



HARVARD

School of Engineering
and Applied Sciences

Investigating Differentiable Neural Architecture Search

Harvard Data Science Capstone (Fall 2019)
Final Presentation

Team

Michael S. Emanuel

Julien Laasri

Dylan Randle

Jiawei Zhuang

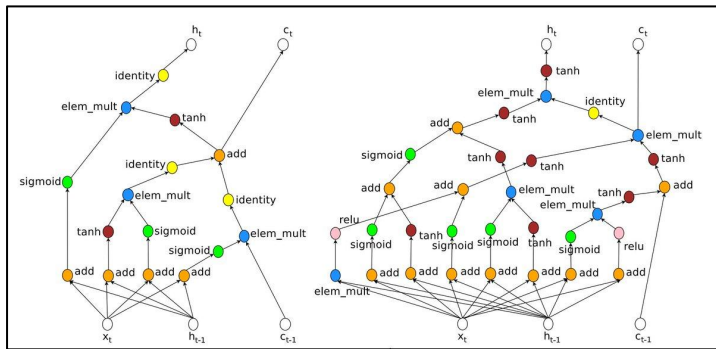
Outline

1. Problem & Motivation
 - a. Neural Architecture Search
 - b. Datasets
2. Introduction to DARTS (Differentiable Neural Architecture Search)
3. Results
 - a. MNIST
 - b. Graphene
 - c. Galaxy Zoo
 - d. Chest X-Ray
4. Conclusions & Future Work

Problem & Motivation

What is Neural Architecture Search?

- Deep learning frees us from feature engineering, but has led us to spend valuable time on **architecture engineering**
- Today: designed by experts
 - Labor-intensive
- Tomorrow: **Neural Architecture Search (NAS)**
 - Automatically find best architecture
- Interest in NAS is increasing rapidly: there is now far more demand for neural network models than available experts who can design model architectures

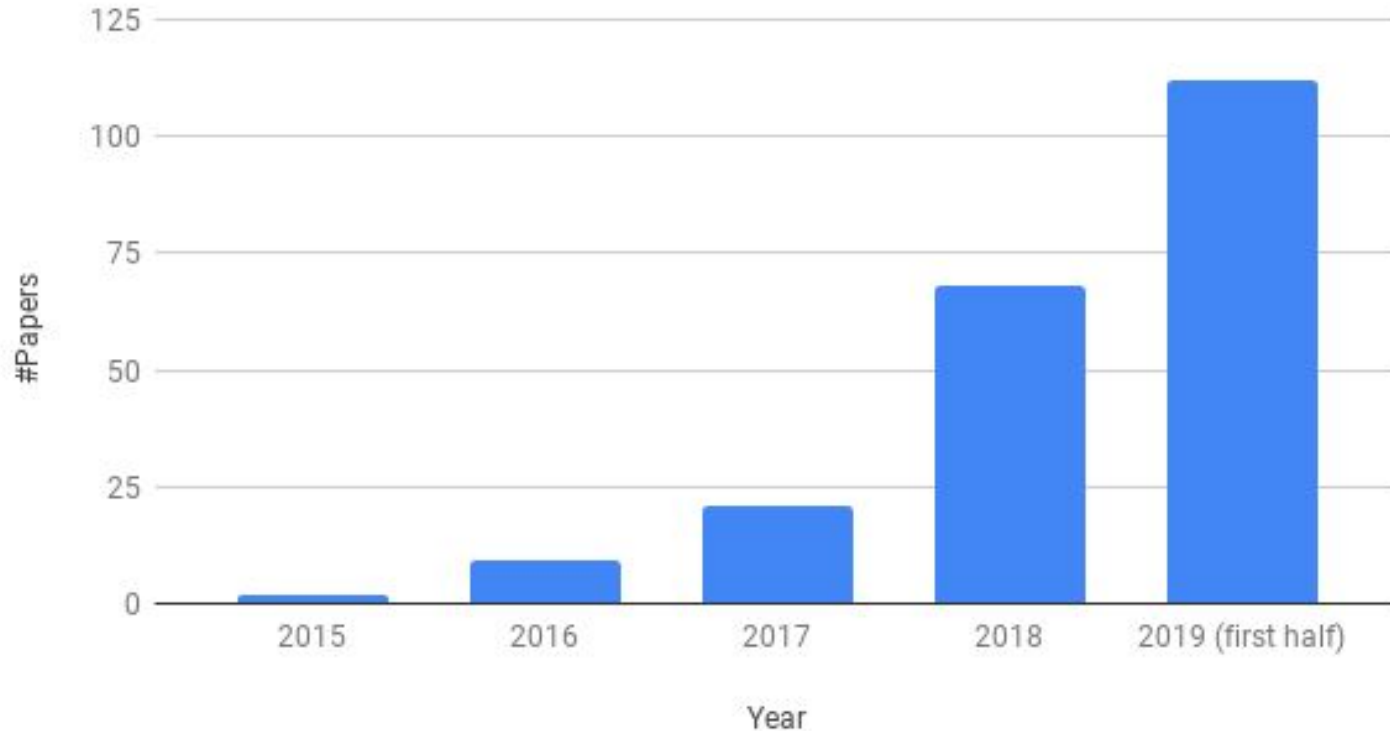


LSTM vs learned recurrent cell using reinforcement learning approach

Importance of Neural Network Architectures

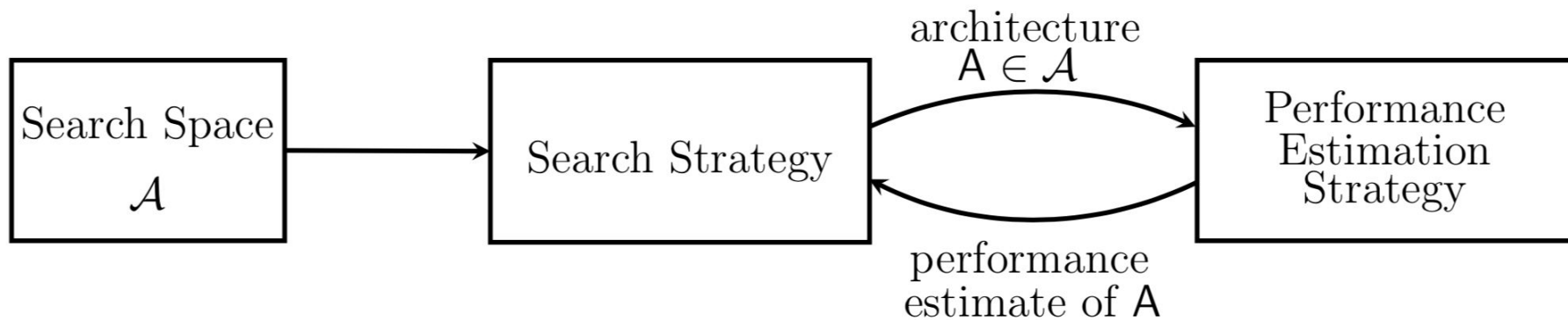
- The choice of neural network architecture can encode a prior belief or inductive bias for a given problem
- Recent advances in architecture design (e.g. ResNet, Transformer) attest to this fact

