

---

# Speeding up NAS with Adaptive Subset Selection

---

Vishak Prasad<sup>1</sup> Colin White<sup>2</sup> Paarth Jain<sup>1</sup> Sibasis Nayak<sup>1</sup> Rishabh Iyer<sup>3</sup>  
Ganesh Ramakrishnan<sup>1</sup>

<sup>1</sup>Indian Institute of Technology, Bombay

<sup>2</sup>Abacus.AI

<sup>3</sup>The University of Texas at Dallas

---

**Abstract** The majority of recent developments in neural architecture search (NAS) have been aimed at decreasing the computational cost of various techniques without affecting their final performance. Towards this direction, many low-fidelity and performance prediction methods have been considered, including using subsets of the training data. In this work, we initiate the study of \*adaptive\* subset selection for NAS and present it as complementary to state-of-the-art NAS approaches. We uncover a natural connection between one-shot NAS algorithms and adaptive subset selection and devise an algorithm that makes use of state-of-the-art techniques from both areas. We use these techniques to substantially reduce the runtime of DARTS-PT, a leading one-shot NAS algorithm, without sacrificing accuracy. Our results are consistent across multiple datasets, and our code and all materials needed to reproduce our results are available at [https://anonymous.4open.science/r/SubsetSelection\\_NAS-2BE4](https://anonymous.4open.science/r/SubsetSelection_NAS-2BE4).

---

## 1 Introduction

Recent developments in neural architecture search (NAS) have focused on decreasing the cost of NAS without sacrificing performance. Towards this direction, “one-shot” methods improve the search efficiency by training just a single over-parameterized neural network (supernet) [18, 1]. For example, the popular DARTS [18] algorithm applies a continuous relaxation to the architecture parameters, allowing the architecture parameters and the weights to be simultaneously optimized via gradient descent. While many follow-up works have improved the performance of DARTS [28, 13], the algorithms are still slow and are expensive in terms of budget and CO2 emissions [27].

On the other hand, the field of subset selection for efficient machine learning model training has seen a flurry of activity. In this area of study, facility location [20], clustering [2], and other subset selection algorithms are used to select a representative subset of the training data, substantially reducing the runtime of model training. Recently, adaptive subset selection algorithms have been used to speed up model training even further [12, 11]. Adaptive subset selection is a powerful technique which regularly updates the current subset of the data as the search progresses, to ensure that the performance of the model is maintained.

In this work, we uncover a natural connection between one-shot NAS algorithms and adaptive subset selection to devise an algorithm that makes use of state-of-the-art techniques from both areas. Specifically, DARTS-PT [28] (a leading one-shot algorithm) and GLISTER [12] (a leading adaptive subset selection algorithm) are both cast as bi-level optimization problems on the training and validation sets, allowing us to formulate a combined approach, ADAPTIVE-NAS, as a mixed discrete and continuous bi-level optimization problem (see Figure 1 for an overview). Across several search spaces, we show that the resulting algorithm achieves significantly improved runtime compared to DARTS-PT, without sacrificing performance. Specifically, due to the use of adaptive subset selection, the training data can be reduced to 10% of the full training set size, resulting in an order of magnitude decrease in runtime, without sacrificing accuracy. To validate this approach, we compare to baselines such as facility location, entropy-based subset selection [21], and random subset selection.

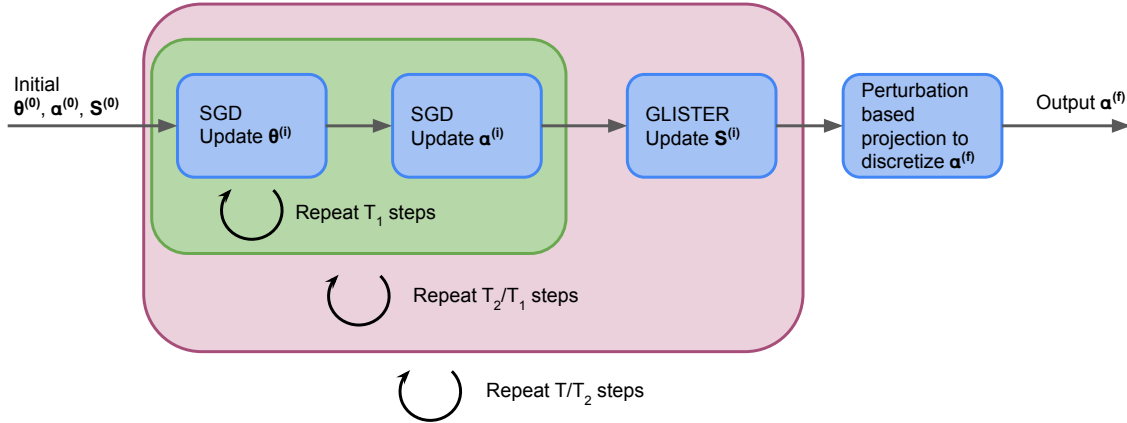


Figure 1: Overview of ADAPTIVE-NAS. The algorithm starts with initial weights  $\theta^{(0)}$ , architecture hyperparameters  $\alpha^{(0)}$ , and subset of the training data  $S^{(0)}$ . Throughout the search, the weights  $\theta^{(i)}$  and architecture weights  $\alpha^{(i)}$  are updated with SGD, and the subset  $S^{(i)}$  is updated with GLISTER, according to different time schedules. Then the final architecture  $\alpha^{(f)}$  is discretized and returned.

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

- We introduce ADAPTIVE-NAS, the first NAS algorithm to make use of adaptive subset selection. This technique is used to substantially reduce the training time needed for one-shot NAS algorithms. We also add facility location as a novel baseline for subset selection applied to NAS.
- Through extensive experiments, we show that ADAPTIVE-NAS substantially reduces the runtime needed for DARTS-PT, with no decrease in final accuracy of the returned architecture. We release all of our code and materials needed to reproduce our results.

