
NAS-Bench-360: Benchmarking Diverse Tasks for Neural Architecture Search

Renbo Tu
Carnegie Mellon University
renbo@cmu.edu

Mikhail Khodak
Carnegie Mellon University
khodak@cmu.edu

Nicholas Roberts
Carnegie Mellon University
ncrobert@andrew.cmu.edu

Maria-Florina Balcan
Carnegie Mellon University
ninamf@cs.cmu.edu

Ameet Talwalkar
Carnegie Mellon University and Determined AI
talwalkar@cmu.edu

Abstract

1 Most existing neural architecture search (NAS) benchmarks and algorithms priori-
2 tize performance on well-studied tasks, focusing on computer vision datasets such
3 as CIFAR and ImageNet. However, the applicability of NAS approaches in other
4 areas is not adequately understood. In this paper, we present **NAS-Bench-360**,
5 a benchmark suite for evaluating state-of-the-art NAS methods on less-explored
6 datasets. To do this, we organize a diverse array of tasks, from classification of
7 simple deformations of natural images to predicting protein folding and partial
8 differential equation (PDE) solving. Our evaluation pipeline compares architecture
9 search spaces of different flavors, and reveals varying performance on different
10 tasks, providing baselines for further use. All data and reproducible evaluation
11 code are open-source and publicly available. The results of our evaluation show
12 that current state-of-the-art NAS methods often struggle to compete with sim-
13 ple baselines and human-designed architectures on the majority of tasks in our
14 benchmark. At the same time, they can be quite effective on a few individual,
15 understudied tasks. This demonstrates the importance of evaluation on diverse
16 tasks to better understand the usefulness of different approaches to architecture
17 search and automation.

18 1 Introduction

19 Neural architecture search (NAS) aims to automate the design of deep neural networks, ensuring
20 performance on par with hand-crafted architectures while reducing human labor devoted to tedious
21 architecture tuning [8]. With the growing number of application areas of ML, and thus of use-cases
22 for automating it, NAS has experienced an intense amount of study, with significant progress in
23 search space design [3, 20, 33], search efficiency [22], and search algorithms [16, 28, 29]. While the
24 use of NAS techniques may be especially impactful in under-explored or under-resourced domains
25 where less expert help is available, the field has largely been dominated by methods designed for
26 and evaluated on benchmarks in computer vision [6, 20, 30]. There have been a few recent efforts
27 to diversify these benchmarks to settings such as vision-based transfer learning [7] and speech and
28 language processing [13, 21]; however, evaluating NAS methods on such well-studied tasks using
29 traditional, domain-specific search spaces does not give a good indication of their utility on more
30 far-afield applications, which have often necessitated the design of custom neural operations [4, 19].

31 We aim to rectify this issue by introducing a suite of diverse benchmark tasks drawn from various
32 data domains that we collectively call **NAS-Bench-360**. This benchmark consists of an organized

33 setup of five suitable datasets that can both (a) be evaluated in a unified way using existing NAS
34 approaches and (b) come from a variety of different application areas, including numerical analysis,
35 organic chemistry, and medical imaging. We also include standard image classification evaluations as
36 a point of comparison, as many new methods continue to be designed for such tasks.

37 Following our construction of this benchmark, we evaluate three different NAS approaches, each
38 characterized by an architecture (search space, search algorithm) pair, and compare the results to
39 expert-driven domain-specific architecture design. As a baseline comparator, the first approach
40 uses a singleton architecture, the Wide ResNet (WRN) [31], as the search space, paired with a
41 hyperparameter tuning algorithm to adjust the training procedure for each task. The other two search
42 spaces are well-studied in modern NAS: DARTS [20] and DenseNAS [9], and we pair them with
43 their respective best-performing search methods. We find that the modern NAS approaches struggle
44 to beat even the simple WRN comparator on the majority of tasks in the benchmark. On two of
45 the tasks—classifying electromyography signals and solving partial differential equations—NAS
46 methods do significantly worse. NAS lags even further behind when we include domain-specific
47 expert-designed architectures, where it lags far behind on even CIFAR-100 when disallowing extra
48 augmentation or pre-training on ImageNet [25]. On the other hand, DARTS cells perform relatively
49 well on two tasks that *a priori* seem more challenging: spherical image classification and protein-
50 distance prediction. These observations and other empirical insights demonstrate the necessity of a
51 benchmark that provides a diverse array of data domains for evaluating NAS methods. Our evaluation
52 results also serve as a baseline for comparison in future development of NAS.

53 To ensure the availability and impact of this benchmark, the associated datasets and eval-
54 uation pipelines will remain open-source and accessible at [https://rtu715.github.io/
55 NAS-Bench-360/](https://rtu715.github.io/NAS-Bench-360/). Reproducibility is assured from open-sourcing all relevant code for the end-to-
56 end procedure, including data processing, architecture search, model retraining, and hyper-parameter
57 tuning frameworks.

58 **2 Related Work**

59 Benchmarks have been very important to the development of NAS in recent years. This includes
60 standard evaluation datasets and protocols, of which the most popular are the CIFAR-10 and ImageNet
61 routines used by DARTS [20]. Another important type of benchmark has been tabular benchmarks
62 such as NAS-Bench-101 [30], NAS-Bench-201 [6], and NAS-Bench-1Shot1 [32]; these benchmarks
63 exhaustively evaluate all architectures in their search spaces, which is made computationally feasible
64 by defining simple searched cells. Consequently, these benchmark cells are less expressive than the
65 DARTS cell [20], often regarded as the most powerful search space in the cell-based regime. Notably,
66 our benchmark is *not* a tabular benchmark, i.e. we do *not* evaluate every architecture from a fixed
67 search space; rather, the focus is on the organization of a suite of tasks to evaluate both NAS methods
68 and search spaces, which would necessarily be restricted if we first fixed a search space to construct a
69 tabular benchmark from.