Number of papers on architecture search



NAS papers per year based on the literature list on automl.org.
The number for 2019 only considers the first half of 2019.
(Lindauer and Hutter, 2019)

Neural Architecture Search Workflow



Credit: Elskin et. al, 2019

NAS can be very expensive

Model	Hardware	Power (W)	Hours	kWh·PUE	CO ₂ e	Cloud compute cost
Transformer _{base}	P100x8	1415.78	12	27	26	\$41–\$140
Transformer _{big}	P100x8	1515.43	84	201	192	\$289–\$981
ELMo	P100x3	517.66	336	275	262	\$433–\$1472
BERT _{base}	V100x64	12,041.51	79	1507	1438	\$3751–\$12,571
BERT _{base}	TPUv2x16	—	96	—	—	\$2074–\$6912
NAS	P100x8	1515.43	274,120	656,347	626,155	\$942,973–\$3,201,722
NAS	TPUv2x1	—	32,623	—	—	\$44,055–\$146,848
GPT-2	TPUv3x32	—	168	—	—	\$12,902–\$43,008

Table 3: Estimated cost of training a model in terms of CO₂ emissions (lbs) and cloud compute cost (USD).⁷ Power and carbon footprint are omitted for TPUs due to lack of public information on power draw for this hardware.

NAS cost is from evolutionary architecture search on Transformer
(Strubell et al. 2019)

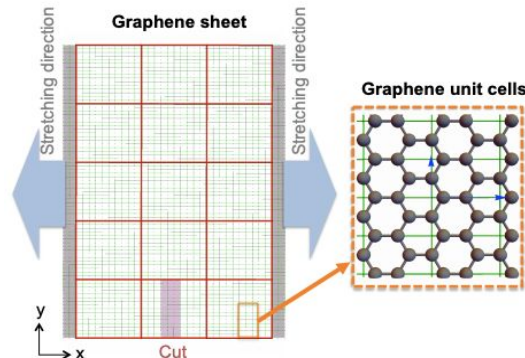
Scientific Datasets

Most NAS studies only use standard image datasets like CIFAR and ImageNet. However, deep learning also shows large potential for various physical sciences. Thus we want to see whether DARTS is useful for scientific datasets.

- ❑ **MNIST:** classifying images of handwritten digits (non-scientific baseline)
- ❑ **Graphene Kirigami:** cutting graphene to optimize stress/strain
- ❑ **Galaxy Zoo:** classifying galaxy morphology from telescope images
- ❑ **Chest X-Ray:** predicting diseases from chest x-rays



Not “scientific” but good
“hello world”



Introduction to DARTS

(Differentiable Neural ARchiTecture Search)

Issues with brute-force or traditional Neural Architecture Search (NAS) approaches

Proposed solutions in DARTS (Differentiable ARchiTecture Search, Liu 2019, ICLR)

Extremely large search space: arbitrary connections and operations between neural network nodes



Only search for the optimal "cell", i.e. a small unit of convolutional layers. Construct the complete model by stacking identical cells. (following NASNet, Zoph 2018, Google Brain)

Every "trial architecture" is **re-trained from scratch**, taking many GPU hours



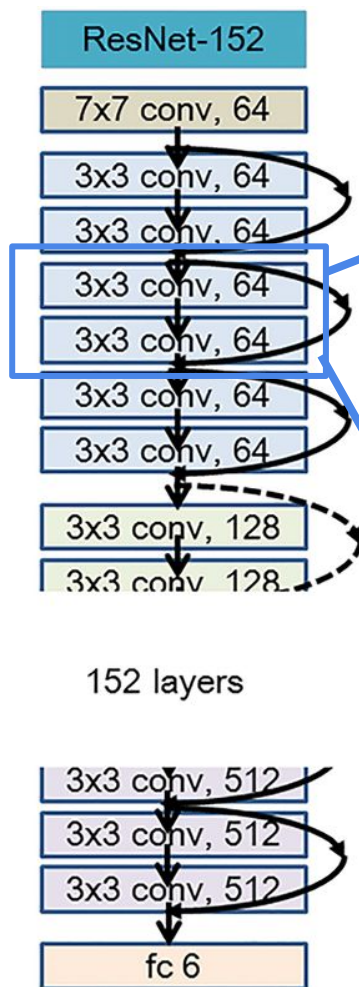
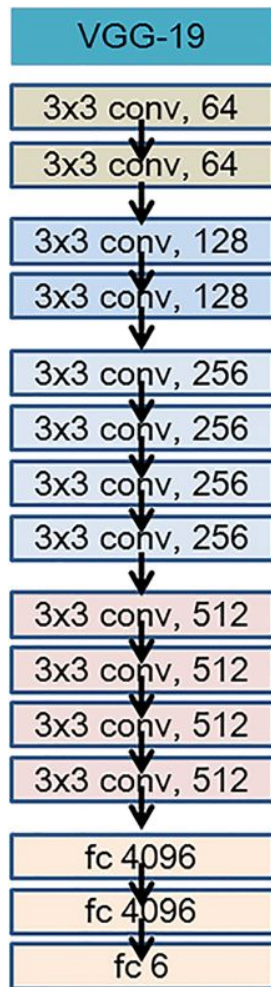
Share weights/parameters between child models, which can be "trained together" (similar to ENAS, Pham 2018, Google Brain)

The choice of operations (e.g. Conv, Pooling) is **discrete**, requiring expensive optimization