## 2 Related Work

**Neural architecture search.** NAS has been studied since the 1980s [5, 26, 19] and has been revitalized in the last few years [35, 18, 8, 23]. One line of work aims to predict the performance of neural architectures before they are fully trained, through low-fidelity estimates such as training for fewer epochs [34, 30] or learning curve extrapolation [3, 32]. Another line of work takes a *one-shot* approach by representing the entire space of neural architectures by a single “supernet”, and then performing gradient descent to efficiently converge to a high-performing architecture [18]. Many follow-up works have improved its performance [17, 31, 13, 15, 33]. Recently, Wang et al. [28] introduced a more reliable perturbation-based operation scoring technique when returning the final architecture, achieving much stronger performance compared to DARTS.

**Subset selection.** Many tools have been developed in the field of subset selection for efficient machine learning model training. Popular fixed subset selection methods include coreset algorithms [6, 20], facility location [20], and entropy-based methods [21]. Recently, GLISTER was proposed as an adaptive subset selection method based on a greedy search [12]; an adaptive gradient-matching algorithm for subset selection was also proposed [11].

**Subset selection in NAS.** A few existing works have applied (offline) subset selection to the field of NAS. Na et al. [21] consider facility location,  $k$ -center, and entropy-based techniques, showing that the latter results in the best speedup for DARTS. Shim et al. [24] consider core-set sampling to speed up PC-DARTS by a factor of 8. Concurrent work [10] uses subset selection algorithms to speed up on multi-fidelity methods such as Hyperband [16] and ASHA [14]. Finally, other recent work uses a generative model to create a small set of *synthetic* training data, used to efficiently train architectures during NAS [25, 22].

### 3 Methodology

**Preliminaries.** In this section, we start by reviewing the ideas behind DARTS and DARTS-PT. The DARTS search space consists of a cell expressed as a directed acyclic graph, where each edge  $(i, j)$  can take on choices of operations  $o^{(i,j)}$  such as `max_pool_3x3` or `sep_conv_5x5`. Denote the full set of possible operations by  $\mathcal{O}$ . Each choice of operation for  $(i, j)$  is given a continuous variable  $\alpha^{(i,j)}$ . Denote  $\mathcal{U}$  and  $\mathcal{V}$  to be the training and validation sets, respectively. Denote  $\mathcal{L}_{\text{train}}$  and  $\mathcal{L}_{\text{val}}$  as the training and validation loss, respectively. These losses are determined by the architecture weights, the architecture itself, and the dataset.

DARTS and DARTS-PT are both gradient-based optimization methods which train a supernet consisting of weights  $\theta$  and architecture parameters  $\alpha$ . DARTS and DARTS-PT both attempt to solve the following expression via alternating gradient updates:

$$\min_{\alpha} \mathcal{L}_{\text{val}} \left( \arg \min_{\theta} \mathcal{L}_{\text{train}} (\theta, \alpha, \mathcal{U}), \alpha, \mathcal{V} \right). \quad (1)$$

**GLISTER.** GLISTER approximately solves the following expression by first approximating the bi-level optimization expression using a single gradient step, and then using a greedy data subset selection procedure [12].

$$\min_{S \subseteq \mathcal{U}, |S| \leq k} \mathcal{L}_{\text{val}} \left( \arg \min_{\theta} \mathcal{L}_{\text{train}} (\theta, \alpha, S), \alpha, \mathcal{V} \right). \quad (2)$$

**ADAPTIVE-NAS.** We present a formulation that organically combines Expressions (1) and (2) into a single mixed discrete and continuous bi-level optimization problem. The inner optimization is the minimization of training loss during architecture training, over a subset of the training data of size  $k$ . In the outer optimization, we minimize the validation loss by simultaneously optimizing for the neural architecture and for the subset of the training data. This allows us to efficiently find the best neural architecture:

$$\arg \min_{S \subseteq \mathcal{U}, |S| \leq k, \alpha} \mathcal{L}_{\text{val}} \left( \arg \min_{\theta} \mathcal{L}_{\text{train}} (\theta, \alpha, S), \alpha, \mathcal{V} \right). \quad (3)$$

Evaluating this expression is computationally prohibitive because of the expensive inner optimization. Instead, we iteratively perform a joint optimization of the inner weights, and the outer training subset and architecture. In order to iteratively update the training subset and architecture, we compute meta-approximations of the inner optimization. For the architecture, we compute

$$\nabla_{\alpha} \mathcal{L}_{\text{val}} \left( \arg \min_{\theta} \mathcal{L}_{\text{train}} (\theta, \alpha, S), \alpha, \mathcal{V} \right) \quad (4)$$

$$\approx \nabla_{\alpha} \mathcal{L}_{\text{val}} (\theta - \zeta \nabla_{\theta} \mathcal{L}_{\text{train}} (\theta, \alpha, S), \alpha, \mathcal{V}). \quad (5)$$

For the subset, following Killamsetty et al. [12], we run a greedy procedure using a similar approximation to the inner optimization:

$$\mathcal{L}_{\text{val}} \left( \arg \min_{\theta} \mathcal{L}_{\text{train}} (\theta, \alpha, S), \alpha, \mathcal{V} \right) \quad (6)$$

$$\approx \mathcal{L}_{\text{val}} (\theta - \zeta \nabla_{\theta} \mathcal{L}_{\text{train}} (\theta, \alpha, S), \alpha, \mathcal{V}). \quad (7)$$

Then we can iteratively update the outer parameters (architecture and subset), and the inner parameters (weights). Following prior work [12, 18], we only update the architecture and subset every  $t_1$  and  $t_2$  steps, respectively, for efficiency.

Table 1: Performance of one-shot NAS algorithms on NAS-Bench-201 CIFAR-10.

Performance on NAS-Bench-201 CIFAR-10			
Algorithm	Test accuracy	GPU hours	% Data used
DARTS-PT	88.21 (88.11)	7.50	100
DARTS-PT-ENTROPY	86.31 $\pm$ 4.66	<b>0.62</b>	10
DARTS-PT-RAND	86.94 $\pm$ 3.58	<b>0.62</b>	10
DARTS-PT-FL	89.27 $\pm$ 1.09	1.60	10
ADAPTIVE-NAS	<b>92.22 <math>\pm</math> 1.76</b>	0.83	10

## 4 Experiments

In this section, we describe our experimental setup and results. We run experiments on NAS-Bench-201 with CIFAR-10, CIFAR-100, and ImageNet16-120, DARTS with CIFAR-10, and DARTS-S4 with CIFAR-10. See Appendix B for details.