70 While NAS methods and benchmarks have generally been focused on computer vision, recent work
71 such as AutoML-Zero [23] and XD operations [24] has started moving towards a more generically
72 applicable set of tools for AutoML. However, even more recent benchmarks that do go beyond the
73 most popular vision datasets have continued to focus on well-studied tasks, including vision-based
74 transfer learning [7], speech recognition [21], and natural language processing [13]. Our aim is to
75 go beyond such areas in order to evaluate the potential of NAS to automate the application of ML
76 in truly under-explored domains. One analogous work to ours in the field of meta-learning is the
77 Meta-Dataset benchmark of few-shot tasks [27], which similarly aimed to establish a wide-ranging
78 set of evaluations for that field.

79 **3 NAS-Bench-360: A Suite of Diverse and Practical Tasks**

80 In this section, we introduce the NAS setting being targeted by our benchmark, our motivation for
81 organizing a new set of diverse tasks as a NAS evaluation suite, and our task-selection methodology.
82 We report evaluations of specific algorithms on this new benchmark in the next section.

83 3.1 Neural Architecture Search: Problem Formulation and Baselines

84 For completeness and clarity, we first formally discuss the architecture search problem itself, starting
85 with the extended hypothesis class formulation [16]. Here the goal is to use a dataset of points $x \in \mathcal{X}$
86 to find parameters $\mathbf{w} \in \mathcal{W}$ and $a \in \mathcal{A}$ of a parameterized function $f_{\mathbf{w},a} : \mathcal{X} \mapsto \mathbb{R}_{\geq 0}$ that minimize
87 the expectation $\mathbb{E}_{x \sim \mathcal{D}} f_{\mathbf{w},a}(x)$ for some test distribution \mathcal{D} over \mathcal{X} ; here \mathcal{X} is the input space, \mathcal{W} is
88 the space of model weights, and \mathcal{A} is the set of architectures. For generality, we do not require the
89 training points to be drawn from \mathcal{D} to allow for domain adaptation, as is the case for one of our tasks,
90 and we do not require the loss to be supervised. Note also that the goal here does not depend on the
91 issue of computational or memory efficiency, which we do not focus on in our evaluations; there our
92 restriction is only that the entire pipeline can be run on an NVIDIA V100 GPU.

93 Notably, this formulation makes no distinction between the model weights \mathbf{w} and architectures a ,
94 treating both as parameters of a larger model. Indeed, the goal of NAS may be seen as similar to
95 model design, except now we include the design of an (often-discrete) *architecture space* \mathcal{A} such that
96 it is easy to find an architecture $a \in \mathcal{A}$ and model weights $\mathbf{w} \in \mathcal{W}$ whose test loss $\mathbb{E}_{\mathcal{D}} f_{\mathbf{w},a}$ is low
97 using a search algorithm. This can be done in a one-shot manner—simultaneously optimizing a and
98 \mathbf{w} —or using the standard approach of first finding an architecture a and then keeping it fixed while
99 training model weights \mathbf{w} for it using a pre-specified algorithm such as tuned stochastic gradient
100 descent (SGD).

101 This formulation also includes non-NAS methods by allowing the architecture search space to be a
102 singleton. When the sole architecture is a standard and common network such as WRN [31], this yields
103 a natural baseline with an algorithm searching for training hyperparameters, not architectures. On the
104 other hand, when \mathcal{A} contains a single domain-specific architecture, such as a spherical convolutional
105 neural network (CNN) [4], it yields the “human baseline” competitor approach without search. For
106 our empirical investigation, we compare the performance of state-of-the-art NAS approaches against
107 that of the two singleton baselines.

108 3.2 Motivation and Task Selection Methodology

109 Curating a diverse, practical set of tasks for the study of NAS is our primary motivation behind this
110 work. We observe that past NAS benchmarks focused on the creation of larger search spaces and
111 more sophisticated search methods for neural networks. However, the utility of these search spaces
112 and methods are only evaluated on canonical computer vision datasets. Whether these new methods
113 can improve upon non-NAS baselines remains an open question. This calls for the introduction of
114 new datasets lest NAS research overfits to the biases of CIFAR-10 and ImageNet. By identifying
115 these possible biases, future directions in NAS research can be better primed to suit the needs of
116 practitioners, thereby incentivizing the deployment of NAS techniques on real applications.

117 NAS-Bench-360 comprises tasks from existing datasets their variants as summarized in Table 1. This
118 work focuses exclusively on datasets with 2d input data including images, wave spectra, differential
119 equations, and protein sequence features. Although in practice neural networks are employed to
120 analyze different data modalities, the most well-studied NAS approaches only accept 2d inputs
121 and therefore we study tasks within this scope. During the selection of tasks, breadth is our main
122 consideration. First, we formalize the categorization of tasks into **point prediction (point)** and
123 **dense prediction (dense)** [26], respectively referring to tasks with scalar outputs and 2d matrix
124 outputs. In other words, point prediction tasks are classification tasks, and dense prediction tasks are
125 element-wise prediction tasks, which is a specific form of regression. The heavy bias of previous
126 NAS research towards point prediction tasks motivates the inclusion of dense prediction tasks in our
127 benchmark. Second, breadth is achieved by selecting tasks from various subjects and applications of
128 deep learning, where introducing NAS could improve upon the performance of handcrafted neural
129 networks.

130 3.3 List of Tasks from Diverse Data Sources

131 In lieu of providing raw data, we perform data pre-processing locally and store the processed data on
132 a public Amazon Web Service’s S3 data bucket with download links available on our website. Our
133 data treatment largely follows the procedure defined by the researchers who provided them. This
134 would enhance the reproducibility of results by ensuring the uniformity of input data for different
135 pipelines. Specific pre-processing and augmentation steps are described below.