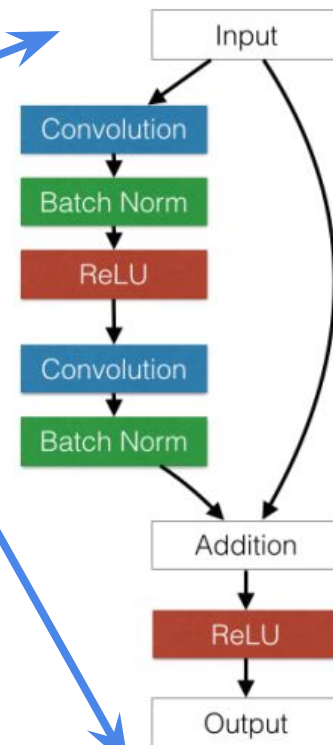


Continuous relaxation: parameterize the choice of operation by "architecture parameter" α , which can be **optimized by gradient descent**

Key observation:
popular CNN
architectures often
contain **repeating
blocks**, stacked
sequentially



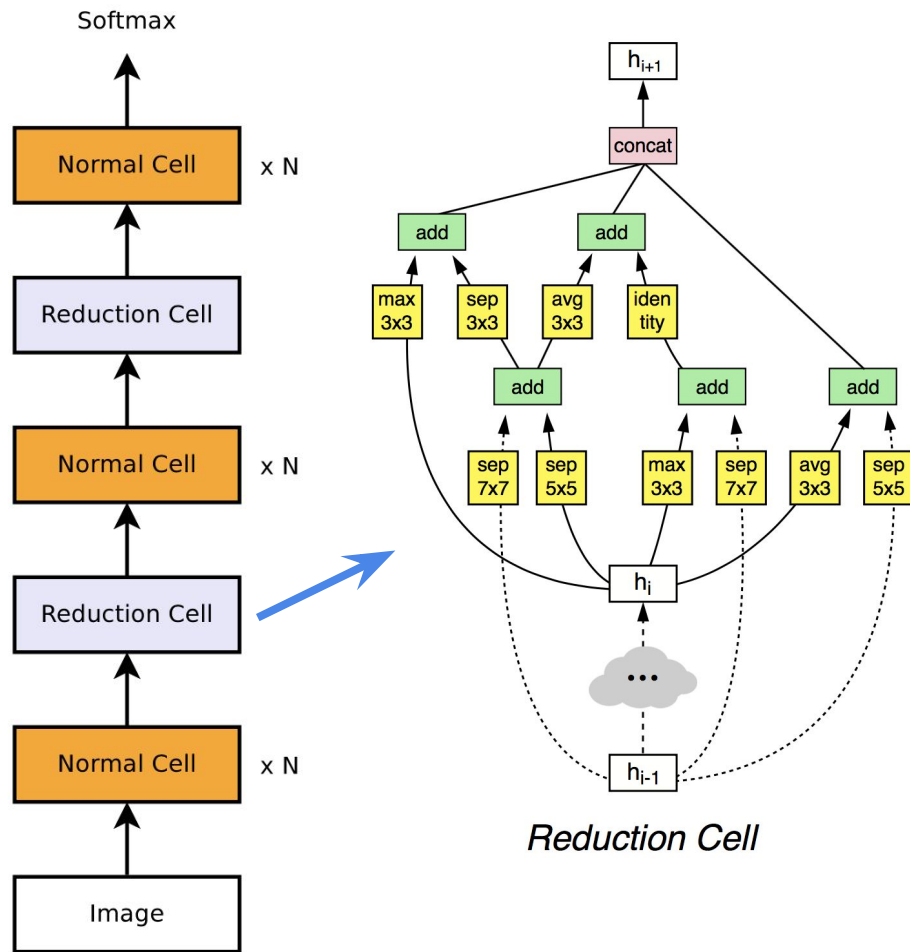
One residual block



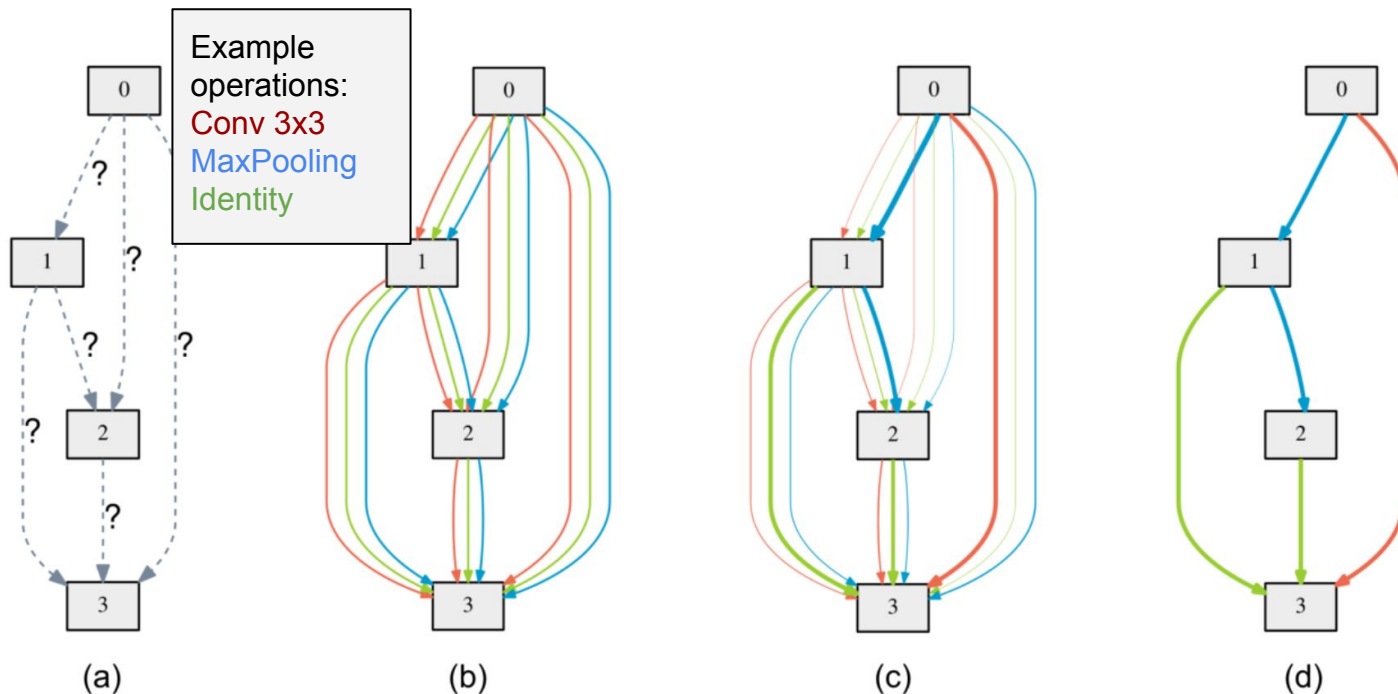
DARTS searches for the optimal "cell", not whole model