We run experiments with DARTS-PT, ADAPTIVE-NAS, and three other (non-adaptive) data subset selection methods added to DARTS-PT: DARTS-PT-RAND (random subset selection), DARTS-PT-FL (facility location), and DARTS-PT-ENTROPY (the entropy-based method taken by Na et al. [21]).

**Experimental setup.** Following Wang et al. [28], we use 50% of the full training dataset for supernet training and 50% for validation, and we report the accuracy of the final returned architecture on the held-out test set. In our main experiments, for all of the (adaptive and non-adaptive) subset selection methods, we set the subset size to 10% of the training dataset. We run the same experimental procedure for each method: select a size-10% subset of the full training dataset, train and discretize the supernet on the subset, and train the final architecture using the full training dataset. For DARTS-PT, we run the same procedure using the full training dataset at each step. We otherwise use the exact same training pipeline as in Wang et al. [28]: batch size of 64, learning rate of 0.025, momentum of 0.9, and cosine annealing.

We run all experiments on a NVIDIA Tesla V100 GPU. We run each algorithm with 5 random seeds, reporting the mean and standard deviation of each method, except DARTS-PT: due to the extreme runtime and availability of existing results, we run once and verify that the result is nearly identical to published results [28]. We also report the time it takes to output the final returned architecture.

**Experimental results and discussion.** In Tables 1, 2, and 3, we report the results on NAS-Bench-201. On CIFAR-10 and ImageNet16-120, ADAPTIVE-NAS achieves significantly higher accuracy than all other algorithms tested. On CIFAR-100, ADAPTIVE-NAS is essentially tied with DARTS-PT-FL for the highest accuracy. Furthermore, all NAS algorithms that use subset selection have significantly decreased runtime – ADAPTIVE-NAS sees a factor of 9 speedup compared to DARTS-PT. Note that DARTS-PT-FL takes more time when the number of examples per class in the dataset is higher, so it sees comparatively higher runtimes on CIFAR-10.

In Tables 8 and 9, we report the results on S4 CIFAR-10 and DARTS CIFAR-10. Once again, the runtime of ADAPTIVE-NAS is significantly faster than DARTS-PT – a factor of 7 speedup. On these search spaces, the performance of the subset-based methods are more similar when compared to NAS-Bench-201, and on the DARTS search space, ADAPTIVE-NAS does not outperform DARTS-PT. A possible explanation is that S4 and DARTS are significantly larger search spaces than NAS-Bench-201 and require more training data to distinguish between architectures. To test this, we included one more experiment in Table 9, giving ADAPTIVE-NAS 20% training data instead of 10%. We find that the accuracy significantly increases, moving within one standard deviation of the accuracy of DARTS-PT.

Table 2: Performance of one-shot NAS algorithms on NAS-Bench-201 CIFAR-100.

Performance on NAS-Bench-201 CIFAR-100			
Algorithm	Test accuracy	GPU hours	%Data used
DARTS-PT	61.650	8.00	100
DARTS-PT-ENTROPY	56.79 $\pm$ 7.63	<b>0.58</b>	10
DARTS-PT-RAND	56.95 $\pm$ 4.54	<b>0.58</b>	10
DARTS-PT-FL	<b>64.28 <math>\pm</math> 3.10</b>	0.67	10
ADAPTIVE-NAS	64.27 $\pm$ 3.37	0.87	10

Table 3: Performance of one-shot NAS algorithms on NAS-Bench-201 ImageNet16-120.

Performance on NAS-Bench-201 Imagenet16-120			
Algorithm	Test accuracy	GPU hours	%Data used
DARTS-PT	35.00	33.50	100
DARTS-PT-ENTROPY	26.52 $\pm$ 3.73	<b>1.58</b>	10
DARTS-PT-RAND	27.04 $\pm$ 5.53	<b>1.58</b>	10
DARTS-PT-FL	29.30 $\pm$ 5.35	1.90	10
ADAPTIVE-NAS	<b>36.10 <math>\pm</math> 6.96</b>	2.60	10

Overall, ADAPTIVE-NAS achieves the highest average performance across all search spaces. Furthermore, ADAPTIVE-NAS achieves no less than a seven-fold increase in runtime compared to DARTS-PT, on all search spaces.

In Appendix B, we give ablation studies on the percentage of data used.

## 5 Conclusions, Limitations, and Broader Impact

In this work, we used a connection between one-shot NAS algorithms and adaptive subset selection to devise an algorithm that makes use of state-of-the-art techniques from both areas. Specifically, DARTS-PT and GLISTER, leading techniques in one-shot NAS and adaptive subset selection, respectively, are both cast as bi-level optimization problems on the training and validation sets, which allowed us to formulate a combined approach, ADAPTIVE-NAS, as a mixed discrete and continuous bi-level optimization problem. We showed that the resulting algorithm is able to train on an (adaptive) dataset that is 10% of the size of the full training set, without sacrificing accuracy, resulting in an order of magnitude decrease in runtime.

**Limitations.** While ADAPTIVE-NAS uses a subset of the data when training and discretizing the supernet, the full dataset is used for training the final architecture. An interesting direction for future work is to use an adaptive subset of the data even when training the final architecture, which may lead to even faster runtime, perhaps at a small cost to performance. Another interesting direction for future work is to apply adaptive subset selection to other non supernet-based NAS algorithms such as regularized evolution [23] or BANANAS [29].

**Broader impact.** Our work combines techniques from two different areas: adaptive subset selection for machine learning, and neural architecture search. The goal of our work is to make it easier and quicker to develop high-performing architectures on new datasets. Our work also helps to unify two sub-fields of machine learning that had thus far been disjoint. Since the end product of our work is a NAS algorithm, it is not itself meant for one application but can be used in any end-application. For example, it may be used to more efficiently find deep learning architectures for applications that help to reduce CO2 emissions, or for creating large language models. Our hope is that future AI models discovered by our work will have a net positive impact, due to the push for the AI community to be more conscious about the societal impact of its work [7].