Table 1: Information of tasks in NAS-Bench-360

Task name	Dataset size	Type	Learning objective	New to NAS
CIFAR-100	60K	Point	Classify natural images into 100 classes	
Spherical	60K	Point	Classify spherically projected images into 100 classes	✓
NinaPro	3956	Point	Classify sEMG signals into 18 classes corresponding to hand gestures	✓
Darcy Flow	1100	Dense	Predict the final state of a fluid from its initial conditions	✓
PSICOV	3606	Dense	Predict pairwise distances between residuals from 2d protein sequence features	✓

136 3.3.1 CIFAR-100: Standard Image Classification

137 As a starting point of comparison to existing benchmarks, we include the **CIFAR-100** task [14], which
 138 contains RGB images from natural settings to be classified into 100 fine-grained categories. CIFAR-
 139 100 is preferred over CIFAR-10 because it is more challenging and suffers less from over-fitting in
 140 previous research.

141 **Data pre-processing:** while the 10,000 testing images are kept aside only for evaluating architec-
 142 tures, the 50,000 training images are randomly partitioned into 40,000 for architecture search and
 143 10,000 for validation. On all of the 50,000 training images, we apply standard CIFAR augmentations
 144 including random crops and horizontal flipping, and finally normalize them using a pre-calculated
 145 mean and standard deviation of this set. On the 10,000 testing images, we only apply normalization
 146 with the same constants.

147 3.3.2 Spherical: Classifying Spherically Projected CIFAR-100 Images

148 To test NAS methods applied to natural-image-like data, we consider the task of classifying spherical
 149 projections of the CIFAR-100 images, which we call the **Spherical** task. In addition to scientific
 150 interest, spherical image data is also present in a variety of applications, such as omnidirectional
 151 vision in robotics and weather modeling in meteorology, as sensors usually produce distorted image
 152 signals in real-life settings. To create a spherical variant of CIFAR, we project the planar signals of
 153 the CIFAR images to the northern hemisphere and add a random rotation to produce spherical signals
 154 for each individual channel following the procedure specified in [4]. The resulting images are 60*60
 155 pixels with RGB channels.

156 **Data pre-processing:** with the same split ratios CIFAR-100, the generated spherical image data is
 157 directly used for training and evaluation without data augmentation and pre-processing.

158 3.3.3 NinaPro: Classifying Electromyography Signals

159 Our final classification task, **NinaPro**, moves away from the image domain to classify hand gestures
 160 indicated by electromyography signals. For this, we use a subset of the NinaPro DB5 dataset [2]
 161 in which two thalamic Myo armbands collect EMG signals from 10 test individuals who hold 18
 162 different hand gestures to be classified. These armbands leverage data from muscle movement, which
 163 is collected using electrodes in the form of wave signals. Each wave signal is then sampled using a
 164 wavelength and frequency prescribed in [5] to produce 2d signals.

165 **Data pre-processing:** Containing less than 4,000 samples, the data is comprised of single-channel
 166 signals with an irregular shape of 16*52 pixels. This task also differs from CIFAR for its class
 167 imbalance, as over 65% of all gestures are the neutral position. We split the data using the same ratio
 168 as CIFAR, resulting in 2638 samples for training, and 659 samples for validation and testing each.
 169 No additional pre-processing is performed.

170 3.3.4 Darcy Flow: Solving Partial Differential Equations (PDEs)

171 Our first regression task, **Darcy Flow**, focuses on learning a map from the initial conditions of a PDE
172 to the solution at a later timestep. This application aims to replace traditional solvers with learned
173 neural networks, which can output a result in a single forward pass. The input is a 2d grid specifying
174 the initial conditions of a fluid and the output is a 2d grid specifying the fluid state at a later time,
175 with the ground truth being the result computed by a traditional solver.

176 **Data pre-processing:** we use scripts provided by [19] to generate the PDEs and their solutions,
177 for a total of 900 data points for training, 100 for validation, and 100 for testing. All input data
178 is normalized with constants calculated on the training set before fed into the neural network and
179 de-normalized following an encode-decode scheme. The solutions, or labels, for the training set are
180 also encoded and decoded this way. The test labels are not processed. We report the mean square
181 error (MSE or ℓ_2).

182 3.3.5 PSICOV: Protein Distance Prediction

183 Our final task, **PSICOV**, studies the use of neural networks in the protein folding prediction pipeline,
184 which has recently received significant attention to the success of methods like AlphaFold [12]. While
185 the dataset and method they use are too large-scale for our purposes, we consider a smaller set of
186 protein structures to tackle the specific problem of inter-residual distance predictions outlined in [1].
187 2d large-scale features are extracted from protein sequences, resulting in input feature maps with a
188 massive number of channels. Correspondingly, the labels are pairwise-distance matrices with the
189 same spatial dimension.

190 **Data pre-processing:** we adopt the chosen subset of DeepCov proteins in [1], consisting of 3,456
191 proteins each with 128*128 feature maps across 57 channels. 100 proteins from this set are used for
192 validation and the rest for training. Test data for final evaluation is gathered from another set of 150
193 proteins, PSICOV. Since these produce feature maps that are larger (512*512), we run the prediction
194 network over all of its non-overlapping 128*128 patches. The evaluation metric is mean absolute
195 error (MAE or ℓ_1) computed on distances below 8 Å, referred to as MAE₈.

196 3.4 Ethics and Responsible Use

197 Within our array of tasks, the only dataset containing human-derived data is NinaPro. Our chosen
198 subset of NinaPro contains only muscle movement data from 10 healthy individuals, without any
199 exposure of personal information from clinical data. The original experiments to acquire NinaPro
200 data are approved by the ethics commission of the state of Valais, Switzerland [2]. For other datasets,
201 we have listed the data licenses in the appendix for responsible usages of data. While we do not
202 view the specific datasets we use in this benchmark as potential candidates for misuse, the broader
203 goal of applying NAS to new domains comes with inherent risks that may require mitigation on an
204 application-by-application basis.