- Two types of cells:
 - Normal Cell:** output same dimension
 - Reduction Cell:** output half dimension
- Stack cells sequentially to form model
- Each cell type **share the same architecture but have independent weights**

(following NASNet, Zoph 2018, Google Brain)



Continuous relaxation of discrete operations enables gradient descent



Goal: Find the optimal cell, by placing proper operations (e.g. conv, pooling) at edges

Superpose: each edge is the sum over the outputs of multiple operations, weighted by continuous "architecture parameters" α

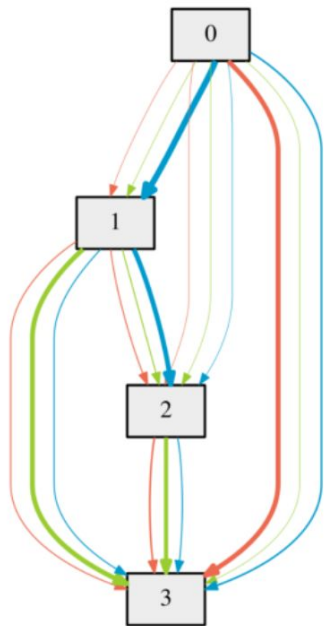
Search: Optimize the architecture weights α , using gradient descent on validation loss

Discretize: select the operation with the highest architecture weight, to be the final architecture

Gradient-based optimization for architecture parameter α

Example
operations:

Conv 3x3
MaxPooling
Identity



The actual operation at edge (i, j) is the average of all candidate operations $o(x)$, **weighted by α** :

$$\bar{o}^{(i,j)}(x) = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_o^{(i,j)})}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{o'}^{(i,j)})} o(x)$$

With a certain choice of **architecture weight α** , the corresponding architecture can be in principle trained to convergence, leading to the optimal **model weights $w^*(\alpha)$** and the **final validation loss $L_{\text{val}}(w^*(\alpha), \alpha)$** .

The gradient of L_{val} w.r.t to α gives the direction for gradient descent!

One-shot evaluation to avoid re-training

Computing the **true loss** \mathcal{L}_{val} by training w to the end is too expensive; thus DARTS just trains w for one step to get a proxy loss:

$$\begin{aligned} & \nabla_{\alpha} \mathcal{L}_{val}(w^*(\alpha), \alpha) \\ & \approx \nabla_{\alpha} \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha) \end{aligned}$$

where the optimal **model weights** $w^*(\alpha)$ is approximated by the one-step training

The training of α and w is performed in an alternate way:

while *not converged* **do**

- 1. Update architecture α by descending $\nabla_{\alpha} \mathcal{L}_{val}(w - \xi \nabla_w \mathcal{L}_{train}(w, \alpha), \alpha)$
($\xi = 0$ if using first-order approximation)
- 2. Update weights w by descending $\nabla_w \mathcal{L}_{train}(w, \alpha)$

However, DARTS doesn't consider traditional hyper-parameters

The parameters below are all fixed & chosen by human in the DARTS paper

- Number of cells to stack together (8 for search and 20 for final model)
- Batch size for SGD (64 / 96)
- Number of channels (16 / 36 for first layer)
- Training epochs (50 / 600)
- Learning rates (0.025 for model weights w and 0.0003 for architecture weight α)

We find that learning rate has crucial impact when applying DARTS to other datasets.

This problem has been pointed out by many ICLR reviewers:

<https://openreview.net/forum?id=S1eYHoC5FX¬elId=r1ekErZ53Q>

Results

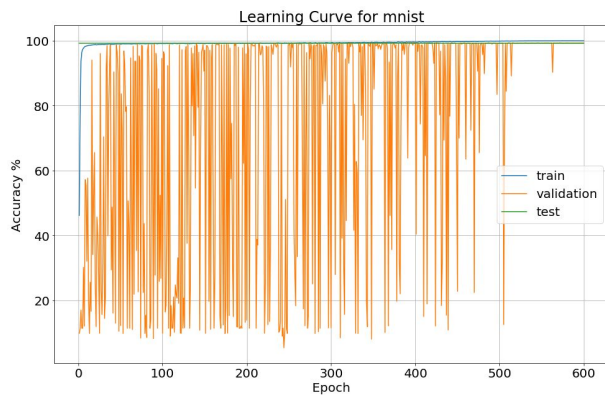
Overview

Model	MNIST	Graphene	Galaxy Zoo	Chest X-Ray
DARTS (Continuous)	99.07	0.89	0.094	0.157
DARTS (Discrete)	99.27	0.92	0.114	0.163
Random Search	99.31	0.90	0.098	0.169
ResNet	99.40	0.92	0.095	0.163
Metric	Accuracy	R^2	RMSE	BCE

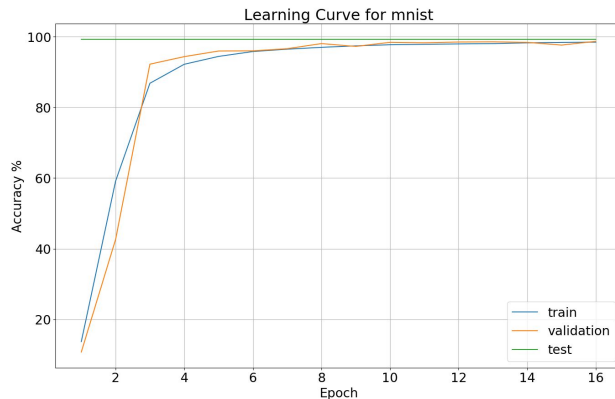
Results: MNIST

- Training MNIST with default learning rate fails
- Had to tune learning rate
- Key point: even on easy problem, DARTS is sensitive to hyperparameters

Default Hyper-parameters



Tuned Hyper-parameters



Model	MNIST
DARTS (Continuous)	99.07
DARTS (Discrete)	99.27
Random Search	99.31
ResNet	99.40
Metric	Accuracy