## 6 Reproducibility Checklist

### 1. For all authors...

- (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [\[Yes\]](#) [The main claims in the abstract and introduction reflect the paper's contributions and scope.]
- (b) Did you describe the limitations of your work? [\[Yes\]](#) [See Section 5.]
- (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) [See Section 5.]
- (d) Have you read the ethics author's and review guidelines and ensured that your paper conforms to them? <https://automl.cc/ethics-accessibility/> [\[Yes\]](#) [We read the ethics and accessibility guidelines and ensured our paper conforms to them.]

### 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, including all requirements (e.g., requirements.txt with explicit version), an instructive README with installation, and execution commands (either in the supplemental material or as a URL)? [\[Yes\]](#) [We include the code, data, and instructions to reproduce the results here: [https://anonymous.4open.science/r/SubsetSelection\\_NAS-2BE4](https://anonymous.4open.science/r/SubsetSelection_NAS-2BE4).]
- (b) Did you include the raw results of running the given instructions on the given code and data? [\[Yes\]](#) [We include our raw results; see [https://anonymous.4open.science/r/SubsetSelection\\_NAS-2BE4/README.md](https://anonymous.4open.science/r/SubsetSelection_NAS-2BE4/README.md).]
- (c) Did you include scripts and commands that can be used to generate the figures and tables in your paper based on the raw results of the code, data, and instructions given? [\[Yes\]](#) [We include scripts to generate our exact results. See [https://anonymous.4open.science/r/SubsetSelection\\_NAS-2BE4](https://anonymous.4open.science/r/SubsetSelection_NAS-2BE4).]
- (d) Did you ensure sufficient code quality such that your code can be safely executed and the code is properly documented? [\[Yes\]](#) [We included multiple documentation files, and put in a reasonable effort to make our code as easy to use as possible.]
- (e) Did you specify all the training details (e.g., data splits, pre-processing, search spaces, fixed hyperparameter settings, and how they were chosen)? [\[Yes\]](#) [See Sections 3 and 4.]
- (f) Did you ensure that you compared different methods (including your own) exactly on the same benchmarks, including the same datasets, search space, code for training and hyperparameters for that code? [\[Yes\]](#) [We compared to different methods in exactly the same setting. See Section 4.]
- (g) Did you run ablation studies to assess the impact of different components of your approach? [\[Yes\]](#) [See Section 4.]
- (h) Did you use the same evaluation protocol for the methods being compared? [\[Yes\]](#) [See Section 4.]

- (i) Did you compare performance over time? [No] [Since we only compared one-shot NAS methods, we did not compare performance over time.]
  - (j) Did you perform multiple runs of your experiments and report random seeds? [Yes] [We ran 5 random seeds for each method.]
  - (k) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] [We reported standard deviation. See Section 4.]
  - (l) Did you use tabular or surrogate benchmarks for in-depth evaluations? [Yes] [We used NAS-Bench-201.]
  - (m) 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] [See Section 4.]
  - (n) Did you report how you tuned hyperparameters, and what time and resources this required (if they were not automatically tuned by your AutoML method, e.g. in a NAS approach; and also hyperparameters of your own method)? [N/A] [We did not tune any hyperparameters.]
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 cited all assets we used.]
  - (b) Did you mention the license of the assets? [N/A] [All assets were public.]
  - (c) Did you include any new assets either in the supplemental material or as a URL? [No] [We did not include new assets.]
  - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [N/A] Our experiments were conducted only on publicly available datasets.
  - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A] Our experiments were conducted only on publicly available datasets.
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 conduct research 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 conduct research 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 conduct research with human subjects.]

## References

- [1] Gabriel Bender, Pieter-Jan Kindermans, Barret Zoph, Vijay Vasudevan, and Quoc Le. Understanding and simplifying one-shot architecture search. In *ICML*, 2018.
- [2] Kevin Clark, Minh-Thang Luong, Quoc V Le, and Christopher D Manning. Electra: Pre-training text encoders as discriminators rather than generators. *arXiv preprint arXiv:2003.10555*, 2020.
- [3] Tobias Domhan, Jost Tobias Springenberg, and Frank Hutter. Speeding up automatic hyperparameter optimization of deep neural networks by extrapolation of learning curves. In *IJCAI*, 2015.

- [4] Xuanyi Dong and Yi Yang. Nas-bench-201: Extending the scope of reproducible neural architecture search. In *ICLR*, 2020.
- [5] WB Dress. Darwinian optimization of synthetic neural systems. Technical report, Oak Ridge National Lab., TN (United States), 1987.
- [6] Sariel Har-Peled and Soham Mazumdar. On coresets for k-means and k-median clustering. In *Proceedings of the thirty-sixth annual ACM symposium on Theory of computing*, pages 291–300, 2004.
- [7] Brent Hecht, Lauren Wilcox, Jeffrey P Bigham, Johannes Schöning, Ehsan Hoque, Jason Ernst, Yonatan Bisk, Luigi De Russis, Lana Yarosh, Bushra Anjum, Danish Contractor, and Cathy Wu. It’s time to do something: Mitigating the negative impacts of computing through a change to the peer review process. *ACM Future of Computing Blog*, 2018.
- [8] Kirthevasan Kandasamy, Willie Neiswanger, Jeff Schneider, Barnabas Poczos, and Eric P Xing. Neural architecture search with bayesian optimisation and optimal transport. In *NeurIPS*, 2018.
- [9] Vishal Kaushal, Ganesh Ramakrishnan, and Rishabh Iyer. Submodlib: A submodular optimization library. *arXiv preprint arXiv:2202.10680*, 2022.
- [10] Krishnateja Killamsetty, Guttu Sai Abhishek, Alexandre V Evfimievski, Lucian Popa, Ganesh Ramakrishnan, Rishabh Iyer, et al. Automata: Gradient based data subset selection for compute-efficient hyper-parameter tuning. *arXiv preprint arXiv:2203.08212*, 2022.
- [11] Krishnateja Killamsetty, S Durga, Ganesh Ramakrishnan, Abir De, and Rishabh Iyer. Gradmatch: Gradient matching based data subset selection for efficient deep model training. In *International Conference on Machine Learning*, pages 5464–5474. PMLR, 2021.
- [12] Krishnateja Killamsetty, Durga Sivasubramanian, Ganesh Ramakrishnan, and Rishabh Iyer. Glistar: Generalization based data subset selection for efficient and robust learning. *arXiv preprint arXiv:2012.10630*, 2020.
- [13] Kevin Alexander Laube and Andreas Zell. Prune and replace nas. *arXiv preprint arXiv:1906.07528*, 2019.
- [14] Liam Li, Kevin Jamieson, Afshin Rostamizadeh, Ekaterina Gonina, Moritz Hardt, Benjamin Recht, and Ameet Talwalkar. A system for massively parallel hyperparameter tuning. In *MLSys Conference*, 2020.
- [15] Liam Li, Mikhail Khodak, Maria-Florina Balcan, and Ameet Talwalkar. Geometry-aware gradient algorithms for neural architecture search. In *Proceedings of the International Conference on Learning Representations (ICLR)*, 2021.
- [16] Lisha Li, Kevin Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar. Hyperband: A novel bandit-based approach to hyperparameter optimization. In *JMLR*, 2018.
- [17] Hanwen Liang, Shifeng Zhang, Jiacheng Sun, Xingqiu He, Weiran Huang, Kechen Zhuang, and Zhenguo Li. Darts+: Improved differentiable architecture search with early stopping. *arXiv preprint arXiv:1909.06035*, 2019.
- [18] Hanxiao Liu, Karen Simonyan, and Yiming Yang. Darts: Differentiable architecture search. In *ICLR*, 2019.