205 4 Using NAS-Bench-360 to Study Architecture Search Methods

206 Having detailed our construction of NAS-Bench-360, we now demonstrate its usefulness on (a)
207 comparing and evaluating state-of-the-art architecture search methods on powerful search spaces and
208 (b) discovering new insights on their performance on under-explored domains. In this section, we
209 first specify the different NAS algorithms and baselines we compare, followed by the experimental
210 and reproducibility setup we follow. Finally, we report our main comparisons and analyze the results.

211 4.1 Baselines and Search Procedures

212 From the discussion in Section 3, the two non-NAS baseline methods we consider—applying a tuned
213 WRN to all tasks and using a fixed, domain-specific architecture—can be viewed via the NAS setup
214 as having a singleton architecture search space. As for NAS algorithms themselves, we focus on
215 two well-known paradigms for search: cell-based NAS (using DARTS [20]) and macro NAS (using
216 DenseNAS [9]). We detail these four approaches below.

217 **Wide ResNet with Hyperparameter Tuning** The residual network (ResNet) and its derivative
218 architectures are canonical for classic computer vision, and we investigate their ability to generalize to
219 our selection of tasks. A more powerful adaptation of ResNet, the Wide ResNet [31] is chosen as the
220 backbone architecture. For automated training, we wrap the training procedure with a hyperparameter
221 tuning algorithm, ASHA [15], an asynchronous version of Hyperband [18]. Given a range for each
222 hyperparameter, either discrete or continuous, ASHA uniformly samples configurations and uses
223 brackets of elimination: at each round, each configuration is trained for some epochs, before the
224 algorithm selects the best-performing portion based on validation metrics. Since we use the Wide
225 ResNet backbone for all tasks, our tuning budget is fixed and uniform.

226 **Expert-Designed Networks** We also include expert-driven design of architectures in specific
227 domains as a more rigorous comparator for NAS methods on our tasks. Frequently this includes
228 not only hand-designed topologies and operation patterns but custom neural operations themselves,
229 which are often crucial for success on domains beyond computer vision. Below we briefly summarize
230 the architectures chosen for each task.

- 231 • **CIFAR-100:** While this task is very heavily studied and one can achieve very high accuracies
232 using optimization tricks and transfer from ImageNet, we restrict our selection to existing
233 results that use only the simple (standard) data augmentation we allow for the evaluation
234 phase. Here the best result found is using DenseNet-BC [10].
- 235 • **Spherical:** This task is often regarded as a canonical example where a specific neural
236 operation, specifically spherical convolutions, are the “right” operation to substitute for
237 the convolution due to data-specific properties. Our result is from a wide variation of the
238 spherical CNN in [4], with a max width of 256 channels from 64.
- 239 • **NinaPro:** As the original paper studying NinaPro used fairly weak networks that achieve a
240 much higher error, here we simply report the performance of our tuned WRN baseline.
- 241 • **Darcy Flow:** Here we report the performance of a four-layer network that replaces convolu-
242 tions with Fourier Neural Operators (FNOs) [19], which were specially designed solving
243 partial differential equations. Note that our reproduced result attains slightly better MSE
244 than the numbers reported by the authors.
- 245 • **PSICOV:** We report the reproduced performance of the ResNet-256 network used by the
246 PDNET, a deeper, narrow, and dilated version of the standard ResNet used for ImageNet;
247 note our reproduction attains much better MAE₈ than the authors report [1].

248 **Cell-based Search Using DARTS** The first state-of-the-art NAS paradigm within our consideration
249 is cell-based NAS. Cell-based methods first search for a genotype, which is a cell containing neural
250 operations such as convolution and pooling. During evaluation, a neural network is constructed by
251 replicating the searched cell and stacking them together. The most popular search space for this
252 approach is the one used by the DARTS space [20], consisting of assigning one of eight operations to
253 six edges in two types of cells: “normal” cells preserve the shape of the input representation while
254 “reduction” cells downsample it. Note that for the dense tasks we do not use the reduction cell so as
255 to not introduce a bottleneck.

256 Finally, to adhere to standard ML practices we do *not* adapt the standard DARTS pipeline, which
257 uses test performance to select from multiple random seeds. This, in addition to not using other
258 evaluation-time enhancements such—specifically auxiliary towers and the cutout data augmentation—
259 leads to lower performance on CIFAR-100 than is reported in the literature. As this search space has
260 been heavily studied since its introduction, we use as a search routine a recent approach—GAEA
261 PC-DARTS—that achieves some of the best-known results on CIFAR-10 and ImageNet for this
262 benchmark [16].

263 **Macro NAS Using DenseNAS** The second NAS paradigm we consider is macro NAS. Instead of
264 building from a fixed cell, macro NAS requires the specification of a super network with different
265 inter-connected network blocks. These blocks and connections are then pruned during the search
266 phase to construct the output neural net for evaluation. For this benchmark, we also choose a recent
267 search space in this NAS paradigm, DenseNAS [9], which similarly to the DARTS space has near
268 state-of-the-art results on ImageNet.

Table 2: Comparing NAS methods with baseline and expert-designed methods on NAS-Bench-360. All automated results (WRN, DenseNAS, and GAEA PC-DARTS) are averages of three random seeds. See Appendix for standard deviations.