Large number of hyperparameters (~2x standard)

```
[--layers LAYERS] [--learning_rate LEARNING_RATE]
[--learning_rate_min LEARNING_RATE_MIN]
[--weight_decay WEIGHT_DECAY] [--L1_lambda L1_LAMBDA]
[--arch_learning_rate ARCH_LEARNING_RATE]
[--arch_weight_decay ARCH_WEIGHT_DECAY] [--cell_steps CELL_STEPS]
[--cell_multiplier CELL_MULTIPLIER] [--epochs EPOCHS]
[--batch_size BATCH_SIZE] [--optimizer OPTIMIZER]
[--momentum MOMENTUM] [--gz_dtree] [--fc1_size FC1_SIZE]
[--fc2_size FC2_SIZE] [--primitives PRIMITIVES]
[--train_portion TRAIN_PORTION] [--grad_clip GRAD_CLIP]
[--unrolled] [--cutout] [--cutout_length CUTOUT_LENGTH]
```

Train DARTS search

```
[--init_channels INIT_CHANNELS] [--layers LAYERS]
[--learning_rate LEARNING_RATE] [--drop_path_prob DROP_PATH_PROB]
[--weight_decay WEIGHT_DECAY] [--arch ARCH] [--epochs EPOCHS]
[--batch_size BATCH_SIZE] [--optimizer OPTIMIZER]
[--momentum MOMENTUM] [--fc1_size FC1_SIZE] [--fc2_size FC2_SIZE]
[--gz_dtree] [--primitives PRIMITIVES]
[--train_portion TRAIN_PORTION] [--grad_clip GRAD_CLIP]
[--auxiliary] [--auxiliary_weight AUXILIARY_WEIGHT] [--cutout]
[--cutout_length CUTOUT_LENGTH]
```

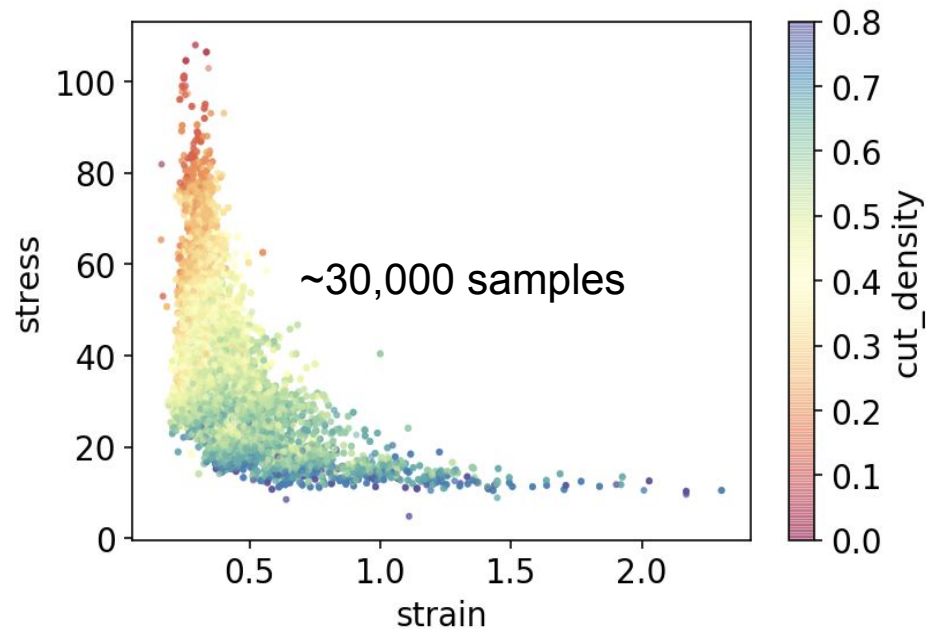
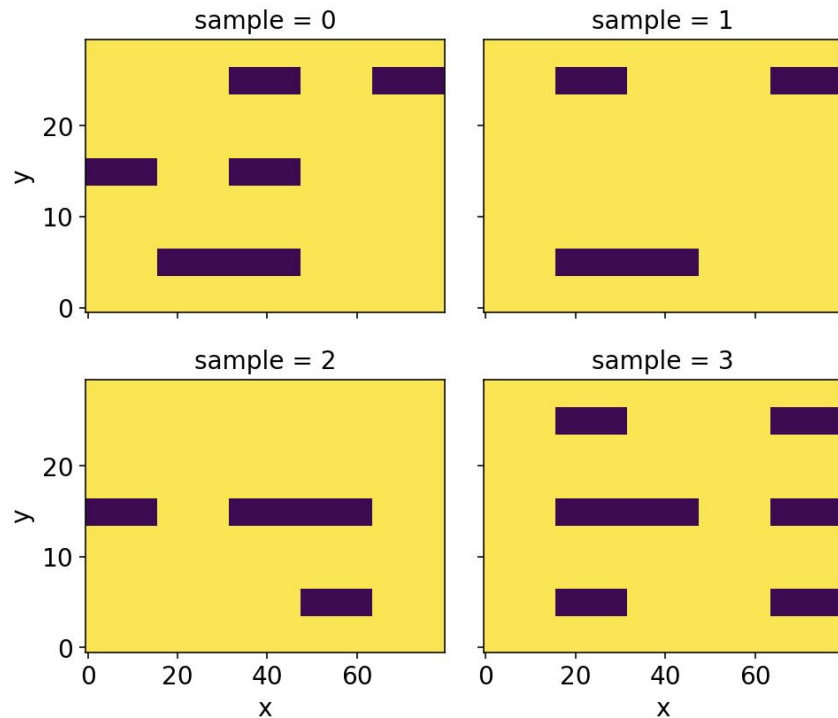
Train discretized