- [19] Geoffrey F Miller, Peter M Todd, and Shailesh U Hegde. Designing neural networks using genetic algorithms. In *ICGA*, volume 89, pages 379–384, 1989.
- [20] Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of machine learning models. In *International Conference on Machine Learning*, pages 6950–6960. PMLR, 2020.
- [21] Byunggook Na, Jisoo Mok, Hyeokjun Choe, and Sungroh Yoon. Accelerating neural architecture search via proxy data. *arXiv preprint arXiv:2106.04784*, 2021.
- [22] Aditya Rawal, Joel Lehman, Felipe Petroski Such, Jeff Clune, and Kenneth O. Stanley. Synthetic petri dish: A novel surrogate model for rapid architecture search, 2020.
- [23] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V Le. Regularized evolution for image classifier architecture search. In *AAAI*, 2019.
- [24] Jae-hun Shim, Kyeongbo Kong, and Suk-Ju Kang. Core-set sampling for efficient neural architecture search. *arXiv preprint arXiv:2107.06869*, 2021.
- [25] Felipe Petroski Such, Aditya Rawal, Joel Lehman, Kenneth O. Stanley, and Jeff Clune. Generative teaching networks: Accelerating neural architecture search by learning to generate synthetic training data, 2019.
- [26] Manoel Tenorio and Wei-Tsih Lee. Self organizing neural networks for the identification problem. *Advances in Neural Information Processing Systems*, 1, 1988.
- [27] Tanja Tornede, Alexander Tornede, Jonas Hanselle, Marcel Wever, Felix Mohr, and Eyke Hüllermeier. Towards green automated machine learning: Status quo and future directions. *arXiv preprint arXiv:2111.05850*, 2021.
- [28] Ruo Chen Wang, Minhao Cheng, Xiangning Chen, Xiaocheng Tang, and Cho-Jui Hsieh. Rethinking architecture selection in differentiable nas. *arXiv preprint arXiv:2108.04392*, 2021.
- [29] Colin White, Willie Neiswanger, and Yash Savani. Bananas: Bayesian optimization with neural architectures for neural architecture search. In *AAAI*, 2021.
- [30] Colin White, Arber Zela, Binxin Ru, Yang Liu, and Frank Hutter. How powerful are performance predictors in neural architecture search? *arXiv preprint arXiv:2104.01177*, 2021.
- [31] Yuhui Xu, Lingxi Xie, Xiaopeng Zhang, Xin Chen, Guo-Jun Qi, Qi Tian, and Hongkai Xiong. Pc-darts: Partial channel connections for memory-efficient architecture search. In *ICLR*, 2019.
- [32] Shen Yan, Colin White, Yash Savani, and Frank Hutter. Nas-bench-x11 and the power of learning curves. In *NeurIPS*, 2021.
- [33] Arber Zela, Thomas Elsken, Tonmoy Saikia, Yassine Marrakchi, Thomas Brox, and Frank Hutter. Understanding and robustifying differentiable architecture search. In *ICLR*, 2020.
- [34] Dongzhan Zhou, Xinchu Zhou, Wenwei Zhang, Chen Change Loy, Shuai Yi, Xuesen Zhang, and Wanli Ouyang. Econas: Finding proxies for economical neural architecture search. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 11396–11404, 2020.
- [35] Barret Zoph and Quoc V. Le. Neural architecture search with reinforcement learning. In *ICLR*, 2017.

---

**Algorithm 1** ADAPTIVE-NAS

---

**Require:** Training data  $\mathcal{U}$ , Validation data  $\mathcal{V}$ , Initial subset  $S^{(0)}$  of size  $k$ , Initial parameters  $\theta^{(0)}$  and  $\alpha^{(0)}$ , steps  $T_1, T_2$ , and  $T$ .  
for all steps  $t$  in  $T$  do  
  if  $t \bmod T_1 == 0$ :  
     $S^{(t)} = \text{GreedyDSS}(\mathcal{U}, \mathcal{V}, \theta^{(t-1)}, \alpha^{(t-1)})$   
  else:  
     $S^{(t)} = S^{(t-1)}$   
  if  $t \bmod T_2 == 0$ :  
    Perform one step of SGD to update  $\alpha^{(t)}$  using  $V$   
  else:  
     $\alpha^{(t)} = \alpha^{(t-1)}$   
    Perform one step of SGD to update  $\theta^{(t)}$  using  $S^{(t)}$  and  $\alpha^{(t)}$ .  
Run discretization of the supernet as in DARTS-PT, using  $S^{(T)}$ , to return final architecture  $\alpha$   
Train  $\alpha$  using SGD with the full training set  $\mathcal{U}$   
**Return:** Final trained architecture  $\alpha$

---

## A Additional Details from Section 3

For the full ADAPTIVE-NAS procedure, see Algorithm 1.

