of the surrogate model is high). The next candidate points to evaluate is then selected by maximizing the acquisition function [64]. The general procedures of BO is summarized in Algorithm 1.

We adopted the widely used acquisition function, expected improvement (EI) [65], in our BO strategy. EI evaluates the expected amount of improvement of a candidate point \mathbf{x} over the minimal value f' observed so far. Specifically, denote the improvement function as $I(\mathbf{x}) = \max(0, f' - f(\mathbf{x}))$, the EI

¹Theoretically, this is not possible with CFGs. However, we can extend the notion of substitution by substituting a string representation of a Python (float) variable for the placeholder variables Pt, Pm, and Pb.

Algorithm 2 Kriging Believer algorithm to select one batch of points.

Input: Observation data \mathcal{D}_t , batch size b
Output: The batch points $\mathcal{B}_{t+1} = \{\mathbf{x}_{t+1}^{(1)}, \dots, \mathbf{x}_{t+1}^{(b)}\}$
 $\tilde{\mathcal{D}}_t = \mathcal{D}_t \cup \tilde{\mathcal{D}}_p$
for $j = 1, \dots, b$ **do**
 Select the next $\mathbf{x}_{t+1}^{(j)}$ by maximizing acquisition function $\alpha(\mathbf{x}|\tilde{\mathcal{D}}_t)$
 Compute the predictive posterior mean $\mu(\mathbf{x}_{t+1}^{(j)}|\tilde{\mathcal{D}}_t)$
 $\tilde{\mathcal{D}}_t \leftarrow \tilde{\mathcal{D}}_t \cup (\mathbf{x}_{t+1}^{(j)}, \mu(\mathbf{x}_{t+1}^{(j)}|\tilde{\mathcal{D}}_t))$
end for

Algorithm 3 Weisfeiler-Lehman subtree kernel computation [19].

Input: Graphs G_1, G_2 , maximum iterations H
Output: Kernel function value between the graphs
Initialize the feature vectors $\phi(G_1) = \phi_0(G_1), \phi(G_2) = \phi_0(G_2)$ with the respective counts of original node labels (i.e., the $h = 0$ WL features)
for $h = 1, \dots, H$ **do**
 Assign a multiset $M_h(v) = \{l_{h-1}(u) | u \in \mathcal{N}(v)\}$ to each node $v \in G$, where l_{h-1} is the node label function of the $h - 1$ -th WL iteration and \mathcal{N} is the node neighbor function
 Sort elements in multiset $M_h(v)$ and concatenate them to string $s_h(v)$
 Compress each string $s_h(v)$ using the hash function f s.t. $f(s_h(v)) = f(s_h(w)) \iff s_h(v) = s_h(w)$
 Add l_{h-1} as prefix for $s_h(v)$
 Concatenate the WL features $\phi_h(G_1), \phi_h(G_2)$ with the respective counts of the *new* labels:
 $\phi(G_1) = [\phi(G_1), \phi_h(G_1)], \phi(G_2) = [\phi(G_2), \phi_h(G_2)]$
 Set $l_h(v) := f(s_h(v)) \forall v \in G$
end for
Compute inner product $k = \langle \phi_h(G_1), \phi_h(G_2) \rangle$ between WL features $\phi_h(G_1), \phi_h(G_2)$ in RKHS \mathcal{H}

acquisition function has the form

$$\begin{aligned} \alpha_{EI}(\mathbf{x}|\mathcal{D}_t) &= \mathbb{E}[I(\mathbf{x})|\mathcal{D}_t] = \int_{-\infty}^{f'} (f' - f) \mathcal{N}(f; \mu(\mathbf{x}|\mathcal{D}_t), \sigma^2(\mathbf{x}|\mathcal{D}_t)) df \\ &= (f' - f) \Phi(f'; \mu(\mathbf{x}|\mathcal{D}_t), \sigma^2(\mathbf{x}|\mathcal{D}_t)) + \sigma^2(\mathbf{x}|\mathcal{D}_t) \phi(f'; \mu(\mathbf{x}|\mathcal{D}_t), \sigma^2(\mathbf{x}|\mathcal{D}_t)) \quad , \end{aligned}$$

where $\mu(\mathbf{x}|\mathcal{D}_t)$ and $\sigma^2(\mathbf{x}|\mathcal{D}_t)$ are the mean and variance of the predictive posterior distribution at a candidate point \mathbf{x} , and $\phi(\cdot)$ and $\Phi(\cdot)$ denote the PDF and CDF of the standard normal distribution, respectively.

To make use of ample distributed computing resource, we adopted Kriging Believer [66] which uses the predictive posterior of the surrogate model to assign hallucinated function values $\{\tilde{f}_p\}_{p \in \{1, \dots, P\}}$ to the P candidate points with pending evaluations $\{\tilde{\mathbf{x}}_p\}_{p \in \{1, \dots, P\}}$ and perform next BO recommendation in the batch by pseudo-augmenting the observation data with $\tilde{\mathcal{D}}_p = \{(\tilde{\mathbf{x}}_p, \tilde{f}_p)\}_{p \in \{1, \dots, P\}}$, namely $\tilde{\mathcal{D}}_t = \mathcal{D}_t \cup \tilde{\mathcal{D}}_p$. The algorithm of Kriging Believer at one BO iteration to select a batch of recommended candidate points is summarized in Algorithm 2.

H.2 Weisfeiler-Lehman kernel

Inspired by Ru et al. [18], we adopted the Weisfeiler-Lehman (WL) graph kernel [19] in the GP surrogate model to handle the graph nature of neural architectures. The basic idea of the WL kernel is to first compare node labels, and then iteratively aggregate labels of neighboring nodes, compress them into a new label and compare them. Algorithm 3 summarizes the WL kernel procedure.

Ru et al. [18] identified three reasons for using the WL kernel: (1) it is able to compare labeled and directed graphs of different sizes, (2) it is expressive, and (3) it is relatively efficient and scalable. Our search space design can afford a diverse spectrum of neural architectures with very heterogeneous

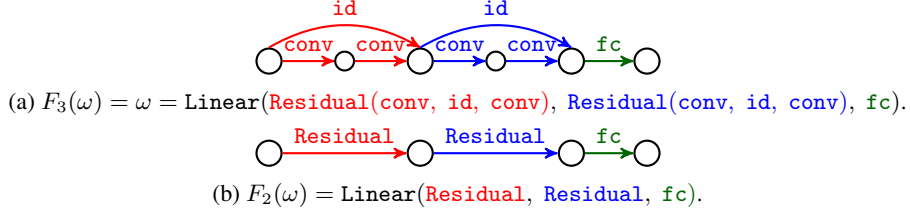


Figure 6: Labeled graphs $\Phi(F_2)$ and $\Phi(F_3)$ of the folds F_2 and F_3 .

topological structure. Therefore, reason (1) is a very important property of the WL kernel to account for the diversity of neural architectures. Moreover, if we allow many hierarchical levels, we can construct very large neural architectures. Therefore, reasons (2) and (3) are essential for accurate and fast modeling. However, neural architectures in our search spaces may be significantly larger, which makes it difficult for a single WL kernel to capture the more global topological patterns. Moreover, modeling solely based on the final neural architecture ignores the useful macro-level information from earlier hierarchical levels. In our experiments (Sec. 5 and J), we have found stronger neural architectures by incorporating the hierarchical information in the kernel design, which provides experimental support for above arguments.

However, modeling solely based on the (standard) WL graph kernel neglects the useful hierarchical information from our assembly process. Moreover, the large size of neural architectures make it still challenging to capture the more global topological patterns. We therefore propose to use hierarchical information through a hierarchy of WL graph kernels that take into account the different granularities of the architectures and combine them in a weighted sum. To obtain the different granularities, we use the fold operators F_l that removes algebraic terms beyond the l -th hierarchical level. Thereby, we obtain the folds

$$F_3(\omega) = \omega = \text{Linear}(\text{Residual}(\text{conv}, \text{id}, \text{conv}), \text{Residual}(\text{conv}, \text{id}, \text{conv}), \text{fc}), \quad (12)$$

$$F_2(\omega) = \text{Linear}(\text{Residual}, \text{Residual}, \text{fc}) \quad , \quad F_1(\omega) = \text{Linear} \quad ,$$

for the algebraic architecture term ω . Note that we ignore the first fold since it does not represent a labeled DAG. Fig. 6 visualizes the labeled graphs $\Phi(F_2)$ and $\Phi(F_3)$ of the folds F_2 or F_3 , respectively. These graphs can be fed into (standard) WL graph kernels. Therefore, we can construct a hierarchy of WL graph kernels k_{WL} as follows:

$$k_{hWL}(\omega_i, \omega_j) = \sum_{l=2}^L \lambda_l \cdot k_{WL}(\Phi(F_l(\omega_i)), \Phi(F_l(\omega_j))) \quad , \quad (13)$$

where ω_i and ω_j are two algebraic architecture terms. Note that λ_l govern the importance of the learned graph information across the hierarchical levels and can be optimized through the marginal likelihood.

H.3 Evolutionary operations in detail

For the evolutionary operations, we adopted ideas from grammar-based genetic programming [67, 68]. In the following, we will show how these evolutionary operations manipulate algebraic terms, e.g.,

$$\text{Linear}(\text{Residual}(\text{conv}, \text{id}, \text{conv}), \text{Residual}(\text{conv}, \text{id}, \text{conv}), \text{fc}) \quad , \quad (14)$$

from the search space

$$S ::= \text{Linear}(S, S, S) \mid \text{Residual}(S, S, S) \mid \text{conv} \mid \text{id} \mid \text{fc} \quad , \quad (15)$$

to generate evolved algebraic terms. Fig. 1 shows how we can derive the algebraic term in Eq. 14 from the search space in Eq. 15. For mutation operations, we first randomly pick a subterm of the algebraic term, e.g., $\text{Residual}(\text{conv}, \text{id}, \text{conv})$. Then, we randomly sample a new subterm with the same nonterminal symbol S as start symbol, e.g., $\text{Linear}(\text{conv}, \text{id}, \text{fc})$, and replace the previous subterm, yielding

$$\text{Linear}(\text{Linear}(\text{conv}, \text{id}, \text{fc}), \text{Residual}(\text{conv}, \text{id}, \text{conv}), \text{fc}) \quad . \quad (16)$$

For (self-)crossover operations, we swap two subterms, e.g., `Residual(conv, id, conv)` and `Residual(conv, id, conv)` with the same nonterminal `S` as start symbol, yielding

$$\text{Linear}(\text{Residual}(\text{conv}, \text{id}, \text{conv}), \text{Residual}(\text{conv}, \text{id}, \text{conv}), \text{fc}) \quad . \quad (17)$$

Note that unlike the commonly used crossover operation, which uses two parents, self-crossover has only one parent. In future work, we could also add a *self-copy* operation that copies a subterm to another part of the algebraic term, explicitly regularizing diversity and thus potentially speeding up the search.

I Extended related work

Neural Architecture Search Neural Architecture Search (NAS) aims to automatically discover architectural patterns (or even entire architectures) [69]. Previous approaches, e.g., used reinforcement learning [20, 21], evolution [22], gradient descent [23], or Bayesian Optimization (BO) [18, 24, 25]. To enable the effective use of BO on graph-like inputs for NAS, previous works have proposed to use a GP with specialized kernels [18, 24], encoding schemes [25, 26], or graph neural networks as surrogate model [27–29]. Different to prior works, we explicitly leverage the hierarchical construction of architectures for modeling.