Problem: Graphene Kirigami

cut configuration

predict

stretch property



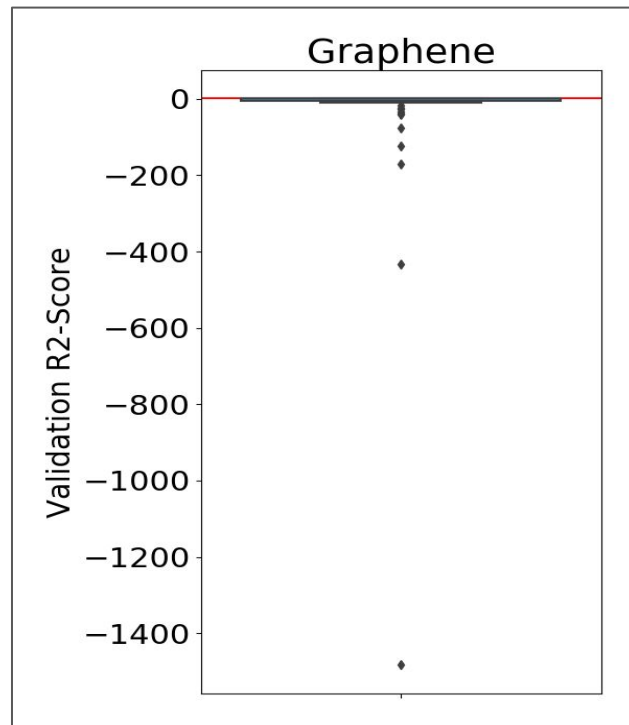
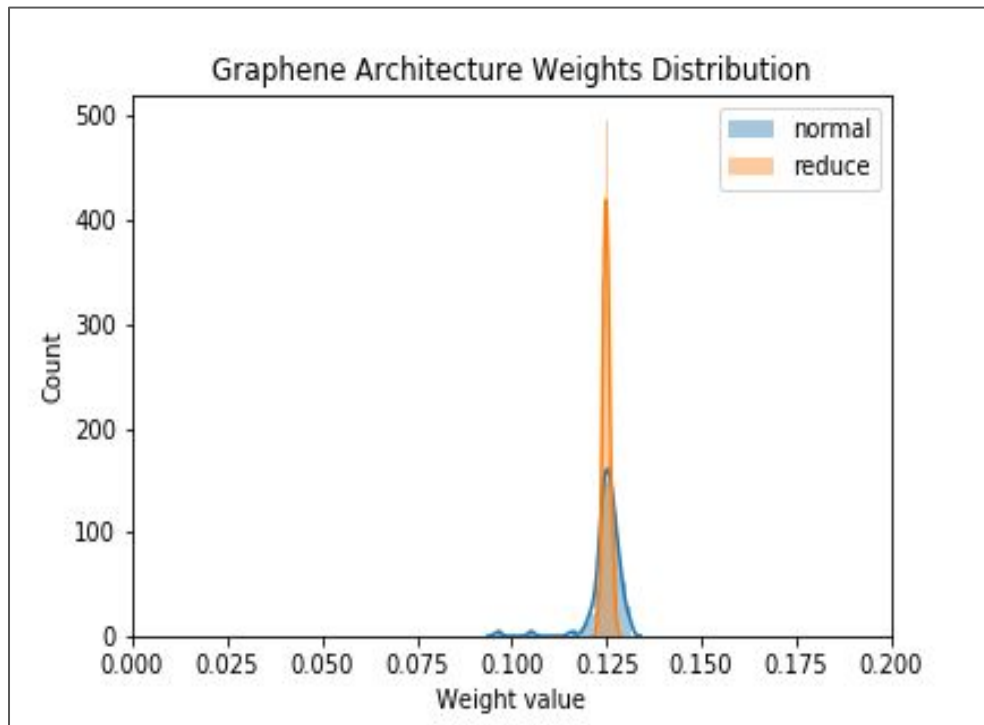
Results: Graphene Kirigami

Method / model	No. Parameters	Training time* (minutes)	Test R ²
DARTS	195,236	1042	0.92
ResNet-18	11,168,193	11	0.92
"Tiny ResNet" (10 layers, 8x less filters)	77,273	4	0.92

1. Time reported to train for 30 epochs on a single GPU.

Conclusion: Graphene Kirigami dataset is too simple for DARTS to be useful.

Graphene Architecture & Random Search



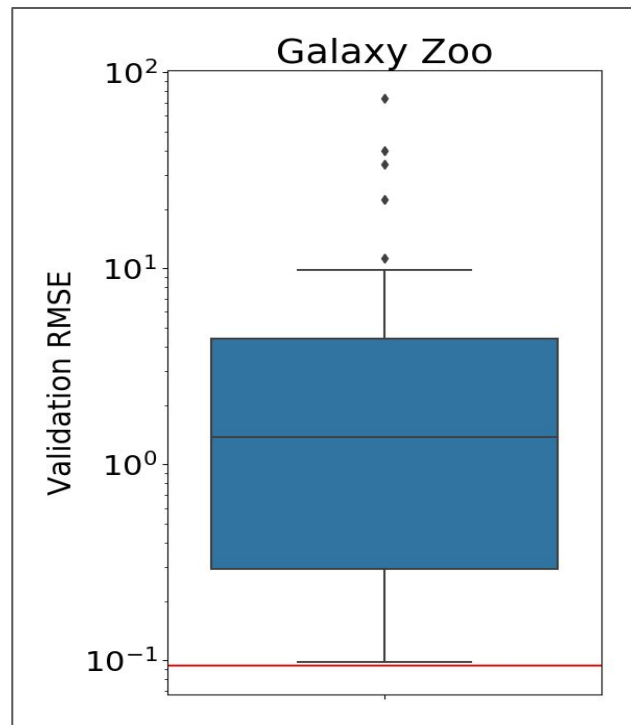
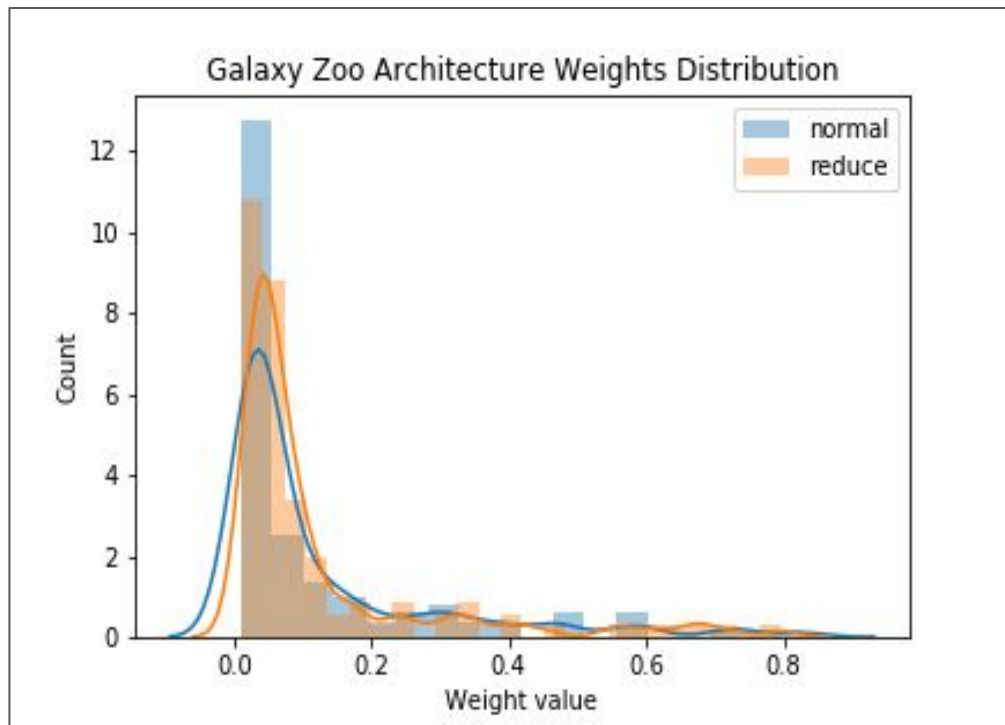
Simple problem → many good architectures → low sparsity in architecture weights