## B Additional Results and Analyses

In this section, we give additional results and analyses to supplement Section 4.

**Search spaces.** We run experiments on NAS-Bench-201 with CIFAR-10, CIFAR-100, and ImageNet16-120, DARTS with CIFAR-10, and DARTS-S4 with CIFAR-10.

NAS-Bench-201 [4] is a cell-based search space which contains 15 625 architectures, or 6 466 architectures that are unique up to isomorphisms. Each cell is a directed acyclic graph consisting of four nodes and six edges. Each of the six edges have five choices of operations. The cell is then stacked several times to form the final architecture.

The DARTS search space [18] is a cell-based search space containing  $10^{18}$  architectures. It consists of a normal cell and a reduction cell, each of which is represented as a directed acyclic graph with four nodes and two incoming edges per node. Each edge has eight choices of operations, and a choice of input node. Similar to NAS-Bench-201, the cells are stacked several times to form the final architecture.

Zela et al. [33] proposed a variant of the DARTS search space, S4, which replaces the original set of eight choices of operations with just two operations:  $3 \times 3 \text{ SepConv}$ , and *Noise*, where *Noise* replaces the feature map values by noise drawn from  $\epsilon \sim \mathcal{N}(0, 1)$ . This search space was designed to test the failure modes of one-shot NAS methods such as DARTS, since *Noise* should not be chosen as it actively hurts performance. S4 and DARTS have no differences other than the operation set.

**Methods tested.** We run experiments with DARTS-PT, ADAPTIVE-NAS, and three other (non-adaptive) data subset selection methods added to DARTS-PT. We describe the details of each approach below.

- DARTS-PT: see Section 3. we use the original implementation [28].
- DARTS-PT-RAND: similar to DARTS-PT, but the supernet is trained and discretized using a random subset of the training data.
- DARTS-PT-FL: similar to DARTS-PT, but the supernet is trained and discretized using a subset of the training data, selected using facility location. The facility location algorithm was implemented using the naive greedy algorithm and run on each class separately, using a dense Euclidean metric. We used the `submodlib` library [9].

- DARTS-PT-ENTROPY [21]: similar to DARTS-PT, but the supernet is trained and discretized using a subset of the training data, selected using a combination of high-, and low-entropy datapoints. The cost of NAS is reduced by selecting a representative set of the original training data. Unlike the other existing zero cost subset selection methods for NAS, this approach is specifically tailored for NAS and accelerates neural architecture search using proxy data. The entropy of a datapoint is calculated by training a base neural architecture from the search space, and computing whether the output probability is low or high. This approach was taken by Na et al. [21] to speed up DARTS.
- ADAPTIVE-NAS: described in the previous section; see Algorithm 1.

### B.1 Details of GLISTER

The optimization that we are trying to solve for GLISTER, equation (2), can be written as

$$S^{t+1} = \operatorname{argmin}_{S \subseteq U, |S| \leq k} G_{\theta^t}(S) \quad (8)$$

where  $G_{\theta^t}(S)$  is equation ???. Since equation (8) is an instance of submodular optimization (as proven in Theorem 1 of [12]), it can be regularized using another function such as supervised facility location. The regularized objective can be written as

$$S^{t+1} = \operatorname{argmin}_{S \subseteq U, |S| \leq k} G_{\theta^t}(S) + \lambda R(S) \quad (9)$$

where  $\lambda$  is a trade-off parameter. GreedyDSS refers to the set of greedy algorithms and approximations that solves 9. Greedy Taylor Approximation algorithm (GreedyTaylorApprox( $U, V, \theta^0, \eta, k, r, \lambda, R$ ), described as Algorithm 2 in [12]) is used as GreedyDSS in our work. Here,  $U$  and  $V$  are the training and validation set respectively.  $\theta^t$  is the current set of parameters,  $\eta$  is the learning rate,  $k$  is the budget, parameter  $r$  governs the number of times we perform the Taylor series approximation, and  $\lambda$  is the regularization constant.

### B.2 Details of facility location

Intuitively, facility location, attempts to model representation of the datapoints. If  $s_{ij}$  is the similarity between datapoints  $i$  and  $j$ , define  $f(X)$  such that

$$f(X) = \sum_{i \in V} \max_{j \in X} s_{ij} \quad (10)$$

where  $V$  is the ground set. If the ground set of items are assumed clustered, an alternative clustered implementation of facility location is computed over the clusters as

$$f(X) = \sum_{l \in 1..K} \sum_{i \in C_l} \max_{j \in X \cap C_l} s_{ij} \quad (11)$$

### B.3 Details of DARTS-PT-ENTROPY

DARTS-PT-ENTROPY is the implementation of [21] where the cost of NAS is reduced by selecting a representative set of the original training data. The entropy of a datapoint is defined as

$$\operatorname{Entropy}(x : M) = - \sum_{\tilde{y}} p(\tilde{y}|x, M) \log p(\tilde{y}|x, M) \quad (12)$$

where  $\tilde{y} = M(x)$  is the predictive distribution of  $x$  w.r.t. a pre-trained baseline model  $M$ .

The selection method selects datapoints from both the high entropy and low entropy regions.

Table 4: Performance of DARTS-PT + GLISTER on DARTS and S4.

Performance on CIFAR-10			
Search Space	Test accuracy	GPU hours	% Data used
DARTS	97.06 ± 1.09	2.65	10
S4	97.28 ± 0.03	1.18	10

Table 5: Ablation study of DARTS-PT-GLISTER on NAS-Bench-201 search space CIFAR-10.

Performance on NAS-Bench-201 CIFAR-10		
Test accuracy	GPU hours	%Data used
57.91 ± 43.77	0.30	1
89.21 ± 3.12	0.35	2
91.95 ± 1.56	0.58	5
92.22 ± 1.76	0.83	10
<b>92.24 ± 1.65</b>	1.58	20
88.23 ± 0.06	3.63	50

If  $h_x$  is a bin of the data point  $x$  on data entropy histogram  $H$ ,  $|h_x|$  is the height of  $h_x$  (number of examples in  $h_x$ ), three probabilities are defined as

$$P_{\{1,2,3\}}(x; H) = \text{norm}(W_{\{1,2,3\}}(h_x; H) / |h_x|) \quad (13)$$

where  $\text{norm}()$  is a normalizer such that the probability terms add to 1.  $W_{\{1,2,3\}}$  are selected such that the tail end entropy data are likely to be selected over the middle entropy data points.