Searching for novel architectural patterns Previous works mostly focused on finding a shared cell [30] with a fixed macro architecture while only few works considered more expressive hierarchical search spaces [4, 31, 33]. The latter works considered hierarchical assembly [31], combination of a cell- and network-level search space [4, 32], evolution of network topologies [34], factorization of the search space [33], parameterization of a hierarchy of random graph generators [9], a formal language over computational graphs [35], or a hierarchical construction of TensorFlow programs [36]. Similarly, our formalism allows to design search spaces covering a general set of architecture design choices, but also permits the search for macro architectures with spatial resolution changes and multiple branches. We also handle spatial resolution changes without requiring post-hoc testing or resizing of the feature maps unlike prior works [34, 58, 70]. Other works proposed approaches based on string rewriting systems [37, 38], cellular (or tree-structured) encoding schemes [39–42], hyperedge replacement graph grammars [43, 44], attribute grammars [45], CFGs [46–53], or And-Or-grammars [54]. Different to these prior works, we construct entire architectures with spatial resolution changes across multiple branches, and propose techniques to incorporate constraints and foster regularity.

Related work beyond NAS Optimizer search is a closely related field to NAS, where we automatically search for an optimizer (i.e., an update function for the weights) instead of an architecture. Initial works used learnable parametric or non-parametric optimizers. While the former approaches [71–74] have poor scalability and generality, the latter works overcome those limitations. Bello et al. [75] searched for an instantiation of hand-crafted patterns via reinforcement learning, while Wang et al. [76] proposed a tree-structured search space² and searched for optimizers via a modified Monte Carlo sampling approach. AutoML-Zero [77] took an even more general approach by searching over entire machine learning algorithms, including optimizers, from a generic search space built from basic mathematical operations with an evolutionary algorithm. Chen et al. [78] used RE to discover optimizers from a generic search space (inspired by AutoML-Zero) for training vision transformers [79].

Complementary to the above, there is recent interest in automatically synthesizing programs from domain-specific languages. Gaunt et al. [80] proposed a hand-crafted program template and simultaneously optimized the parameters of the differentiable program with gradient descent. The HOUDINI framework [81] proposed type-directed (top-down) enumeration and evolution approaches over differentiable functional programs. Shah et al. [82] hierarchically assembled differentiable programs and used neural networks for the approximation of missing expression in partial programs. Cui and Zhu [83] treated CFGs stochastically with trainable production rule sampling weights, which were optimized with a gradient-based approach [23]. However, naively applying gradient-based approaches does not work in our search spaces due to the exponential explosion of supernet weights, but still renders an interesting direction for future work.

²Note that the tree-structured search space can equivalently be described with a CFG (with a constraint on the number of maximum depth of the syntax trees).

Compared to these lines of work, we extended CFGs to handle changes in spatial resolution, promote regularity, and (compared to most of them) incorporate constraints, the latter two of which could also be applied in those domains. We also proposed a BO search strategy to search efficiently with a tailored kernel design to handle the hierarchical nature of the search space (i.e., the architectures).

J Details for Sec. 5

In this section, we provide the evaluation details (App. J.1), implementation details (App. J.2) as well as training details (App. J.3). We also present further search results and conduct analyses on the architectures observed during all of our search runs (App. J.4).

J.1 Evaluation details

For all search experiments, we compared the search strategies BANAT, Random Search (RS), Regularized Evolution (RE) [31, 84], and BANAT (WL) [18]. For implementation details of the search strategies, please refer to App. J.2. We ran search for a total of 100 evaluations with a random initial design of 10 on three seeds {777, 888, 999} on the hierarchical NAS-Bench-201 search space using 8 asynchronous workers each with a single NVIDIA RTX 2080 Ti GPU. In each evaluation, we fully trained the architectures and recorded their last validation error.

To assess the modeling performance of our surrogate, we compared regression performance of GPs with different kernels, i.e., our hierarchical WL kernel (hWL), (standard) WL kernel [18], and NASBOT’s kernel [24]. We also tried the GCN encoding [28] but it could not capture the mapping from the complex graph space to performance, resulting in constant performance predictions. Further, note that the adjacency encoding [26] and path encoding [25] cannot be used in our hierarchical search spaces since the former requires the same amount of nodes across graphs and the latter scales exponentially in the number of nodes. We ran 20 trials over the seeds {0, 1, ..., 19} and re-used the data from the search runs. In every trial, we sampled a training and test set of 700 or 500 architecture and validation error pairs, respectively. We fitted the surrogates with a varying number of training samples by randomly choosing samples from the training set without replacement, and recorded Kendall’s τ rank correlation between the predicted and true validation error.

J.2 Implementation details

BANAT & BANAT (WL) The only difference between BANAT and BANAT (WL) is that the former uses our proposed hierarchy of WL kernels (hWL), whereas the latter only uses a single WL kernel (WL) for the entire architecture (c.f., [18]). We ran BANAT asynchronously in parallel throughout our experiments with a batch size of $B = 1$, i.e., at each BO iteration a single architecture is proposed for evaluation. For the evolutionary acquisition function optimization, we used a pool size of $P = 200$, where the initial population consisted of the current ten best-performing architectures and the remainder were randomly sampled architectures to encourage exploration in the huge search spaces. During evolution, the mutation probability was set to $p_{mut} = 0.5$ and crossover probability was set to $p_{cross} = 0.5$. From the crossovers, half of them were self-crossovers of one parent and the other half were common crossovers between two parents. The tournament selection probability was set to $p_{tour} = 0.2$. We evolved the population at least for ten iterations and a maximum of 50 iterations using an early stopping criterion based on the fitness value improvements over the last five iterations.

Regularized Evolution (RE) RE [31, 84] iteratively mutates the best architectures out of a sample of the population. We reduced the population size from 50 to 30 to account for fewer evaluations, and used a sample size of 10. We also ran RE asynchronously for better comparability.

J.3 Training details

Training protocols For training of architectures on CIFAR-10/100 and ImageNet-16-120, we followed Dong and Yang [55]. We trained architectures with SGD with learning rate of 0.1, Nesterov momentum of 0.9, weight decay of 0.0005 with cosine annealing [85], and batch size of 256 for 200 epochs. The initial channels were set to 16. For both CIFAR-10 and CIFAR-100, we used random flip with probability 0.5 followed by a random crop (32x32 with 4 pixel padding) and normalization.

Table 2: Licenses for the datasets we used in our experiments.

Dataset	License	URL
CIFAR-10 [87]	MIT	https://www.cs.toronto.edu/~kriz/cifar.html
CIFAR-100 [87]	MIT	https://www.cs.toronto.edu/~kriz/cifar.html
ImageNet-16-120 [88]	MIT	https://patrykchrabaszcz.github.io/Imagenet32/
CIFARTile [86]	GNU	https://github.com/RobGeada/cvpr-nas-datasets
AddNIST [86]	GNU	https://github.com/RobGeada/cvpr-nas-datasets

For ImageNet-16-120, we used a 16x16 random crop with 2 pixel padding instead. For training of architectures on AddNIST and CIFARTile, we followed the training protocol from the CVPR-NAS 2021 competition [86]. We trained architectures with SGD with learning rate of 0.01, momentum of 0.9, and weight decay of 0.0003 with cosine annealing, and batch size of 64 for 64 epochs. We set the initial channels to 16 and did not apply any further data augmentation.

Dataset details In Tab. 2, we provide the licenses for the datasets used in our experiments. For training of architectures on CIFAR-10, CIFAR-100 [87], and ImageNet-16-120 [88], we followed the dataset splits and training protocol of NAS-Bench-201 [55]. For CIFAR-10, we split the original training set into a new training set with 25k images and validation set with 25k images following [55]. The test set remained unchanged. For evaluation, we trained architectures on both the training and validation set. For CIFAR-100, the training set remained unchanged, but the test set was partitioned in a validation set and new test set with each 5K images. For ImageNet-16-120, all splits remained unchanged. For AddNIST and CIFARTile, we used the training, validation, and test splits as defined in the CVPR-NAS 2021 competition [86].

J.4 Further search results and analyses

Supplementary to Fig. 2 (top), Fig. 7 compares the cell-based vs. hierarchical NAS-Bench-201 search space from Section 6.1 using RS, RE, and BANAT (WL). The cell-based search space design shows on par or stronger performance on all datasets except for CIFARTile for the three search strategies. In contrast, for our proposed search strategy BANAT we find on par (CIFAR-10/100) or superior (ImageNet-16-120, CIFARTile, and AddNIST) performance using the hierarchical search space design. This clearly shows that the increase of the search space does not necessarily yields the discovery of stronger neural architectures. Further, it exemplifies the importance of a strong search strategy to search effectively and efficiently in huge hierarchical search spaces (Q2), and provides further evidence that the incorporation of hierarchical information is a key contributor for search efficiency (Q3). Based on this, we believe that future work using, e.g., graph neural networks as a surrogate, may benefit from the incorporation of hierarchical information.

We report the test errors of our best found architectures in Tab. 3. We observe that our search strategy BANAT finds the strongest performing architectures across all dataset (Q2, Q3). Also note that we achieve better (validation and) test performance on ImageNet-16-120 on the hierarchical than the state-of-the-art search strategy on the cell-based NAS-Bench-201 search space (i.e., +0.37%p compared to Shapley-NAS [56]) (Q1).

Search costs Search time varied across datasets from ca. 0.5 days (CIFAR-10) to ca. 1.8 days (ImageNet-16-120) using eight asynchronous workers, each with an NVIDIA RTX 2080 Ti GPU, for ca. 4 to ca. 14.4 GPU days in total.

Is our search strategy BANAT exploring well-performing architectures during search? To investigate the question, we studied density estimates of the validation error of proposed candidates for all search strategies across our experiments from Sec. 5. This provides a better view for whether search strategies are exploring well-performing architectures or wasting computational resources on low-performing architectures. Fig. 8 shows that our proposed search strategy BANAT explored better architecture candidates across all the datasets, i.e., it has smaller median validation errors and the distributions are further shifted towards smaller validation errors than for the other search strategies.