Search space	Search method	CIFAR-100 (0-1 err.)	Spherical (0-1 err.)	NinaPro (0-1 err.)	Darcy Flow MSE	PSICOV MAE ₈
WRN baseline expert design*	ASHA	24.89	88.45	6.88	0.041	5.71
	hand-tuning	17.17	64.42	6.88	0.0096	3.50
DenseNAS-R1	DenseNAS	27.44	72.99	10.17	0.10	3.84
DARTS Cell	GAEA PC-DARTS	24.19	52.90	11.43	0.056	2.80

* Chosen according to best-effort literature search and implementation; c.f. Section 4.1.

269 DenseNAS searches for architectures with densely-connected, customizable routing blocks to emulate
 270 DenseNet [10]. In our experiments, we use the ResNet-based search space, DenseNAS-R1, with all of
 271 WRN’s neural operations for better comparison with the baseline backbone. For point tasks and dense
 272 tasks, we adapt two super networks from the one used for ImageNet as inputs to the search algorithm.
 273 The super network for dense tasks maintains the same spatial dimensions without downsampling to
 274 avoid bottlenecks, and we use a lower learning rate for evaluating architectures on dense tasks to
 275 prevent divergence. Other training and evaluation procedures are identical to those in the original
 276 paper and uniform across all tasks.

277 4.2 Experimental Setup

278 Our main experiments consist of 3 evaluation trials for every combination of method and task, fixing
 279 one random seed for each trial. We present these results in Table 2 and discuss the specific procedure,
 280 reproducibility, and extension experiments in the following subsections.

281 **Using validation data** For best practices in NAS, we argue for the separation of the final testing
 282 set and the validation set, which is specifically for selecting neural architectures and hyperparameters.
 283 After this process, we combine training and validation data to perform retraining and evaluation on
 284 the test set. This result is reported as final and is not used in any way to further optimize the model.

285 **Hyperparameter tuning** In experiments with hyperparameter tuning, we consistently use the
 286 same hyperparameter ranges and fix the tuning budget, in terms of the number of configurations and
 287 maximum training epochs, across all tasks. The tuning budget is selected to be 2.5 to 3 times the
 288 backbone training time. This is to eliminate inductive biases for specific tasks. Details on the tuning
 289 procedure are in the appendix.

290 **Software and hardware** We adopt the free, open-source software *Determined*¹ for experiment
 291 management, hyperparameter tuning, AWS cloud deployment with docker containers. All experi-
 292 ments are performed on a single p3.2xlarge instance with one Nvidia V100 GPU. The computation
 293 cost in GPU hours of individual experiments using this setup can be found in the appendix.

294 **Reproducibility** The following measures in our experimental pipeline are taken to ensure the
 295 reproducibility of our results:

- 296 1. We perform most data pre-processing steps beforehand and store the processed data in the
 297 cloud for download. A data splitting scheme, once randomly selected, is then fixed for all
 298 experiments on that task, i.e. the same training, validation, and testing sets fed into the
 299 dataloader are always the same.
- 300 2. Experimentation code is always executed in a fixed docker container using a pre-built docker
 301 image on Docker Hub. This guarantees a uniform execution environment and saves users
 302 from the manual labor of configuring dependencies.

¹GitHub repository: <https://github.com/determined-ai/determined>

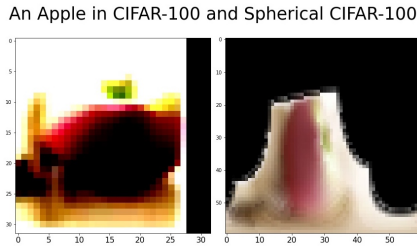


Figure 1: Comparison of the same CIFAR-100 image before and after the spherical transformation.

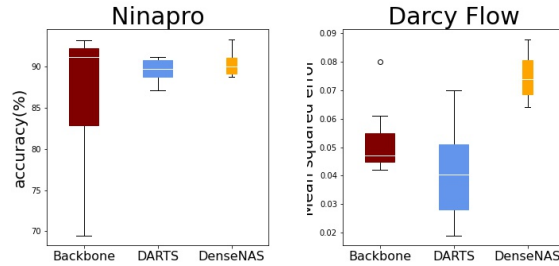


Figure 2: Distribution of random architectures and hyperparameters' performance on NinaPro and Darcy Flow.

Table 3: Experiment runtimes of NAS-Bench-360 (GPU hours)

Task	GAEA PC-DARTS	DenseNAS	WRN
CIFAR-100	33	2	8
Spherical CIFAR-100	39	2.5	8.5
NinaPro	2	0.5	1
Darcy Flow	15	0.5	2
PSICOV	59.5	23	61

- 303 3. Via the specification of a random seed, *Determined* controls several important sources of
 304 randomness during code execution, including hyperparameter sampling and training data
 305 shuffling.
- 306 4. During training, we always validate on the full validation set, not on a mini-batch, to avoid
 307 stochasticity in the results.

308 4.3 Comparing NAS Approaches Using NAS-Bench-360

309 **Generalization of NAS to other domains** Our experiments demonstrate that state-of-the-art NAS
 310 approaches in classic vision are unable to outperform human-designed neural networks on 3 out
 311 of 5 tasks in NAS-Bench-360. They do especially poorly on the Darcy Flow task and fall short
 312 of matching both non-NAS comparators by a large margin. Perhaps most surprisingly, neither the
 313 DARTS space nor DenseNAS, both very recent search spaces with strong results on ImageNet
 314 (and CIFAR-10 for the former) are able to outperform the reported performance of a fairly basic
 315 architecture (DenseNet) on CIFAR-100; this is especially interesting as DenseNAS was built around
 316 this architecture. Overall, our results suggest that modern NAS, despite its promise to automate deep
 317 learning, is not yet well-equipped to handle its various domains of applications studied in this paper.
 318 These empirical results also serve as new baselines for comparison in future research to extend NAS
 319 approaches to generalize to new areas.