In [21],  $P_1(x)$  was the highest performer. We have used  $P_1(x)$  in our experiments.

In Table B.3, we give a summary of the improvements of ADAPTIVE-NAS when compared to DARTS-PT.

Search Space	Dataset	Accuracy	Time reduced	% Data Used
NAS-Bench-201	CIFAR-10	+5.07	8.62 times	10
NAS-Bench-201	CIFAR-100	+2.63	9.20 times	10
NAS-Bench-201	Imagenet-16-120	+1.10	12.80 times	10
S4	CIFAR-10	-0.01	7.76 times	10
DARTS	CIFAR-10	-0.44	7.49 times	10
DARTS	CIFAR-10	-0.20	4.58 times	20

In Table 4, we give the results of ADAPTIVE-NAS with using the full data for the DARTS-PT projection step for search spaces DARTS and S4. Although we were able to get better accuracy (when compared 10% data on projection step) on DARTS space, the accuracy went down a little bit for S4.

In Tables 5 and 6, we give the tables to match the plots in Figure 2.

#### B.4 Results of DARTS-PT-GRAD-MATCH

In Section 3, we introduced GLISTER applied to DARTS-PT as ADAPTIVE-NAS. However, there is another choice of adaptive subset selection algorithm: GRAD-MATCH [11]. In this section, we describe GRAD-MATCH and give results on GRAD-MATCH applied to DARTS-PT, showing that it does not work as well as ADAPTIVE-NAS.

Table 6: Ablation study of DARTS-PT-GLISTER on NAS-Bench-201 search space CIFAR-10 with the full data for the DARTS-PT projection.

Performance on NAS-Bench-201 CIFAR-10		
Test accuracy	GPU hours	%Data used
92.68 ± 1.47	2.28	1
91.30 ± 1.69	2.33	2
92.88 ± 1.51	2.50	5
<b>93.85 ± 0.51</b>	2.66	10
91.81 ± 1.84	3.23	20
88.15 ± 0.14	4.73	50

Table 7: Performance of DARTS-PT + GRAD-MATCH on NAS-Bench-201

Performance on NAS-Bench-201 CIFAR-10			
Dataset	Test accuracy	GPU hours	% Data used
CIFAR-10	88.83 ± 1.09	0.87	10
CIFAR-100	63.70 ± 3.98	0.87	10

The core idea of GRAD-MATCH is to find a subset of the original training set whose gradients match the gradients of the training/validation set. The gradient error term can be given as

$$Err(w^t, X^t, L, L_T, \theta_t) = \left\| \sum_{i \in X^t} w_i^t \nabla_{\theta} L_T^i(\theta_t) - \nabla_{\theta} L(\theta_t) \right\| \quad (14)$$

where  $w^t$  are the weights produced by the adaptive data selection algorithm,  $X^t$  are the subsets selected,  $L_t$  is the training loss,  $L$  is either the training or the validation loss, and  $\theta_t$  are the classifier model parameters. In GRAD-MATCH, minimizing the above equation is reformulated as optimizing an equivalent submodular function with approximations for efficiency.

GRAD-MATCH was integrated with DARTS-PT the same way we described integrating GLISTER with DARTS-PT in Section 3. Furthermore, we used the same hyperparameters for DARTS-PT as with ADAPTIVE-NAS. We denote the algorithm as DARTS-PT-GRAD-MATCH. The results were found to be better than some baselines but still not better than ADAPTIVE-NAS for CIFAR-10 and CIFAR-100 on NAS-Bench-201. While GRAD-MATCH is faster compared to GLISTER, the bottleneck in one-shot NAS is training the supernet, therefore, the improved performance of GLISTER makes it the better fit to be incorporated with DARTS-PT in creating ADAPTIVE-NAS.

## C Details from Section 4

In this section, we give more details for the experiments conducted in Section 4.

### C.1 Experiments on NAS-Bench-201

We used the original code from DARTS-PT [28] and GLISTER [12]. The DARTS-PT code consists of two parts. The supernet training and a perturbation based projection part to discretize  $\alpha$ . The Supernet training is run for 100 epochs and at each 10 epoch interval, we select a new subset of data by passing the model and architecture parameters. At every epoch, we use 10% of the original dataset. We use a batchsize of 64, learning rate of 0.025, momentum of 0.9, and cosine annealing. We use 50% of data for training and 50% for validation, as in the DARTS-PT paper [28]. The last 10% data subset is saved and used for the perturbation based projection part of DARTS-PT. We run the projection part for 25 epochs. For subset selection, we used the same code of GLISTER with selection algorithm run on each class separately.

Table 8: Performance of one-shot NAS algorithms on S4 search space CIFAR-10.

Performance on S4 CIFAR-10			
Algorithm	Test accuracy	GPU hours	%Data used
DARTS-PT	97.31 (97.36)	8.38	100
DARTS-PT-ENTROPY	97.45 $\pm$ 0.10	<b>0.86</b>	10
DARTS-PT-RAND	97.40 $\pm$ 0.06	<b>0.86</b>	10
DARTS-PT-FL	97.34 $\pm$ 0.13	1.08	10
ADAPTIVE-NAS	97.30 $\pm$ 0.12	1.08	10

For DARTS-PT-FL, we used the implementation of Facility Location as present in `submodlib`. This subset selection algorithm was used in the dense Euclidean setting. The algorithm is used separately for each class so as to keep the representation across classes the same as original. It was optimised using the ‘NaiveGreedy’ algorithm. For the experiments, 10% data was used.

For DARTS-PT-RAND and DARTS-PT-ENTROPY, we combined DARTS-PT with proxy data using two methods of subset selection techniques for dataset, one a random subset selection and other an entropy based subset selection technique [21]. For random subset data was chosen randomly from the dataset. For the entropy based selection, we used the entropy files for CIFAR-10 and CIFAR-100 from Na et al. [21] which was obtained by training a baseline network of ResNet20 and ResNet56 respectively. For ImageNet16-120, we trained a ResNet-50 model from the PyTorch model zoo.