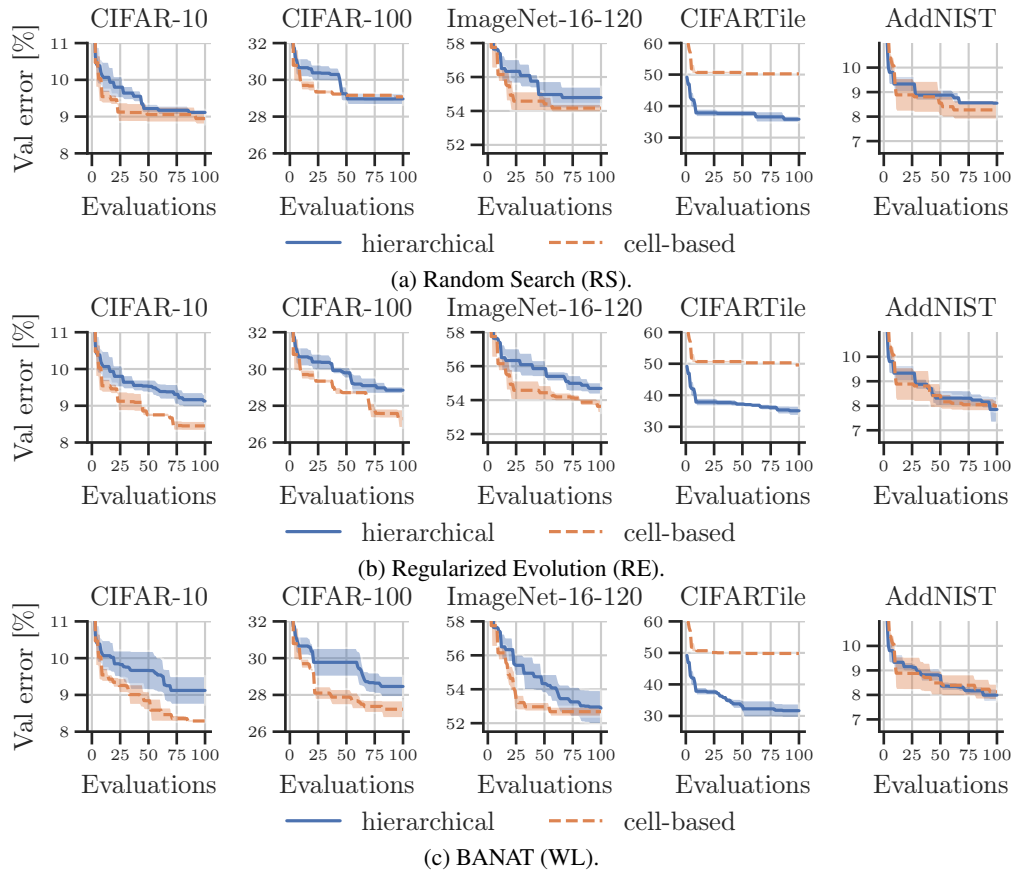


Figure 7: Cell-based vs. hierarchical search spaces. We plot mean and ± 1 standard error of the validation error on the cell-based (dashed orange) and hierarchical (solid blue) NAS-Bench-201 search space using Random Search (RS) (top), Regularized Evolution (RE) (middle), and BANAT (WL) (bottom).

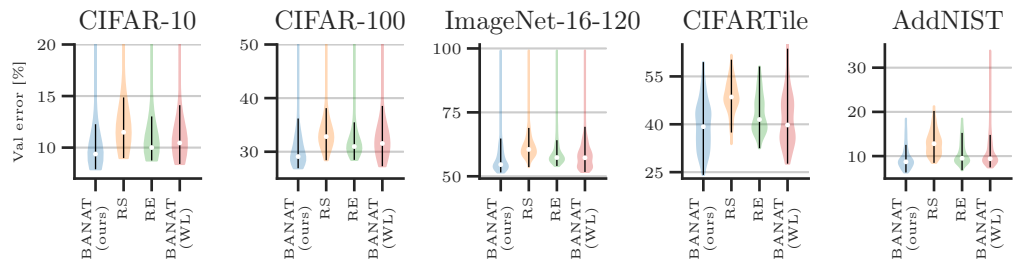


Figure 8: Density estimates for the validation error of all architecture candidates proposed by the search strategies (i.e., BANAT, RS, RE, and BANAT (WL)) and across datasets (i.e., CIFAR-10/100, ImageNet-16-120, CIFARTile, AddNIST) from our experiments in Sec. 5.

Table 3: Test errors (and ± 1 standard error) of popular baseline architectures (e.g., ResNet [89] and EfficientNet [2] variants), and our best found architectures on the cell-based and hierarchical NAS-Bench-201 search space. Note that we picked the ResNet and EfficientNet variant based on the test error, consequently giving an overestimate of their test performance.

[†] optimal numbers as reported in Dong and Yang [55].

(val best) test error (and ± 1 standard error) across three seeds {777, 888, 999} of the best architecture of the three search runs with lowest validation error.

Method	CIFAR-10		CIFAR-100		ImageNet-16-120		CIFAR100		AddNIST	
	cell-based	hierarchical	cell-based	hierarchical	cell-based	hierarchical	cell-based	hierarchical	cell-based	hierarchical
Best ResNet [89]	6.49 \pm 0.24 (32)		27.1 \pm 0.67 (110)		54.83 \pm 0.78	54.7 \pm 0.18 (56)	57.8 \pm 0.57 (18)		7.82 \pm 0.36	8.05 \pm 0.29
Best EfficientNet [2]	11.73 \pm 0.1 (B0)		35.17 \pm 0.42 (B6)		53.92 \pm 0.6	77.73 \pm 0.29 (B0)	61.01 \pm 0.62 (B0)		7.69 0.35	7.56 \pm 0.69
NAS-Bench-201 oracle [†]		5.63	26.49			52.69				
RS	6.39 \pm 0.18	6.77 \pm 0.1	28.75 \pm 0.57	29.49 \pm 0.57	54.83 \pm 0.78	54.7 \pm 0.18	57.8 \pm 0.57	40.93 \pm 0.81	7.82 \pm 0.36	8.05 \pm 0.29
RE [31, 84]	5.76 \pm 0.17	6.88 \pm 0.24	27.68 \pm 0.55	30.0 \pm 0.32	53.92 \pm 0.6	55.39 \pm 0.54	52.79 \pm 0.59	40.99 \pm 2.89	7.69 0.35	7.56 \pm 0.69
BANAT (WL) [18]	5.68 \pm 0.11	6.98 \pm 0.5	27.66 \pm 0.18	28.7 \pm 0.64	53.67 \pm 0.39	53.47 \pm 0.86	52.81 \pm 0.27	35.75 \pm 1.58	7.86 \pm 0.41	8.2 \pm 0.37
BANAT	5.68 \pm 0.11	6.0 \pm 0.16	27.66 \pm 0.18	27.57 \pm 0.46	53.67 \pm 0.39	53.43 \pm 0.61	52.81 \pm 0.27	32.28 \pm 2.39	7.86 \pm 0.41	6.09 \pm 0.34
BANAT (best)	5.64 \pm 0.14	5.65 \pm 0.09	27.03 \pm 0.23	27.63 \pm 0.2	53.54 \pm 0.43	52.78 \pm 0.23	53.18 \pm 0.91	30.33 \pm 0.77	8.04 \pm 0.45	6.33 \pm 0.59

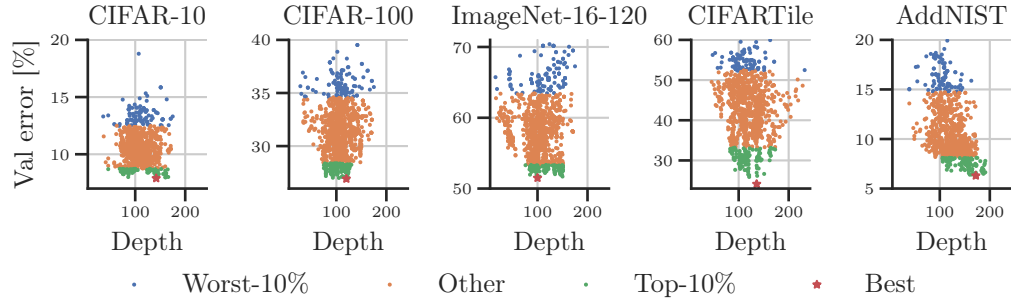


Figure 9: Validation error over maximal depth of all architecture candidates proposed by the search strategies (i.e., BANAT, RS, RE, and BANAT (WL)) and across datasets (i.e., CIFAR-10/100, ImageNet-16-120, CIFAR10Tile, AddNIST) from our experiments in Sec. 5.

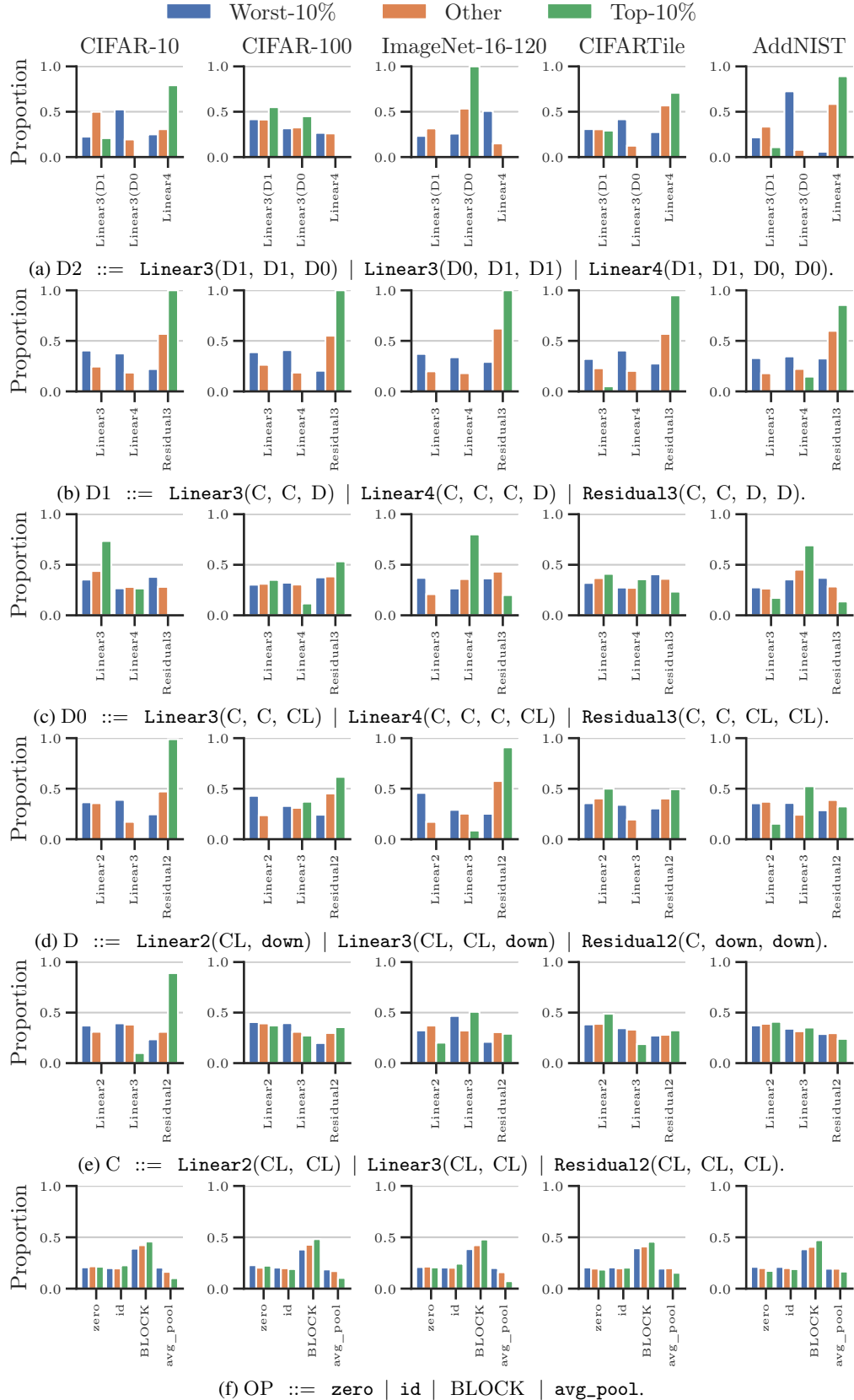
What distinguishes top-performing neural architectures from the other ones? To understand what distinguishes top-performing neural architecture from other ones, we analyze the impact of maximum depth on performance and the frequency of production rules in the worst-10%, top-10%, or other neural architectures, respectively. In another analysis, we marginalize out the validation error of every production rule; thereby relating the contribution of a production rule with the performance of the architecture. Note, however, that both analyses ignore the topological information, i.e., a topological operator or primitive computation may have a different effect at different stages of the architecture.