320 **Computational cost** In some time-sensitive applications of NAS, both efficiency and perfor-
 321 mance are criteria for NAS method selection. Our choice of methods exemplifies a tradeoff
 322 between these two factors. As a more computationally heavy method, GAEA PC-DARTS
 323 beats the more lightweight DenseNAS on most of the tasks except for NinaPro, where they
 324 achieve similar accuracies. On certain tasks, such as NinaPro and PSICOV, DenseNAS would
 325 be the more cost-effective option than GAEA PC-DARTS to have decent performance on par
 326 with handcrafted neural architectures. Note that the computation cost of the WRN baseline can
 327 vary due to randomness inherent in ASHA's asynchrony. We report all experiment runtimes in Table 3.
 328

329 **CIFAR-100 vs. Spherical** The Spherical task can be directly compared to CIFAR-100 to assess
 330 how well NAS methods could handle image distortions. With the same setup across tasks, both

Table 4: ℓ_1 error of supernet and searched architectures (discretized) on grid tasks

Task	DARTS		DenseNAS	
	Supernet	Discretized	Supernet	Discretized
Darcy Flow	0.031 ± 0.001	0.057 ± 0.012	0.041 ± 0.002	0.10 ± 0.010
PSICOV	3.87 ± 0.12	2.80 ± 0.057	7.96 ± 0.20	3.84 ± 0.15

331 the DARTS space and DenseNAS have reasonably good numbers on CIFAR-100, but their results
 332 significantly deteriorate on the spherical variant. Both obtain much worse error when the images
 333 are spherically projected, but a much larger gap emerges between the two methods, with DenseNAS
 334 performing quite badly. On the other hand, the searched DARTS Cell not only performs 20-36% better
 335 than the other convolutional approaches but even beats our best-effort adaptation of the spherical
 336 CNN approach to this task [4], in which we expanded the size of that network. This is surprising
 337 because spherical convolutions were designed specifically for such data. We believe these results
 338 indicate that the spherical dataset may be a useful but simple way for distinguishing NAS approaches
 339 when they are overfitting to standard computer vision domains; Figure 1 provides an example of the
 340 distortion.

341 **WRN as a baseline** Viewing the WRN baseline as a singleton architecture search space, we
 342 compare this baseline to more sophisticated NAS search spaces. On our set of new tasks, NAS does
 343 not perform better than Wide ResNet with hyperparameter tuning on CIFAR-100, NinaPro, and
 344 Darcy Flow but excels on the rest. Hyperparameter optimization can boost the backbone performance
 345 considerably to rival the performance of NAS methods. Most non-Bayesian hyperparameter tuning
 346 algorithms, such as random search [17], population-based training [11], and Hyperband [18], are
 347 also straightforward to apply with any neural network backbone. Therefore, we argue for the use of
 348 hyperparameter-tuned backbones to assess the effectiveness of NAS approaches and encourage their
 349 inclusion in NAS benchmarks.

350 4.4 In-Depth Studies Using NAS-Bench-360

351 **Supernet performance on grid tasks** During architecture search, our NAS methods on the DARTS
 352 and DenseNAS search spaces train the supernets to find optimal neural operations on the validation
 353 set. Surprisingly, the validation error of the supernet is sometimes lower than that of the final searched,
 354 discretized neural network. Therefore, we evaluate the supernet of DARTS and DenseNAS on the
 355 testing set, and we compare its performance with that of the final neural network in Table 4. The
 356 supernet outperforms the final network on Darcy Flow for both methods, but the reverse is true for
 357 the PSICOV task and all point tasks. The supernet is not in the search space and so we report the
 358 discretized result; nevertheless, this fact suggests that performance on a task like Darcy Flow might
 359 benefit from a better search space.

360 **Evaluating random architectures and hyperparameters** The power of an architecture or hyper-
 361 parameter space can also be characterized by the performance of its random elements. We assess both
 362 the average and variance of the results. To do this, we randomly sample 8 network architectures each
 363 from the search spaces of DARTS and DenseNAS, and we test their performance on the NinaPro
 364 and Darcy Flow tasks, one for classification and the other for regression. For comparison, we also
 365 randomly sample 8 hyperparameter configurations to train the backbone Wide-ResNet in Figure 2.
 366 While rather successful on NinaPro, the random architectures have a high average error and vary in
 367 performance on the Darcy Flow task. Random hyperparameters are more unstable on NinaPro, but its
 368 median performance is better than NAS.

369 **Utility of hyperparameter tuning** The final experiment examines whether hyperparameter tuning
 370 improves the performance of WRN on various tasks. During hyperparameter search, we compare the
 371 validation metrics of training using default hyperparameters and using tuned ones from ASHA to
 372 select final hyperparameters for retraining. Despite the small tuning budget allocated to ASHA, tuned
 373 hyperparameters could outperform the default setting on all tasks except for CIFAR-100. Our results
 374 suggest that wide ResNet’s standard set of hyperparameters are only optimized for conventional image
 375 classification. On other tasks, hyperparameter optimization is helpful for boosting performance.

376 5 Conclusion

377 **NAS-Bench-360** is a benchmarking suite with a novel, diverse set of tasks. The tasks are derived from
378 various fields of academic research, leading to different potential applications. Our selection of NAS
379 approaches achieves state-of-the-art performances on most tasks, which points to new possibilities
380 of incorporating NAS into new research domains. All datasets and reproducible experiment code
381 are open-sourced, and we welcome researchers to use these tasks and further iterate on them with
382 new NAS methods. Finally, a possible extension to generalize this set of tasks is datasets with 1d or
383 3d inputs, such as audio. We hope our work can encourage the NAS community to move towards
384 tackling more diverse problems in the real world.