Results: Galaxy Zoo

- Continuous DARTS better than ResNet
- “Discretized” architecture is **worse**
- Shows heuristic discretization step can fail
- Kaggle winner score is ~ 0.075 (with extensive data augmentation and model ensembling)

Model	Galaxy Zoo
DARTS (Continuous)	0.094
DARTS (Discrete)	0.114
Random Search	0.098
ResNet	0.095
Metric	RMSE

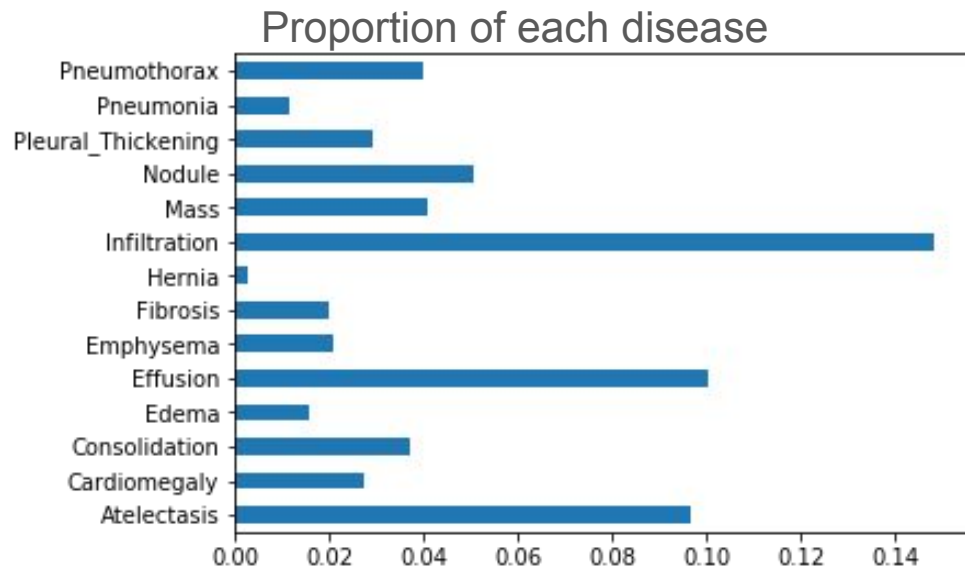
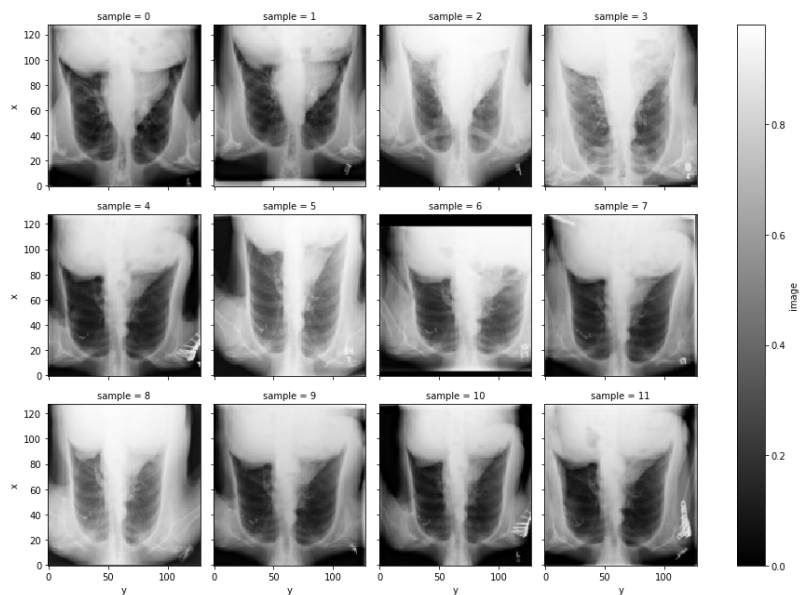
Galaxy Architecture & Random Search



Large variance in architectures \rightarrow sparse cell learned by DARTS

Problem: Chest X-Ray

- 39,589 chest X-rays
- 14 independent disease labels (confirmed diagnoses)
- Models assessed with mean binary cross entropy (BCE)