Fig. 9 shows no particular trend (e.g., more depth yields better performance) across the datasets, indicating that depth may not be the most important factor for performance in our hierarchical search space. In contrast, Fig. 10 and Fig. 11 show that particularly macro-level production rules (i.e., for the nonterminals D2, D1, D0, and D) have a large effect on the performance of an architecture. Interestingly, we find that top-performing architectures (almost exclusively) use the topological operator `Residual3` for derivations from the nonterminal D1 across search spaces. This hints that a residual connection at the macro-level could be a strong topological structure, but remains to be evaluated for a variety of architectures. Cell-level production choices have less effect on performance. However, we hypothesize that this may also be due to the neglect of topological information. We leave further analysis for future work.

What is the impact of flexible parameterization of convolutional blocks? To investigate the impact of the flexible parameterization of the convolutional blocks (i.e., activation functions, normalizations, and type of convolution), we removed the flexible parameterization and allowed only the same primitive computations as in the cell-based NAS-Bench-201 search space, while still searching over the macro architecture. More explicitly, we only allow ReLU non-linearity as the activation function, batch normalization as the normalization, and 1×1 or 3×3 convolutions. Fig. 12 shows that for all datasets except CIFAR-100, flexible parameterization of the convolutional blocks improves performance of the found architectures. Interestingly, we find an architecture on CIFAR-100, which achieves 26.24 % test error with 1.307 MB and 167.172 M number of parameters or FLOPs, respectively. This architecture is superior to the optimal architecture in the cell-based NAS-Bench-201 search space. Note that this architecture is also pareto-optimal for test error vs. number of parameters and test error vs. number of FLOPs.

Test error vs. number of parameters and FLOPs Fig. 13 shows the test error vs. the number of parameters or FLOPs. Our best found architectures fall well within the parameter and FLOPs ranges of the cell-based NAS-Bench-201 search space across all datasets, except for the parameters on CIFAR-10. Note that our best found architecture on ImageNet-16-120 is pareto-optimal for test error vs. number of parameters and test error vs. number of FLOPs.

Best architectures Below we report the best found architecture per dataset on the hierarchical NAS-Bench-201 search space (Sec. 5) for each dataset. Fig. 14 provides a graphical summary of the best architectures. Fig. 15 visualizes the novel and diverse design of the architectures (including stem and classifier head).



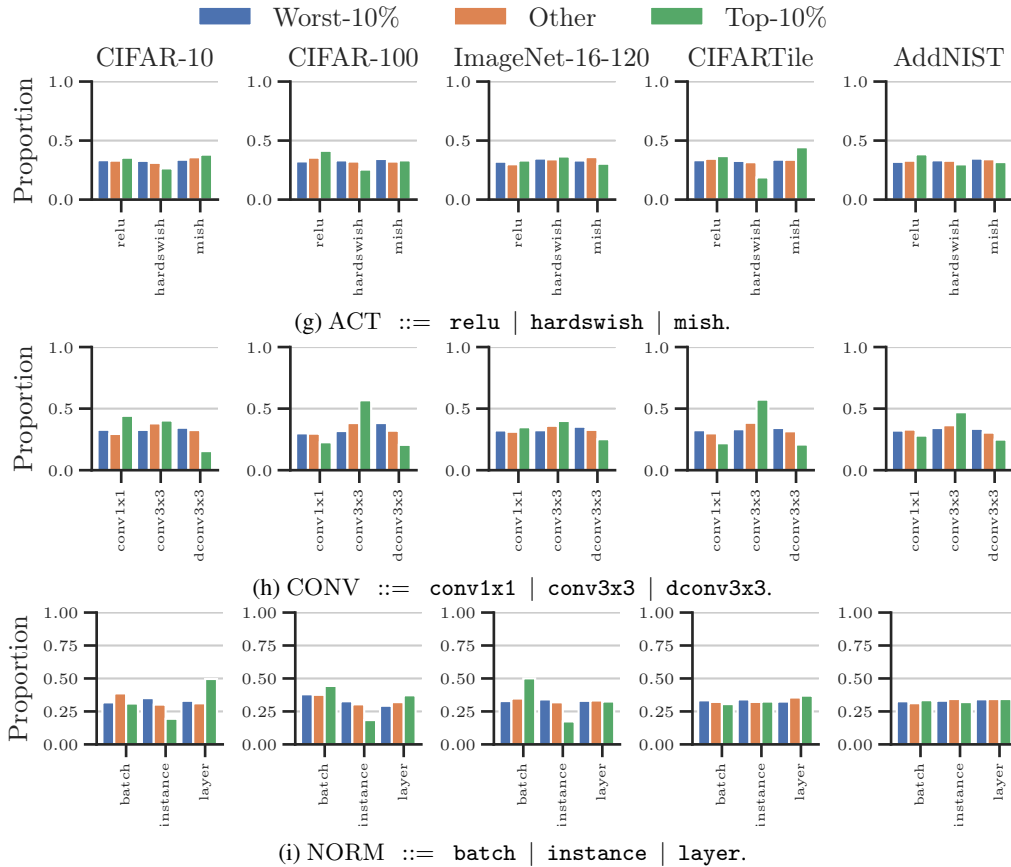


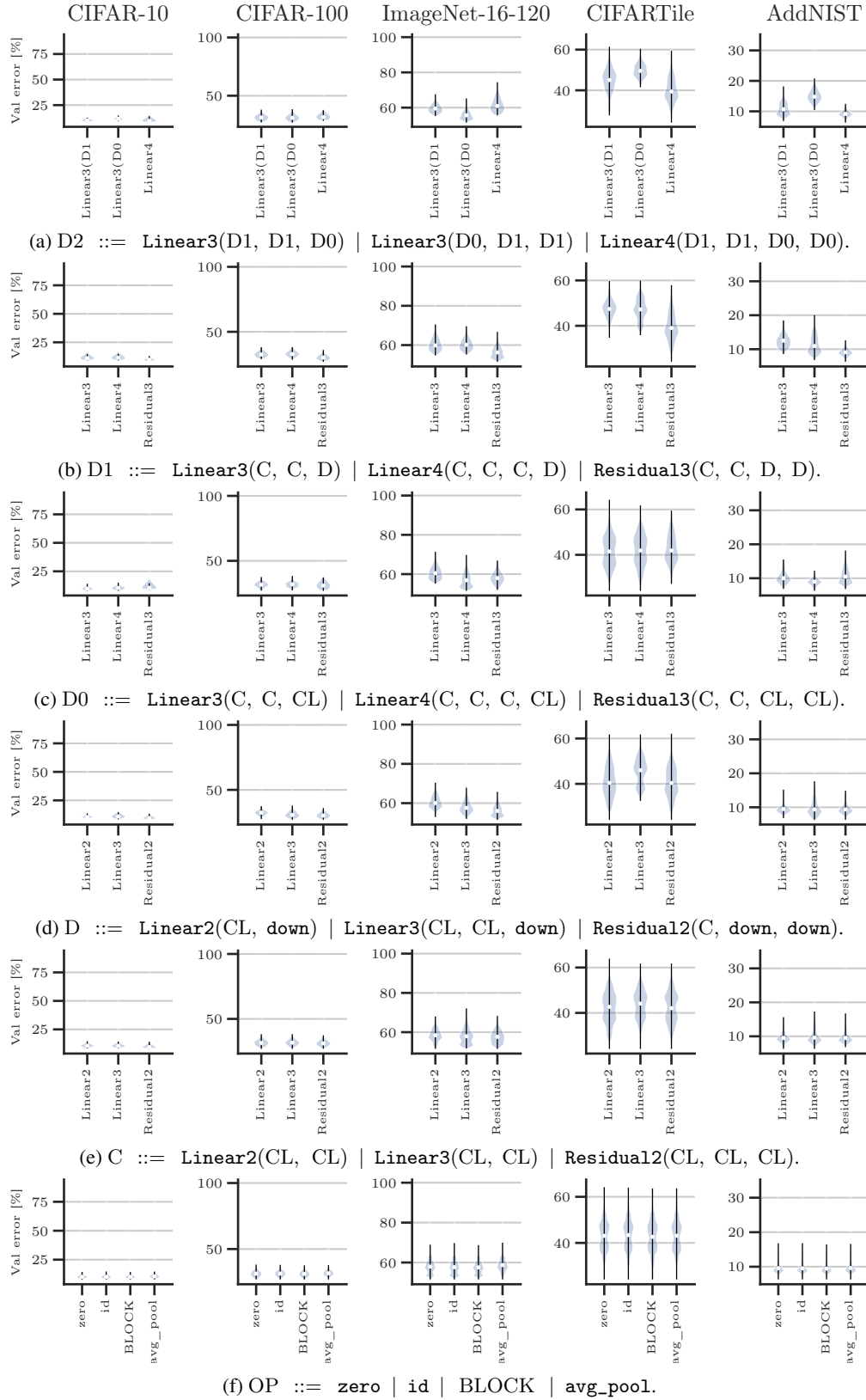
Figure 10: Comparison of the proportion of production rules in the worst-10% (blue), top-10% (green) and other (orange) neural architectures from our experiments in Sec. 5.

CIFAR-10 (mean test error 5.65 %, #params 2.204 MB, FLOPs 127.673 M):

```

Linear4(Residual3(Residual2(Cell(id, zero, Linear1(Linear3(hardswish, conv1x1,
layer)), Linear1(Linear3(hardswish, conv3x3, layer))), zero, Linear1(Linear3(mish,
conv3x3, instance))), Cell(Linear1(Linear3(relu, dconv3x3, layer))), id, avg_pool, Lin-
ear1(Linear3(relu, dconv3x3, layer))), id, zero, Cell(zero, Linear1(Linear3(relu, conv1x1,
layer))), id, Linear1(Linear3(hardswish, conv1x1, instance))), Linear1(Linear3(hardswish,
conv3x3, layer)), Linear1(Linear3(hardswish, dconv3x3, layer))), Residual2(Cell(id,
zero, Linear1(Linear3(relu, conv1x1, layer))), Linear1(Linear3(mish, conv1x1, layer))),
Linear1(Linear3(hardswish, conv3x3, layer)), zero), Cell(id, zero, id, Linear1(Linear3(relu,
conv3x3, batch))), id, id), Cell(Linear1(Linear3(hardswish, conv3x3, layer)), Lin-
ear1(Linear3(hardswish, conv1x1, layer)), Linear1(Linear3(relu, conv1x1, layer)), Lin-
ear1(Linear3(relu, conv3x3, layer)), zero, id)), Residual2(Cell(Linear1(Linear3(hardswish,
conv1x1, instance))), Linear1(Linear3(hardswish, dconv3x3, batch)), Linear1(Linear3(mish,
dconv3x3, instance))), Linear1(Linear3(relu, conv1x1, batch))), id, id), down, down),
Residual2(Cell(Linear1(Linear3(hardswish, conv1x1, layer))), Linear1(Linear3(hardswish,
dconv3x3, batch)), Linear1(Linear3(relu, conv1x1, batch))), Linear1(Linear3(hardswish,
conv3x3, layer))), id, avg_pool), down, down)), Residual3(Residual2(Cell(id, zero,
Linear1(Linear3(hardswish, conv1x1, layer))), Linear1(Linear3(hardswish, conv3x3,
layer))), id, Linear1(Linear3(mish, conv3x3, instance))), Cell(Linear1(Linear3(relu,
dconv3x3, layer))), id, avg_pool, Linear1(Linear3(relu, dconv3x3, layer))), id, zero),
Cell(zero, Linear1(Linear3(relu, conv1x1, layer))), id, Linear1(Linear3(hardswish, conv1x1,
instance))), Linear1(Linear3(hardswish, conv1x1, layer)), Linear1(Linear3(hardswish,
dconv3x3, layer))), Residual2(Cell(id, zero, Linear1(Linear3(mish, conv1x1, layer))),
Linear1(Linear3(mish, conv3x3, layer))), Linear1(Linear3(hardswish, dconv3x3,
batch))), zero), Cell(id, zero, id, Linear1(Linear3(relu, conv3x3, batch))), id, id),
Cell(Linear1(Linear3(hardswish, conv3x3, layer)), Linear1(Linear3(hardswish, conv1x1,
layer))), Linear1(Linear3(hardswish, conv3x3, layer))