For S4 and the DARTS search space, we used the same configuration as for NAS-Bench-201. Since S4 and DARTS are non-tabular, we used a separate evaluation code for computing the performance of the selected architecture. We used the same evaluation code given in DARTS-PT. The code uses a batch size of 96, learning rate of 0.025, momentum of 0.9 and weight decay of 0.025. The architecture is trained for 600 epochs.

**Ablation study.** To explore the effect of the percentage of data used, in Figure 2 (left), we run ADAPTIVE-NAS with different percentages of the training data, ranging from 1% to 50%. In Figure 2 (Right), we run the same experiment using the full training data in the projection step of ADAPTIVE-NAS. Interestingly, we see a definitive U-shape in the first experiment: the highest accuracy with ADAPTIVE-NAS is at 20%, achieving accuracy *higher* than the standard setting of 100% data (i.e., DARTS-PT). Since the supernet is an over-parameterized model of weights and architecture parameters, and ADAPTIVE-NAS regularly updates the training subset to maximize validation accuracy, ADAPTIVE-NAS may help the supernet from overfitting. Furthermore, in the second experiment, we see that the relatively, the accuracies are much more consistent when varying the percentage of the training set used, when the projection step is allowed to use the full training dataset. Therefore, keeping the full training dataset for the projection step leads to higher and more consistent performance, at the expense of more GPU-hours.

Overall, based on the ablation studies in Figure 2, the user may decide their desired tradeoff between performance and accuracy, and choose the subset size in the supernet training accordingly. For example, with a budget of 1 GPU hour, the best approach is to use a 10% subset of the training data for the supernet training and projection, but with a budget of 2.5 GPU hours, the best approach is to use a 10% subset of the training data for the supernet and the full training data for the projection.

Table 9: Performance of one-shot NAS algorithms on DARTS search space CIFAR-10.

Performance on DARTS CIFAR-10			
Algorithm	Test accuracy	GPU hours	%Data used
DARTS-PT	97.17 (97.39)	20.59	100
DARTS-PT-ENTROPY	96.68 $\pm$ 0.26	3.40	10
DARTS-PT-RAND	97.01 $\pm$ 0.32	2.35	10
DARTS-PT-FL	96.91 $\pm$ 0.15	4.00	10
ADAPTIVE-NAS	96.73 $\pm$ 0.29	2.75	10
ADAPTIVE-NAS	96.97 $\pm$ 0.24	4.50	20

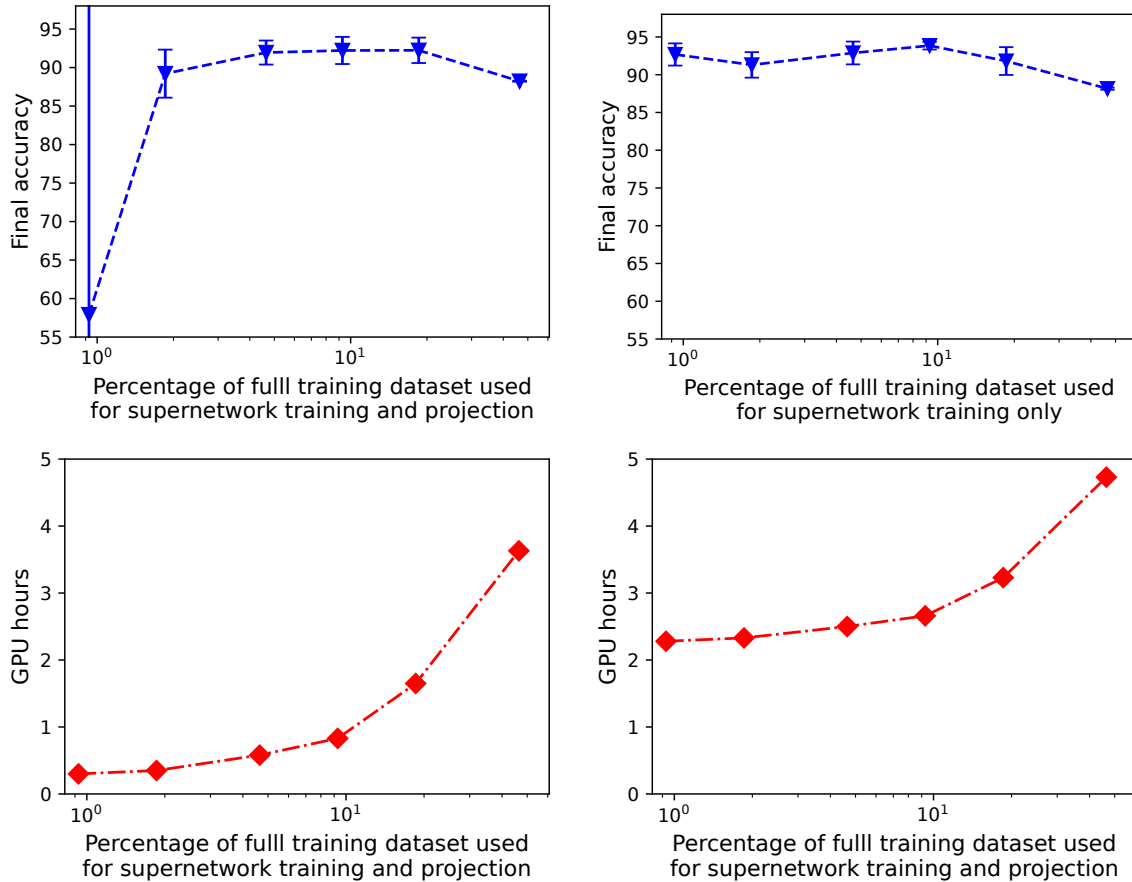


Figure 2: Performance and runtime of ADAPTIVE-NAS varies as the percentage of training data increases. (Left) The supernet training and projection step are given a percentage of the full training dataset. (Right) The supernet training is given a percentage, while the projection step is given the full training dataset.