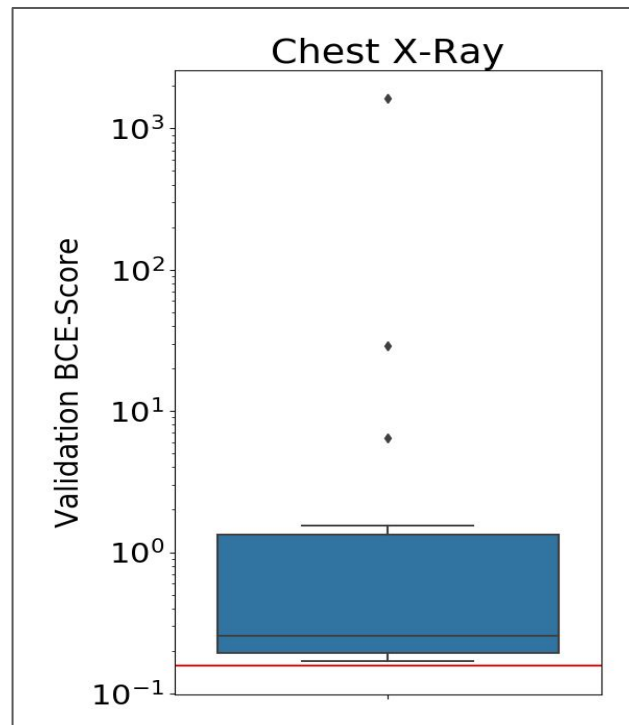
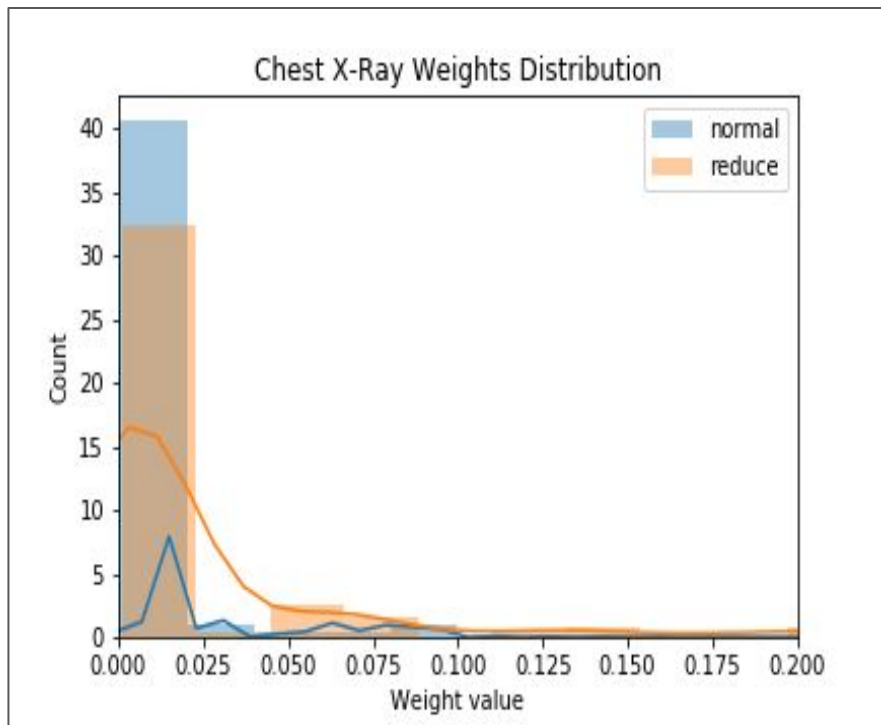


Results: Chest X-Ray

- DARTS performs well after some hyperparameter tuning
- The discretized network was worse (same as ResNet)
- Discretization step failed again

Model	Chest X-Ray
DARTS (Continuous)	0.157
DARTS (Discrete)	0.163
Random Search	0.169
ResNet	0.163
Metric	BCE

Chest X-Ray Architecture & Random Search



Again: large variance in architectures \rightarrow sparse cell learned by DARTS

Conclusions & Future Work

Conclusions

- DARTS a useful tool, but overkill on simple tasks
- ResNet and random search could be good enough
- DARTS introduces additional hyperparameters (~2x regular)
- DARTS "discretization step" is heuristic and can fail
- Computational cost: ~10x more expensive than single (discrete) model, due to overlapping 10 ops. Batch size reduced by ~10x due to memory footprint

Recommendation

- If small increase in performance important, DARTS worthwhile

Future Work

- Automatic tuning and/or better defaults for hyperparameters
- Fix discretization heuristic, a crucially overlooked step (or eliminate it)
 - Encourage sparsity with sparsemax in place of softmax, or L_p regularization on the architecture weights
 - Dynamically prune architecture to remove components during training, eliminate need to re-train

Thank you for a great semester!

DARTS: Useful Tool, but Often Overkill

- While DARTS is powerful, it requires setting many hyperparameters including learning rate, architecture learning rate, and number of layers
- Setting any of these parameters incorrectly can lead to catastrophic failure
- DARTS is at its best on problems like image classification where CNNs of repeated cells such as VGG or ResNet can do well
- Of course, on such a problem, ResNet will often do very well out of the gate!
- We do not see DARTS in its current form as a logical first choice when faced with a novel problem

Advice: First Try ResNet

- A recurring theme as we trained networks on both reference image classification and novel scientific datasets, was...
- Start with ResNet!
- ResNet consistently achieved good to excellent results, with remarkably short development time
- ResNet is the result of many years of effort by top researchers and their GPUs
- Because DARTS is searching for repeated cells and not a truly general architecture, ResNet will often fare well with much less effort

Advice: Next Try Random Search

- We observed that after equalizing for the amount of GPU time spent training, the best random search was often competitive with DARTS
- Random search only requires specifying the search space (and none of the other DARTS hyperparameters); it's also embarrassingly parallelizable
- This suggests that when faced with a novel problem, after specifying a search space, it might be advantageous to run RS before DARTS
- The results of the random search can be used to establish a baseline model, develop intuition about the problem, and perhaps refine the search space

DARTS Discretization May Reduce Performance

- As the ICLR reviewers noted, the discretization step at the end of DARTS is a heuristic with no theoretical guarantees of success
- Indeed our work demonstrated two cases in which the discrete architecture extracted from DARTS was inferior to the result obtained during the architecture search, even when it was trained for many more epochs
- One theoretical guarantee we *can* make is that if the architecture search has produced completely sparse weights where all entries are 1 or 0, the discretization step will not reduce performance
- This suggests an avenue of future work: modifying DARTS to encourage sparsity in the architecture weights
 - Two ideas include using sparsemax in place of softmax, and introducing an L_p regularization penalty on the architecture weights for $p < 1$ to promote sparsity

DARTS is a Powerful Tool Once You Learn to Use It

- Our previous conclusions highlight various limitations of DARTS
- We also advise first trying ResNet and Random Search
- We do not want to leave the impression that DARTS is a poor tool for NAS
- DARTS is a powerful tool—once you learn to use it
- Here are some tips for effective use of DARTS we would share
 - Choose a learning rate that isn't too big to start; 0.001 is a sensible default
 - Use the Adam optimizer instead of the SGD optimizer in the original DARTS paper
 - Start with an architecture learning rate that isn't too big either, e.g. 0.01
 - Monitor the evolution of the architecture weights; if they stay too tightly clustered around the starting weights, try again with a higher architecture LR
 - Don't use too many layers or too many initial channels to start; there's a high risk of overfitting or diverging training
 - We suggest 8 layers and 16 initial channels as defaults
 - Once you get decent results with these, you can add channels or depth later