```



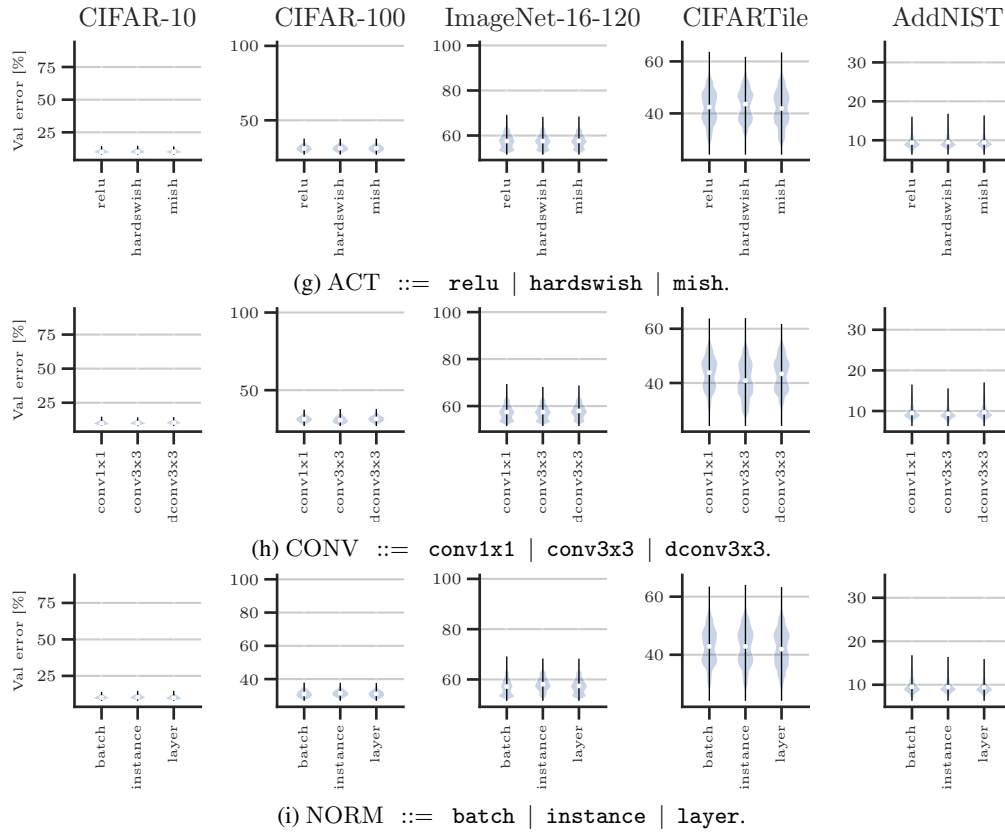


Figure 11: Marginalized performance of every production rule in our hierarchical NAS-Bench-201 search space from Sec. 5.

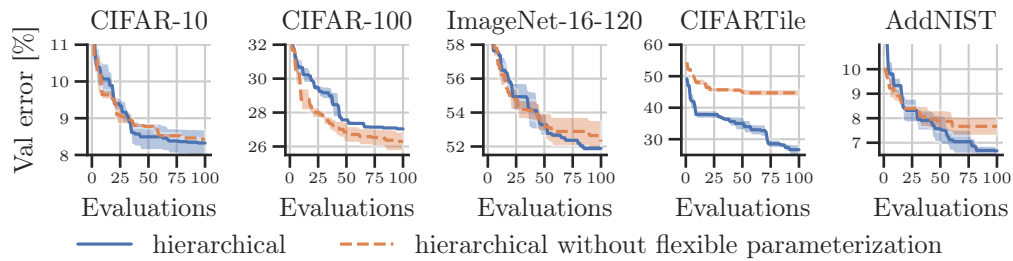


Figure 12: Impact of flexible parameterization of convolutional blocks in the hierarchical NAS-Bench-201 search space.

layer)), Linear1(Linear3(relu, conv3x3, layer)), id), Cell(Linear1(Linear3(relu, conv1x1, batch)), id, id, Linear1(Linear3(relu, conv3x3, batch)), id, id), Cell(id, Linear1(Linear3(relu, conv1x1, instance))), Linear1(Linear3(relu, conv1x1, instance))), Linear1(Linear3(relu, conv1x1, layer)), zero, Linear1(Linear3(hardswish, conv3x3, layer))), Cell(id, Linear1(Linear3(hardswish, conv1x1, layer)), Linear1(Linear3(mish, conv1x1, batch))), id, zero, id)))

CIFAR-100 (mean test error 27.63 %, #params 0.962 MB, FLOPs 115.243 M):

Linear3(Residual3(Linear3(Cell(Linear3(mish, conv3x3, layer), avg_pool, Linear3(hardswish, conv1x1, instance), zero, Linear3(mish, conv3x3, batch), zero), Cell(Linear3(hardswish, dconv3x3, batch), zero, Linear3(hardswish, dconv3x3, batch), Linear3(relu, dconv3x3, batch), id, id), Cell(Linear3(mish, conv3x3, batch), zero, id, zero, Linear3(hardswish, dconv3x3, batch), id)), Linear2(Cell(id, zero, Linear3(mish, conv3x3, batch), zero, zero, Linear3(mish, conv1x1, batch))), Cell(zero, zero, zero, id, zero, avg_pool)), Cell(Linear3(relu, conv3x3, batch), zero, Linear3(hardswish, conv3x3, instance), id, id, avg_pool), Cell(id, id, zero, zero, id, id)), Residual3(Linear3(Cell(Linear3(mish, conv3x3, layer), id, Linear3(hardswish, dconv3x3, layer), Linear3(hardswish, dconv3x3, batch), Linear3(mish, conv3x3, instance), Linear3(mish, conv3x3, batch))), Cell(Linear3(hardswish, conv1x1, layer), id, Linear3(hardswish, dconv3x3, batch), Linear3(relu, conv3x3, layer), id, id), Cell(Linear3(relu, conv3x3, instance), zero, id, zero, Linear3(mish, conv3x3, batch), avg_pool)), Linear3(Cell(zero, id, Linear3(hardswish, conv1x1, layer), Linear3(mish, conv3x3, instance), Linear3(mish, conv3x3, instance), zero), Cell(Linear3(hardswish, conv1x1, layer), id, Linear3(hardswish, dconv3x3, batch), Linear3(relu, conv3x3, batch), id, id), Cell(Linear3(relu, conv3x3, instance), zero, id, zero, Linear3(mish, conv3x3, layer), avg_pool)), Residual2(Cell(zero, id, zero, Linear3(mish, conv3x3, layer), avg_pool, Linear3(mish, conv3x3, layer)), down, down), Residual2(Cell(zero, id, zero, Linear3(mish, conv3x3, batch), avg_pool, Linear3(mish, conv3x3, layer)), down, down)), Residual3(Linear3(Cell(Linear3(mish, conv3x3, layer), id, Linear3(hardswish, dconv3x3, layer), Linear3(hardswish, dconv3x3, batch), Linear3(mish, conv3x3, instance), Linear3(mish, conv3x3, batch))), Cell(Linear3(hardswish, conv1x1, layer), id, Linear3(hardswish, dconv3x3, batch), Linear3(relu, conv3x3, layer), id, id), Cell(Linear3(relu, conv3x3, instance), zero, id, zero, Linear3(mish, conv3x3, batch), avg_pool)), Linear3(Cell(Linear3(mish, conv3x3, batch), id, Linear3(hardswish, conv1x1, layer), Linear3(mish, conv3x3, instance), Linear3(mish, conv3x3, instance), zero), Cell(Linear3(hardswish, conv1x1, layer), id, Linear3(hardswish, dconv3x3, batch), Linear3(hardswish, dconv3x3, batch), id, id), Cell(Linear3(relu, conv3x3, instance), zero, id, zero, Linear3(mish, conv3x3, layer), avg_pool)), Residual2(Cell(zero, id, zero, Linear3(mish, conv3x3, layer), avg_pool, Linear3(mish, conv3x3, layer)), down, down), Residual2(Cell(zero, id, zero, Linear3(mish, conv3x3, batch), avg_pool, Linear3(mish, conv3x3, layer)), down, down)))

ImageNet-16-120 (mean test error 52.78 %, #params 0.626 MB, FLOPs 23.771 M):

Linear3(Linear4(Residual2(Cell(id, avg_pool, id, id, Linear3(relu, dconv3x3, layer), zero), Cell(Linear3(hardswish, conv1x1, batch), zero, zero, Linear3(mish, dconv3x3, layer), zero, zero), Cell(Linear3(relu, dconv3x3, layer), Linear3(mish, dconv3x3, layer), zero, Linear3(hardswish, conv3x3, layer), Linear3(relu, dconv3x3, instance), Linear3(hardswish, conv3x3, instance))), Linear2(Cell(zero, Linear3(relu, conv3x3, layer), Linear3(mish, conv1x1, batch), Linear3(mish, conv1x1, batch), avg_pool, Linear3(relu, conv3x3, layer)), Cell(id, id, Linear3(mish, conv3x3, layer), Linear3(relu, conv3x3, instance), id, id)), Residual2(Cell(zero, avg_pool, Linear3(mish, conv1x1, batch), Linear3(mish, conv1x1, layer), zero, zero), Cell(id, Linear3(relu, dconv3x3, layer), zero, zero, Linear3(relu, dconv3x3, instance), zero), Cell(id, Linear3(relu, conv3x3, layer), id, zero, zero, id)), Cell(zero, Linear3(hardswish, conv3x3, layer), avg_pool, zero, Linear3(hardswish, conv1x1, layer), id)), Residual3(Residual2(Cell(Linear3(relu, conv1x1, instance), Linear3(mish, conv1x1, layer), Linear3(mish, conv1x1, instance), zero, Linear3(hardswish, dconv3x3, layer), id), Cell(id, avg_pool, avg_pool, Linear3(relu, conv1x1, instance), id, zero), Cell(avg_pool, Linear3(mish, conv3x3, instance), Linear3(mish, conv1x1, instance), Linear3(relu, dconv3x3, batch), id, Linear3(hardswish, conv3x3, instance))), Linear2(Cell(zero, Linear3(relu, conv3x3, layer), Linear3(mish, conv1x1, batch), Linear3(mish, conv1x1, batch), avg_pool, Linear3(relu, conv3x3, instance))), Cell(id, zero, Linear3(mish, conv3x3, layer), Linear3(relu, conv3x3, instance), id, id)), Residual2(Cell(Linear3(mish, conv3x3, layer), Linear3(mish, conv1x1, batch), id, Linear3(mish, conv1x1, layer), zero, id), down, down), Residual2(Cell(Linear3(relu, conv3x3, layer), zero, Linear3(relu, dconv3x3, layer), Linear3(mish,

conv1x1, layer), zero, id), down, down)), Residual3(Residual2(Cell(Linear3(mish, conv1x1, instance), Linear3(mish, conv1x1, layer), Linear3(mish, conv1x1, instance), avg_pool, Linear3(hardswish, dconv3x3, layer), id), Cell(id, avg_pool, avg_pool, Linear3(relu, conv1x1, instance), id, zero), Cell(avg_pool, Linear3(mish, conv3x3, instance), Linear3(mish, conv1x1, instance), Linear3(relu, dconv3x3, batch), id, Linear3(hardswish, conv3x3, instance))), Linear2(Cell(zero, Linear3(relu, conv3x3, layer), Linear3(mish, conv1x1, batch), Linear3(mish, conv1x1, batch), avg_pool, Linear3(relu, conv3x3, layer))), Cell(id, zero, Linear3(mish, conv3x3, layer), Linear3(relu, conv3x3, instance), id, id)), Residual2(Cell(Linear3(relu, conv3x3, layer), avg_pool, id, Linear3(mish, conv3x3, instance), zero, id), down, down), Residual2(Cell(Linear3(relu, conv3x3, layer), zero, Linear3(relu, dconv3x3, instance), Linear3(mish, conv1x1, layer), zero, id), down, down)))