385 References

- 386 [1] Badri Adhikari. A fully open-source framework for deep learning protein real-valued distances.
387 *Scientific reports*, 10(1):1–10, 2020.
- 388 [2] Manfredo Atzori, Arjan Gijsberts, Simone Heynen, Anne-Gabrielle Mittaz Hager, Olivier
389 Deriaz, Patrick Van Der Smagt, Claudio Castellini, Barbara Caputo, and Henning Müller.
390 Building the ninapro database: A resource for the biorobotics community. In *2012 4th IEEE*
391 *RAS & EMBS International Conference on Biomedical Robotics and Biomechatronics (BioRob)*,
392 pages 1258–1265. IEEE, 2012.
- 393 [3] Han Cai, Ligeng Zhu, and Song Han. ProxylessNAS: Direct neural architecture search on
394 target task and hardware. In *Proceedings of the 7th International Conference on Learning*
395 *Representations*, 2019.
- 396 [4] Taco S. Cohen, Mario Geiger, Jonas Koehler, and Max Welling. Spherical CNNs. In *Proceedings*
397 *of the 6th International Conference on Learning Representations*, 2018.
- 398 [5] Ulysse Côté-Allard, Cheikh Latyr Fall, Alexandre Drouin, Alexandre Campeau-Lecours, Clément
399 Gosselin, Kyrre Glette, François Laviolette, and Benoit Gosselin. Deep learning for
400 electromyographic hand gesture signal classification using transfer learning. *IEEE Transactions*
401 *on Neural Systems and Rehabilitation Engineering*, 27(4):760–771, 2019.
- 402 [6] Xuanyi Dong and Yi Yang. NAS-Bench-201: Extending the scope of reproducible neural archi-
403 tecture search. In *Proceedings of the 8th International Conference on Learning Representations*,
404 2020.
- 405 [7] Yawen Duan, Xin Chen, Hang Xu, Zewei Chen, Xiaodan Liang, Tong Zhang, and Zhenguo
406 Li. TransNAS-Bench-101: Improving transferability and generalizability of cross-task neural
407 architecture search. In *Proceedings of the IEEE Conference on Computer Vision and Pattern*
408 *Recognition*, 2021.
- 409 [8] Thomas Elsken, Jan Hendrik Metzen, and Frank Hutter. Neural architecture search: A survey.
410 *Journal of Machine Learning Research*, 20(55):1–21, 2019.
- 411 [9] Jiemin Fang, Yuzhu Sun, Qian Zhang, Yuan Li, Wenyu Liu, and Xinggang Wang. Densely
412 connected search space for more flexible neural architecture search. In *Proceedings of the IEEE*
413 *Conference on Computer Vision and Pattern Recognition*, 2020.
- 414 [10] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger. Densely connected
415 convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern*
416 *recognition*, pages 4700–4708, 2017.
- 417 [11] Max Jaderberg, Valentin Dalibard, Simon Osindero, Wojciech M Czarnecki, Jeff Donahue, Ali
418 Razavi, Oriol Vinyals, Tim Green, Iain Dunning, Karen Simonyan, et al. Population based
419 training of neural networks. *arXiv preprint arXiv:1711.09846*, 2017.
- 420 [12] John Jumper, Richard Evans, Alexander Pritzel, Tim Green, Michael Figurnov, Kathryn Tunya-
421 suvunakool, Olaf Ronneberger, Russ Bates, Augustin Židek, Alex Bridgland, Clemens Meyer,
422 Simon A A Kohl, Anna Potapenko, Andrew J Ballard, Andrew Cowie, Bernardino Romera-
423 Paredes, Stanislav Nikolov, Rishub Jain, Jonas Adler, Trevor Back, Stig Petersen, David Reiman,