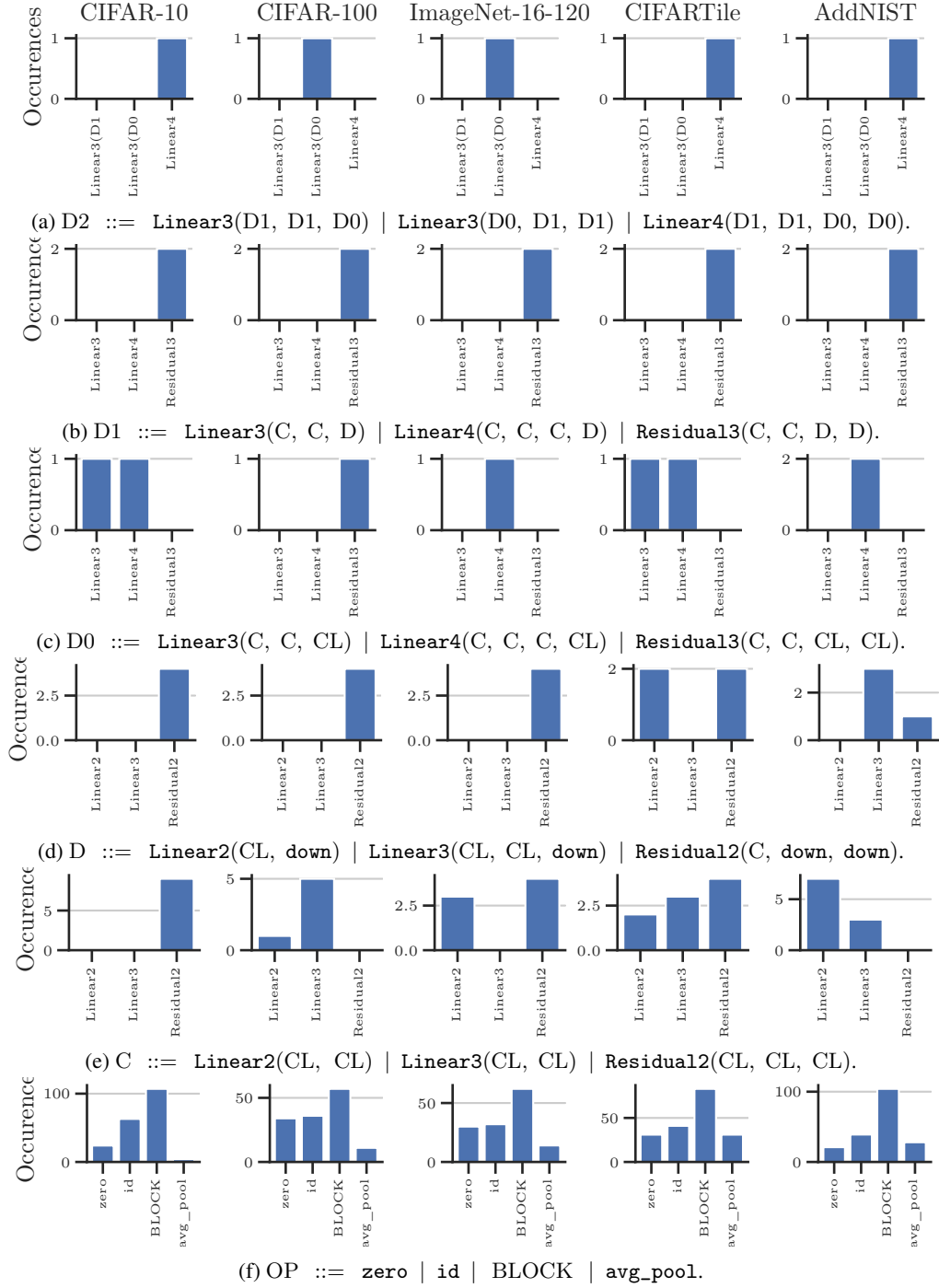
CIFAR10 (mean test error 30.33 %, #params 2.356 MB, FLOPs 372.114 M):

Linear4(Residual3(Residual2(Cell(Linear3(hardswish, conv3x3, instance), id, zero, Linear3(relu, dconv3x3, instance), Linear3(mish, conv1x1, instance), avg_pool), Cell(avg_pool, avg_pool, id, zero, Linear3(hardswish, conv3x3, batch), avg_pool), Cell(Linear3(relu, dconv3x3, instance), zero, id, Linear3(relu, dconv3x3, layer), id, id)), Residual2(Cell(zero, zero, Linear3(mish, conv1x1, instance), Linear3(mish, conv3x3, batch), zero, id), Cell(Linear3(mish, conv3x3, instance), zero, Linear3(relu, dconv3x3, batch), id, Linear3(mish, conv3x3, batch), id), Cell(Linear3(hardswish, dconv3x3, batch), Linear3(relu, conv3x3, batch), Linear3(relu, conv1x1, batch), zero, Linear3(relu, conv3x3, batch), id)), Linear2(Cell(Linear3(relu, dconv3x3, layer), Linear3(mish, conv1x1, layer), id, zero, Linear3(mish, conv3x3, batch), Linear3(relu, dconv3x3, layer))), down), Linear2(Cell(id, Linear3(hardswish, conv1x1, layer), id, Linear3(relu, conv1x1, instance), avg_pool, Linear3(relu, conv1x1, layer))), down)), Residual3(Residual2(Cell(id, avg_pool, avg_pool, Linear3(hardswish, dconv3x3, instance), Linear3(mish, conv1x1, layer), Linear3(hardswish, dconv3x3, instance)), Cell(id, id, Linear3(relu, dconv3x3, layer), id, id, zero), Cell(Linear3(relu, conv3x3, layer), id, avg_pool, Linear3(mish, dconv3x3, instance), Linear3(relu, conv1x1, layer), zero)), Residual2(Cell(Linear3(mish, conv3x3, batch), Linear3(mish, conv3x3, instance), zero, avg_pool, avg_pool, Linear3(mish, conv1x1, batch)), Cell(Linear3(mish, conv1x1, batch), Linear3(relu, dconv3x3, layer), zero, id, avg_pool, avg_pool), Cell(avg_pool, Linear3(hardswish, conv1x1, instance), id, avg_pool, avg_pool, Linear3(hardswish, conv1x1, instance))), Residual2(Cell(Linear3(relu, dconv3x3, batch), Linear3(relu, id, avg_pool, id, zero), down, down), Residual2(Cell(zero, zero, Linear3(relu, dconv3x3, batch), avg_pool, Linear3(hardswish, conv1x1, instance), avg_pool), down, down)), Linear4(Linear3(Cell(Linear3(hardswish, conv3x3, batch), Linear3(hardswish, conv3x3, batch), Linear3(relu, conv1x1, instance), id, Linear3(relu, conv1x1, layer), Linear3(relu, conv3x3, layer))), Cell(id, Linear3(relu, conv3x3, instance), Linear3(hardswish, conv1x1, instance), Linear3(relu, conv3x3, layer), avg_pool, Linear3(mish, conv1x1, layer))), Cell(zero, zero, id, Linear3(relu, conv3x3, batch), id, Linear3(relu, conv1x1, layer))), Linear3(Cell(Linear3(hardswish, conv3x3, batch), Linear3(hardswish, conv3x3, batch), Linear3(relu, conv1x1, instance), Linear3(relu, dconv3x3, layer), Linear3(mish, conv1x1, layer), Linear3(relu, conv3x3, batch)), Cell(id, Linear3(relu, conv3x3, instance), Linear3(hardswish, conv1x1, instance), Linear3(relu, dconv3x3, instance), avg_pool, Linear3(mish, conv1x1, layer))), Cell(zero, zero, id, Linear3(relu, conv3x3, batch), id, avg_pool)), Linear3(Cell(id, id, avg_pool, Linear3(mish, conv1x1, layer), Linear3(mish, conv3x3, batch), zero), Cell(id, Linear3(relu, conv1x1, batch), avg_pool, Linear3(relu, conv1x1, layer), avg_pool, zero), Cell(zero, Linear3(relu, conv1x1, batch), Linear3(mish, dconv3x3, batch), Linear3(mish, conv1x1, batch), id, id)), Cell(id, Linear3(hardswish, conv1x1, layer), zero, id, zero, id)), Linear3(Linear2(Cell(id, zero, Linear3(mish, dconv3x3, instance), Linear3(mish, conv3x3, batch), Linear3(mish, dconv3x3, instance), Linear3(relu, conv1x1, instance))), Cell(Linear3(relu, dconv3x3, instance), avg_pool, Linear3(mish, conv1x1, instance), Linear3(hardswish, dconv3x3, instance), id, Linear3(hardswish, conv1x1, layer))), Linear2(Cell(zero, zero, Linear3(mish, dconv3x3, instance), Linear3(relu, conv3x3, instance), Linear3(hardswish, conv3x3, batch), avg_pool), Cell(id, id, Linear3(hardswish, conv1x1, instance), avg_pool, zero, Linear3(hardswish, conv3x3, batch))), Cell(avg_pool, Linear3(mish, dconv3x3, layer), zero, avg_pool, avg_pool, zero)))

AddNIST (mean test error 6.33 %, #params 2.853 MB, FLOPs 593.856 M):

Linear4(Residual3(Linear3(Cell(id, Linear3(hardswish, dconv3x3, batch), Linear3(relu, conv1x1, layer), Linear3(mish, conv3x3, batch), avg_pool, zero), Cell(zero, zero, avg_pool, id, avg_pool, Linear3(hardswish, conv1x1, instance))), Cell(Linear3(relu, conv3x3, layer), id, zero, Linear3(mish, conv3x3, instance), id, avg_pool)), Linear2(Cell(id, Linear3(relu,

conv3x3, layer), Linear3(relu, conv3x3, layer), Linear3(hardswish, conv3x3, batch), id, Linear3(relu, conv3x3, layer)), Cell(Linear3(mish, conv3x3, instance), id, Linear3(mish, conv3x3, batch), id, avg_pool, id)), Linear3(Cell(zero, id, Linear3(relu, dconv3x3, instance), Linear3(relu, dconv3x3, layer), Linear3(relu, dconv3x3, instance), Linear3(mish, conv3x3, batch))), Cell(Linear3(mish, conv1x1, instance), zero, Linear3(relu, conv3x3, instance), id, zero, Linear3(relu, conv3x3, batch)), down), Linear3(Cell(zero, avg_pool, Linear3(hardswish, dconv3x3, layer), Linear3(relu, conv3x3, layer), Linear3(hardswish, conv1x1, instance), Linear3(hardswish, conv3x3, batch)), Cell(Linear3(hardswish, conv3x3, batch), Linear3(hardswish, conv1x1, layer), Linear3(mish, conv1x1, batch), id, Linear3(hardswish, conv3x3, batch), zero), down)), Residual3(Linear2(Cell(Linear3(mish, conv1x1, layer), avg_pool, Linear3(hardswish, dconv3x3, batch), Linear3(mish, dconv3x3, batch), id, Linear3(mish, conv3x3, layer)), Cell(zero, Linear3(relu, dconv3x3, layer), Linear3(hardswish, conv3x3, instance), avg_pool, avg_pool, zero))), Linear3(Cell(Linear3(relu, conv3x3, batch), id, Linear3(relu, conv3x3, layer), Linear3(mish, conv1x1, instance), id, Linear3(relu, dconv3x3, batch)), Cell(Linear3(mish, conv3x3, batch), Linear3(mish, conv1x1, instance), Linear3(mish, conv3x3, instance), zero, Linear3(mish, dconv3x3, layer), Linear3(relu, conv3x3, batch)), Cell(avg_pool, Linear3(mish, conv1x1, instance), Linear3(relu, conv3x3, batch), avg_pool, id, Linear3(mish, dconv3x3, batch))), Linear3(Cell(zero, avg_pool, Linear3(hardswish, dconv3x3, layer), Linear3(relu, conv3x3, batch), Linear3(hardswish, conv1x1, batch), Linear3(hardswish, conv3x3, batch)), Cell(avg_pool, Linear3(hardswish, dconv3x3, layer), Linear3(mish, conv1x1, batch), id, Linear3(hardswish, conv3x3, batch), zero), down), Residual2(Cell(zero, Linear3(mish, conv1x1, instance), Linear3(hardswish, conv1x1, instance), avg_pool, Linear3(relu, conv1x1, layer), Linear3(hardswish, dconv3x3, batch)), down, down)), Linear4(Linear2(Cell(Linear3(relu, conv3x3, instance), id, Linear3(relu, conv3x3, batch), avg_pool, zero, id), Cell(avg_pool, Linear3(hardswish, conv3x3, layer), avg_pool, Linear3(mish, conv3x3, batch), Linear3(relu, conv3x3, batch), id)), Linear2(Cell(Linear3(mish, conv1x1, layer), avg_pool, Linear3(hardswish, dconv3x3, batch), Linear3(mish, dconv3x3, batch), id, Linear3(mish, conv3x3, layer)), Cell(zero, Linear3(relu, dconv3x3, layer), Linear3(hardswish, conv3x3, instance), avg_pool, avg_pool, zero))), Linear2(Cell(id, Linear3(relu, conv3x3, instance), Linear3(relu, conv3x3, layer), Linear3(hardswish, dconv3x3, batch), id, Linear3(relu, conv3x3, layer)), Cell(Linear3(mish, conv1x1, batch), id, avg_pool, id, avg_pool, id)), Cell(id, Linear3(relu, conv3x3, layer), Linear3(mish, conv1x1, instance), Linear3(hardswish, conv3x3, batch), Linear3(mish, dconv3x3, instance), Linear3(hardswish, conv1x1, instance))), Linear4(Linear2(Cell(Linear3(relu, conv3x3, instance), id, Linear3(relu, conv3x3, batch), avg_pool, zero, id), Cell(zero, Linear3(relu, conv3x3, batch), avg_pool, Linear3(mish, conv3x3, batch), Linear3(relu, dconv3x3, instance), id)), Linear3(Cell(Linear3(relu, conv3x3, batch), id, Linear3(relu, conv3x3, layer), Linear3(mish, conv1x1, layer), id, Linear3(relu, dconv3x3, instance)), Cell(Linear3(mish, conv3x3, batch), Linear3(mish, conv1x1, instance), Linear3(hardswish, dconv3x3, instance), zero, Linear3(mish, dconv3x3, layer), Linear3(relu, conv3x3, batch)), Cell(avg_pool, Linear3(mish, conv1x1, instance), Linear3(relu, conv3x3, batch), avg_pool, id, Linear3(mish, dconv3x3, batch))), Linear2(Cell(id, Linear3(relu, conv3x3, layer), Linear3(hardswish, conv3x3, layer), Linear3(hardswish, dconv3x3, batch), id, Linear3(relu, conv3x3, layer)), Cell(Linear3(mish, conv3x3, batch), id, avg_pool, id, avg_pool, id)), Cell(id, Linear3(relu, conv3x3, layer), Linear3(mish, conv1x1, instance), Linear3(hardswish, conv3x3, batch), Linear3(mish, dconv3x3, instance), Linear3(mish, conv3x3, instance)))) .



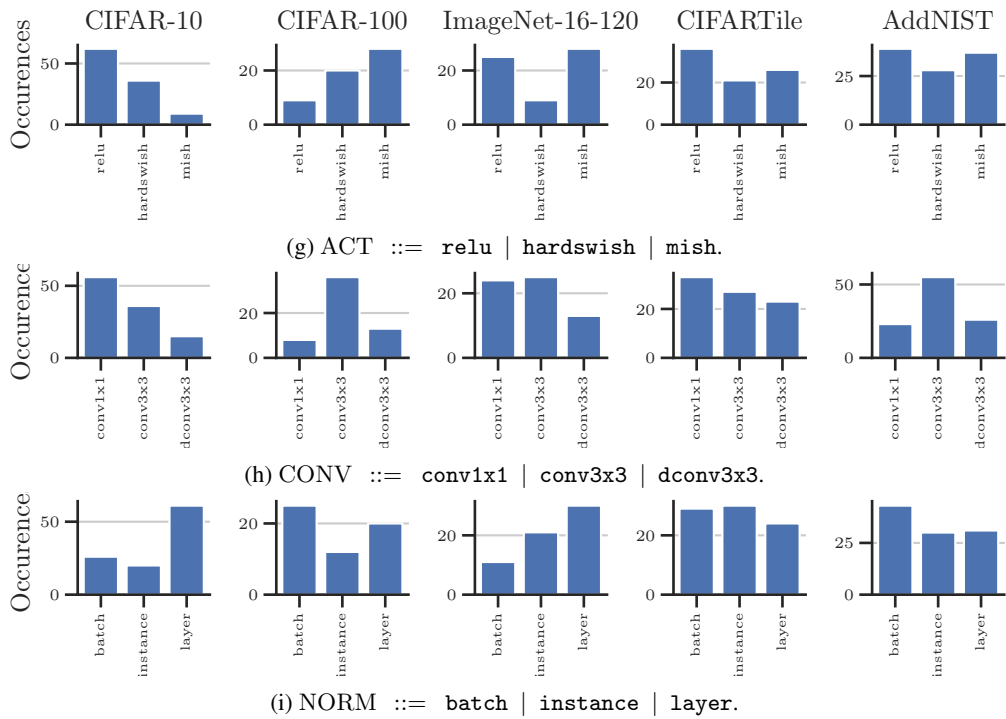
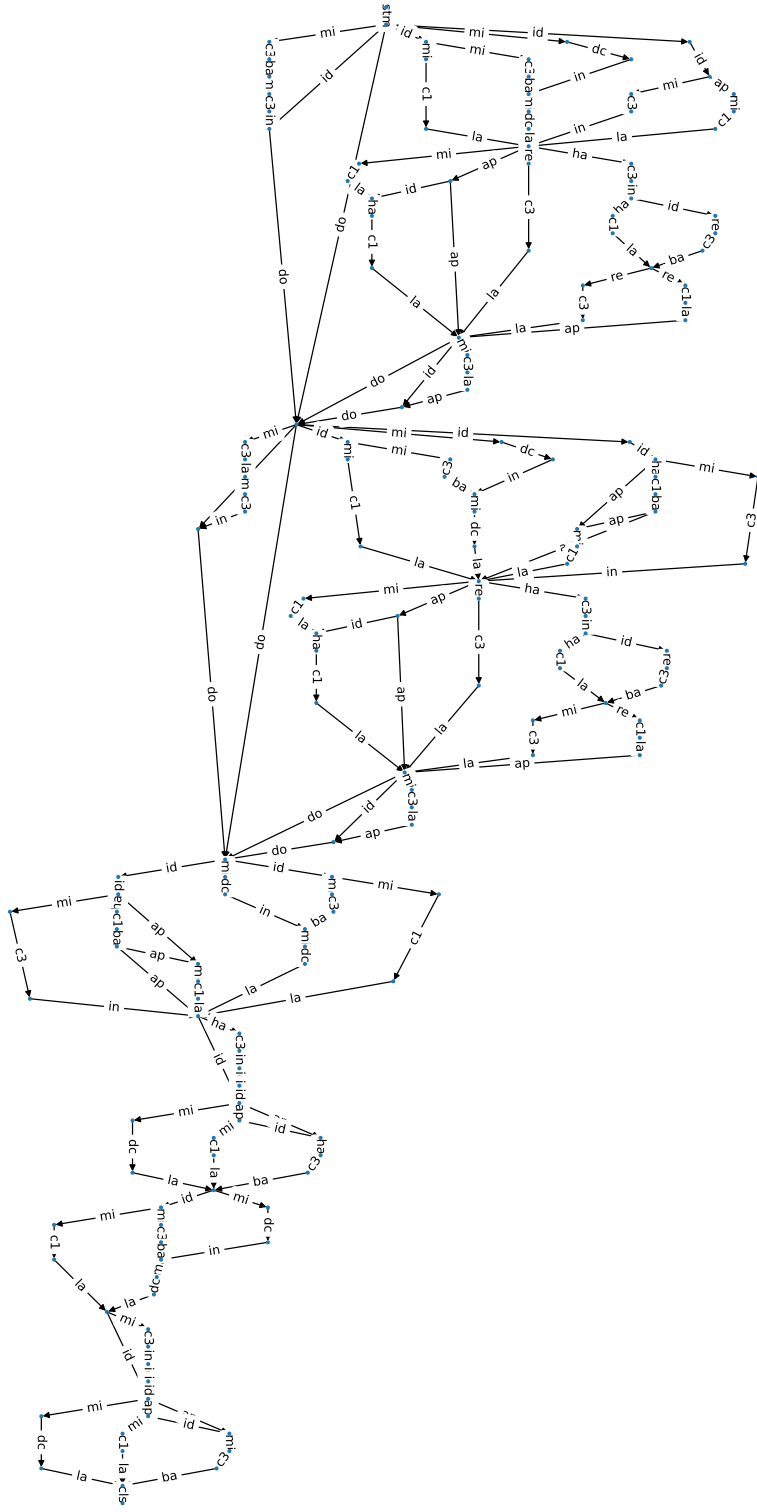
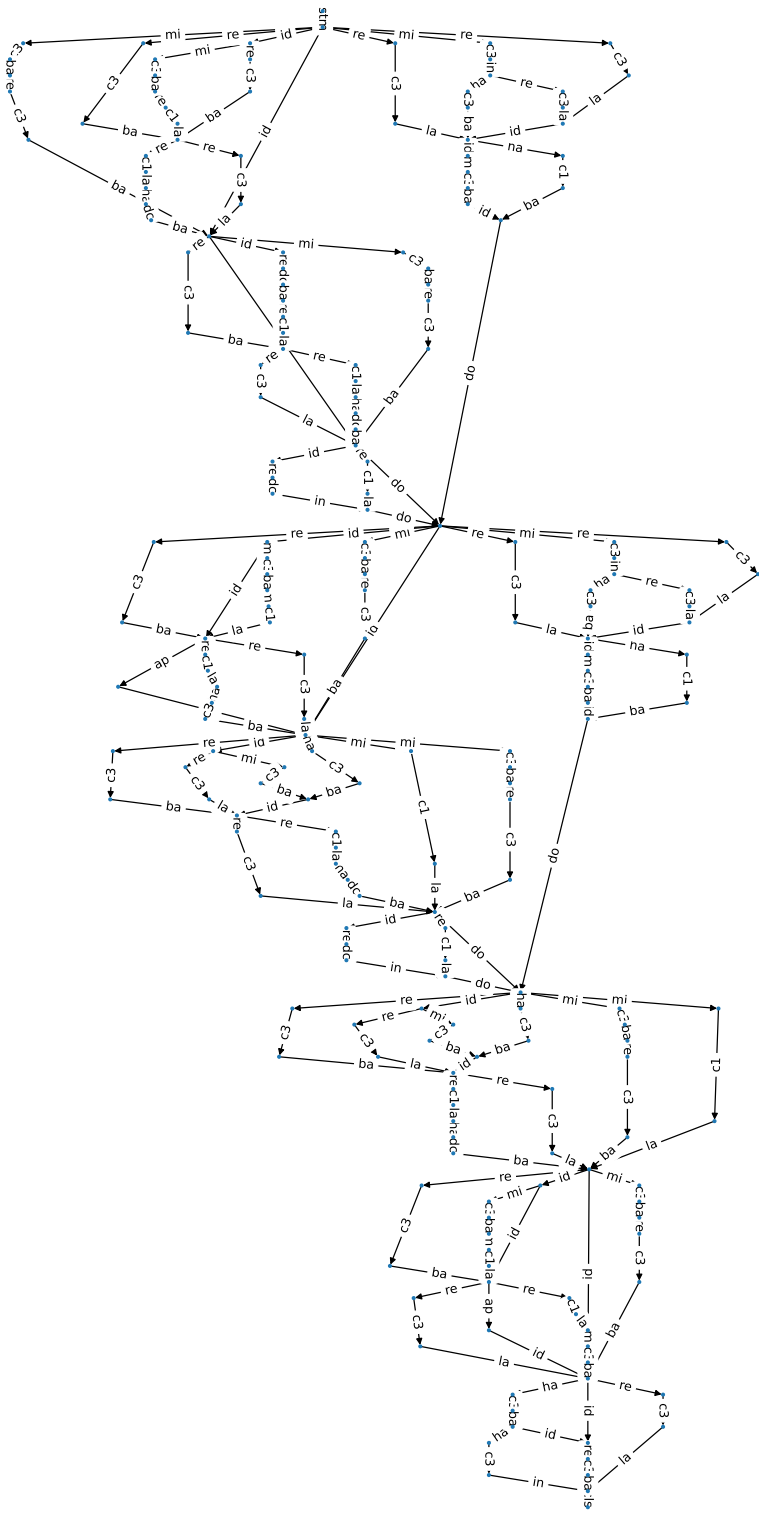