- 424 Martin Steinegger, Michalina Pacholska, David Silver, Oriol Vinyals, Andrew W Senior, Koray
425 Kavukcuoglu, Pushmeet Kohli, and Demis Hassabis. High accuracy protein structure prediction
426 using deep learning. In *Fourteenth Critical Assessment of Techniques for Protein Structure
427 Prediction (Abstract Book)*, 2020.
- 428 [13] Nikita Klyuchnikov, Ilya Trofimov, Ekaterina Artemova, Mikhail Salnikov, Maxim Fedorov,
429 and Evgeny Burnaev. NAS-Bench-NLP: Neural architecture search benchmark for natural
430 language processing. arXiv, 2020.
- 431 [14] Alex Krizhevsky. Learning multiple layers of features from tiny images. Technical report, 2009.
- 432 [15] Liam Li, Kevin Jamieson, Afshin Rostamizadeh, Ekaterina Gonina, Moritz Hardt, Benjamin
433 Recht, and Ameet Talwalkar. A system for massively parallel hyperparameter tuning. *arXiv
434 preprint arXiv:1810.05934*, 2018.
- 435 [16] Liam Li, Mikhail Khodak, Maria-Florina Balcan, and Ameet Talwalkar. Geometry-aware
436 gradient algorithms for neural architecture search. In *Proceedings of the 9th International
437 Conference on Learning Representations*, 2021.
- 438 [17] Liam Li and Ameet Talwalkar. Random search and reproducibility for neural architecture search.
439 In *Proceedings of the Conference on Uncertainty in Artificial Intelligence*, 2019.
- 440 [18] Lisha Li, Kevin G Jamieson, Giulia DeSalvo, Afshin Rostamizadeh, and Ameet Talwalkar.
441 Hyperband: Bandit-based configuration evaluation for hyperparameter optimization. In *ICLR
442 (Poster)*, 2017.
- 443 [19] Zongyi Li, Nikola Borislavov Kovachki, Kamyar Azizzadenesheli, Burigede Liu, Kaushik
444 Bhattacharya, Andrew Stuart, and Anima Anandkumar. Fourier neural operator for parametric
445 partial differential equations. In *Proceedings of the 9th International Conference on Learning
446 Representations*, 2021.
- 447 [20] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS: Differentiable architecture search.
448 In *Proceedings of the 7th International Conference on Learning Representations*, 2019.
- 449 [21] Abhinav Mehrotra, Alberto Gil, C. P. Ramos, Sourav Bhattacharya, Łukasz Dudziak, Ravichan-
450 der Vipera, Thomas Chau, Samin Ishtiaq, Mohamed S. Abdelfattah, and Nicholas D. Lane.
451 NAS-Bench-ASR: Reproducible neural architecture search for speech recognition. In *Proceed-
452 ings of the 8th International Conference on Learning Representations*, 2021.
- 453 [22] Hieu Pham, Melody Y. Guan, Barret Zoph, Quoc V. Le, and Jeff Dean. Efficient neural
454 architecture search via parameter sharing. In *Proceedings of the 35th International Conference
455 on Machine Learning*, 2018.
- 456 [23] Esteban Real, Chen Liang, David R. So, and Quoc V. Le. AutoML-Zero: Evolving machine
457 learning algorithms from scratch. In *Proceedings of the 37th International Conference on
458 Machine Learning*, 2020.
- 459 [24] Nicholas Roberts, Mikhail Khodak, Tri Dao, Liam Li, Chris Ré, and Ameet Talwalkar. Rethink-
460 ing neural operations for diverse tasks. arXiv, 2021.
- 461 [25] Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng
462 Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, Alexander C. Berg, and Li Fei-Fei.
463 ImageNet large scale visual recognition challenge. *International Journal of Computer Vision*,
464 115(3):211–252, 2015.
- 465 [26] Tom Sercu and Vaibhava Goel. Dense prediction on sequences with time-dilated convolutions
466 for speech recognition. *arXiv preprint arXiv:1611.09288*, 2016.
- 467 [27] Eleni Triantafillou, Tyler Zhu, Vincent Dumoulin, Pascal Lamblin, Utku Evci, Kelvin Xu, Ross
468 Goroshin, Carles Gelada, Kevin Swersky, Pierre-Antoine Manzagol, and Hugo Larochelle.
469 Meta-dataset: A dataset of datasets for learning to learn from few examples. In *Proceedings of
470 the 8th International Conference on Learning Representations*, 2020.

- 471 [28] Colin White, Willie Neiswanger, and Yash Savani. BANANAS: Bayesian optimization with
 472 neural architectures for neural architecture search. In *Proceedings of the 35th AAAI Conference*
 473 *on Artificial Intelligence*, 2021.
- 474 [29] Yuhui Xu, Lingxi Xie, Xiaopeng Zhang, Xin Chen, Guo-Jun Qi, Qi Tian, and Hongkai Xiong.
 475 PC-DARTS: Partial channel connections for memory-efficient architecture search. In *Proceed-*
 476 *ings of the 8th International Conference on Learning Representations*, 2020.
- 477 [30] Chris Ying, Aaron Klein, Eric Christiansen, Esteban Real, Kevin Murphy, and Frank Hutter.
 478 NAS-Bench-101: Towards reproducible neural architecture search. In *Proceedings of the 36th*
 479 *International Conference on Machine Learning*, 2019.
- 480 [31] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. In *Proceedings of the British*
 481 *Machine Vision Conference*, 2016.
- 482 [32] Arber Zela, Julien Siems, and Frank Hutter. NAS-Bench-1Shot1: Benchmarking and dissecting
 483 one-shot neural architecture search. In *Proceedings of the 8th International Conference on*
 484 *Learning Representations*, 2020.
- 485 [33] Barret Zoph, Vijay Vasudevan, Jonathon Shlens, and Quoc V. Le. Learning transferable
 486 architectures for scalable image recognition. In *Proceedings of the IEEE Conference on*
 487 *Computer Vision and Pattern Recognition*, 2018.

488 Checklist

- 489 1. For all authors...
- 490 (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s
 491 contributions and scope? [Yes]
- 492 (b) Did you describe the limitations of your work? [Yes] . See conclusion (section 5).
- 493 (c) Did you discuss any potential negative societal impacts of your work? [Yes] , in section
 494 3.4.
- 495 (d) Have you read the ethics review guidelines and ensured that your paper conforms to
 496 them? [Yes]
- 497 2. If you are including theoretical results...
- 498 (a) Did you state the full set of assumptions of all theoretical results? [N/A]
- 499 (b) Did you include complete proofs of all theoretical results? [N/A]
- 500 3. If you ran experiments...
- 501 (a) Did you include the code, data, and instructions needed to reproduce the main experi-
 502 mental results (either in the supplemental material or as a URL)? [Yes] . Instructions
 503 are described in the paper; code and data are available on our dedicated website.
- 504 (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they
 505 were chosen)? [Yes] , in both section 3 and appendix section A.
- 506 (c) Did you report error bars (e.g., with respect to the random seed after running experi-
 507 ments multiple times)? [Yes] , in the main results table and the appendix table.
- 508 (d) Did you include the total amount of compute and the type of resources used (e.g., type
 509 of GPUs, internal cluster, or cloud provider)? [Yes] , in section 4.2.
- 510 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
- 511 (a) If your work uses existing assets, did you cite the creators? [Yes]
- 512 (b) Did you mention the license of the assets? [Yes] , in appendix section B.
- 513 (c) Did you include any new assets either in the supplemental material or as a URL? [Yes]
 514 . New assets are on our website.
- 515 (d) Did you discuss whether and how consent was obtained from people whose data you’re
 516 using/curating? [N/A]
- 517 (e) Did you discuss whether the data you are using/curating contains personally identifiable
 518 information or offensive content? [Yes] , in section 3.4.

519

5. If you used crowdsourcing or conducted research with human subjects...

520

(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]

521

522

(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]

523

524

(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

525