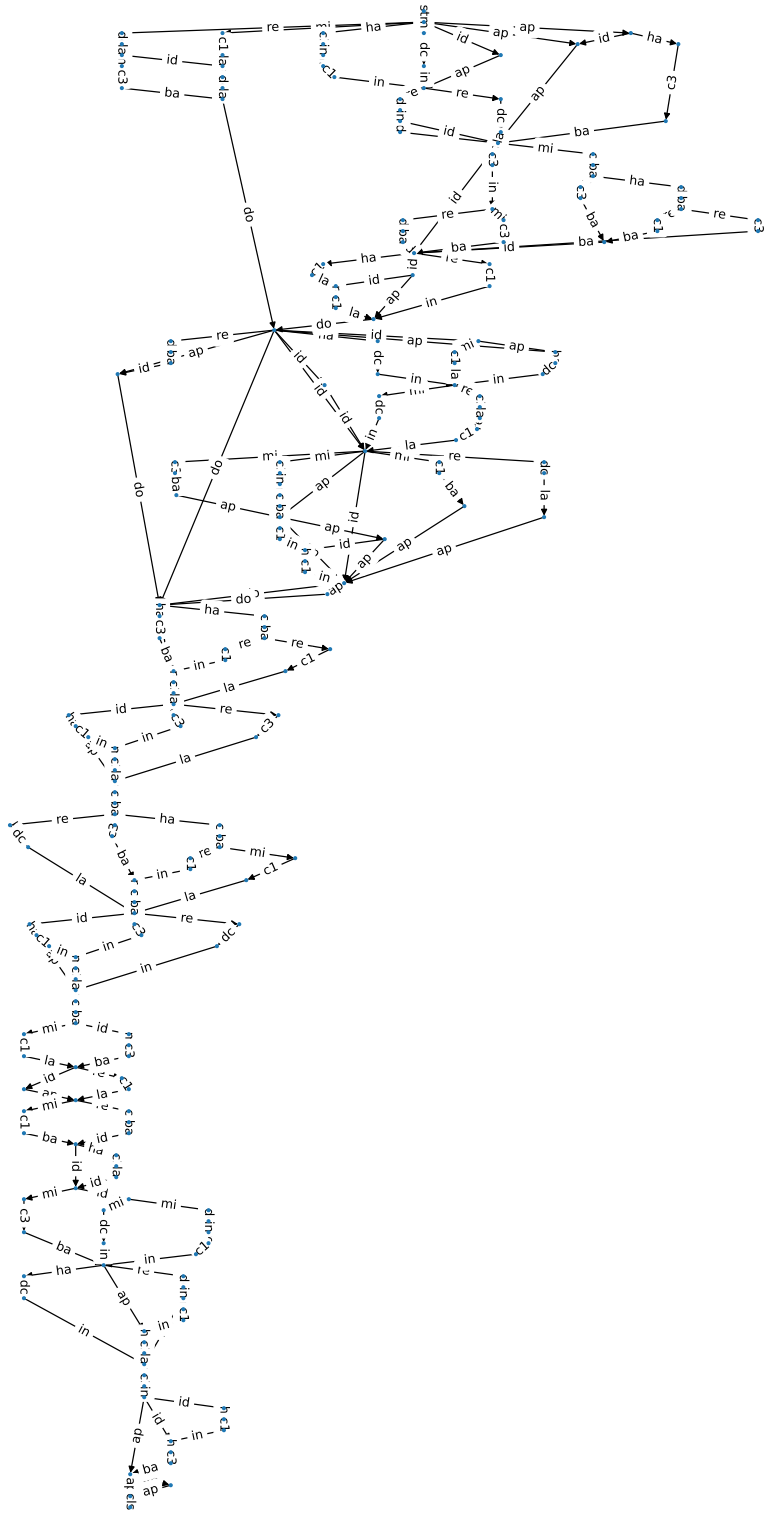
Figure 14: Summary of the distribution of topological operators and primitive operations of the best architectures on the hierarchical NAS-Bench-201 search space.



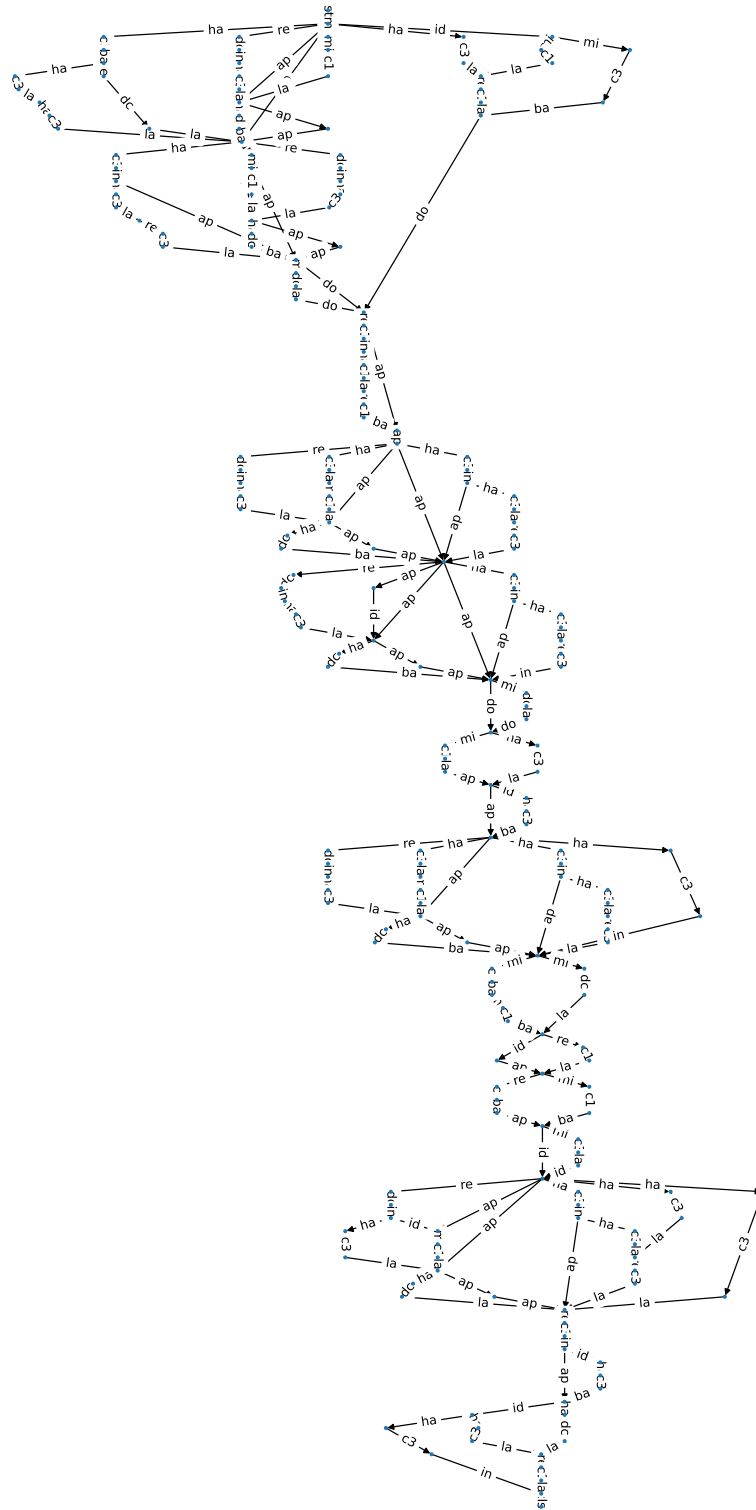
(a) CIFAR-10.



(b) CIFAR-100.



(d) CIFARTile.



(e) AddNIST.

Figure 15: Visualization of the best found architectures in our hierarchical NAS-Bench-201 search space. Abbreviations are defined as follows: ap=avg_pool, ba=batch, c1=conv1x1, c3=conv3x3, cls=classifier, dc=dconv3x3, ha=hardswish, in=instance, la=layer, mi=mish, re=relu, and stm=stem. Best viewed with zoom.