

Efficient Neural Architecture Search

Martin Wistuba, Tejaswini Pedapati

Amazon Web Services, IBM Research

25 July 2022

Efficient NAS



**360
GPU DAYS**

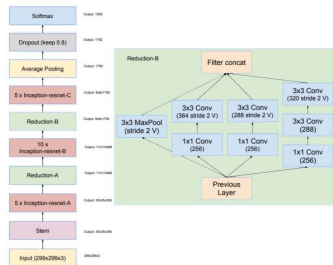
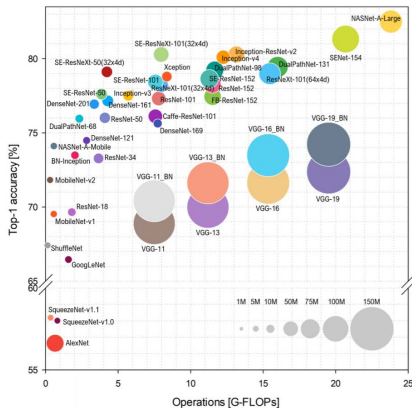


**0.1
GPU DAYS**

Outline

1. Introduction
2. One-Shot Architecture Search
3. Zero-Shot NAS
4. Transfer Learning for NAS
5. Conclusions

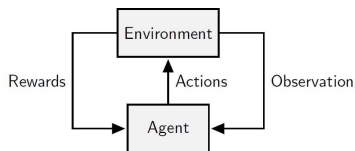
Architecture Design for Image Tasks



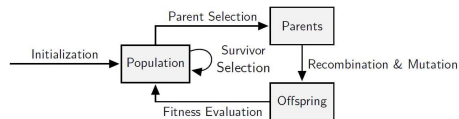
Innovations: Deeper networks, auxiliary classifiers, skip connections, bottlenecks, convolution stacking, global average pooling and many more

Images taken from Simone Bianco et al. "Benchmark Analysis of Representative Deep Neural Network Architectures". In: *IEEE Access* 6 (2018), pp. 64270–64277, Christian Szegedy et al. "Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning". In: *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4–9, 2017, San Francisco, California, USA.* 2017, pp. 4278–4284

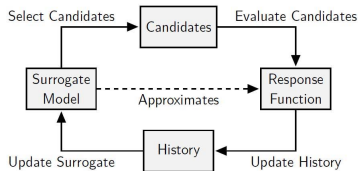
Neural Architecture Search



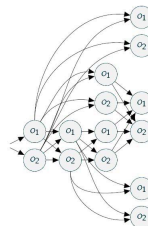
Reinforcement Learning



Evolutionary Algorithm



Surrogate Model-based Optimization



One-Shot Architecture Search

Tutorial Outline

Part 1

- ▶ Formal Definition of NAS and NASNet search space
- ▶ One-shot techniques in NAS
 - ▶ Overview
 - ▶ Shortcomings
 - ▶ Other flavours

Part 2

- ▶ Zero-shot NAS
- ▶ Effective NAS with transfer learning approaches based on
 - ▶ Transfer NAS optimizers
 - ▶ Few-Shot NAS optimizers
 - ▶ Learning Curve Ranking

Problem Definition

Machine Learning Problem

$$\Lambda(\alpha, d) = \arg \min_{m_{\alpha, \theta} \in M_{\alpha}} \mathcal{L}(m_{\alpha, \theta}, d_{\text{train}}) + \mathcal{R}(\theta) . \quad (1)$$

- ▶ m - machine learning model
- ▶ θ - model parameters
- ▶ α - neural architecture
- ▶ d - dataset

NAS Problem

$$\alpha^* = \arg \max_{\alpha \in A} \mathcal{O}(\Lambda(\alpha, d_{\text{train}}), d_{\text{valid}}) = \arg \max_{\alpha \in A} f(\alpha) . \quad (2)$$

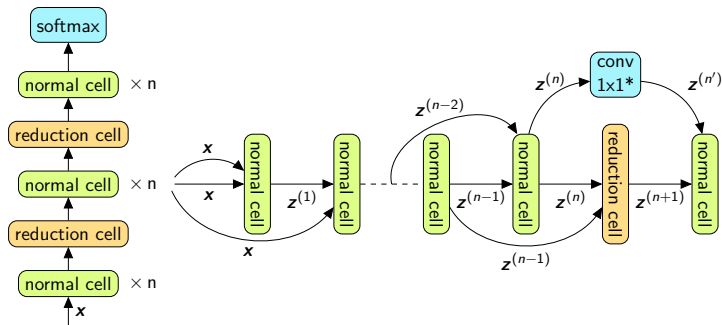
- ▶ f - response function
- ▶ A - search space

Search Space

- ▶ Neural architecture search space: subspace of all possible neural architectures.
- ▶ The limitation to a subspace allows for considering
 - ▶ human expert knowledge,
 - ▶ specific task (e.g. mobile architectures) and
- ▶ We distinguish two types of search spaces:
 - ▶ global search space
 - ▶ cell-based search space

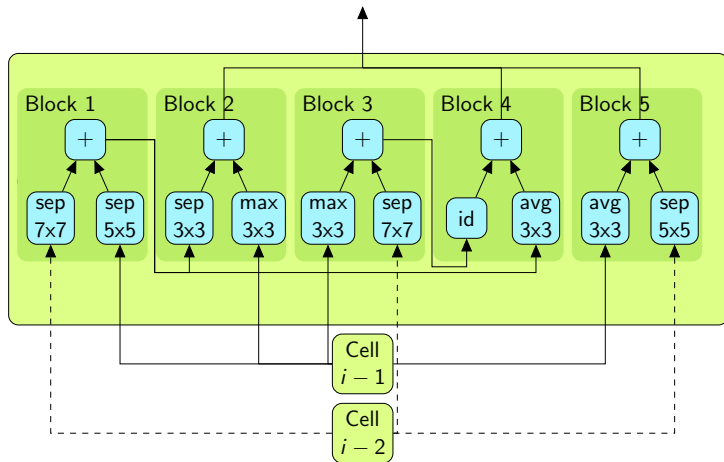
NASNet Search Space

Architectures from a cell-based search space are built by stacking few cells with the same topology.



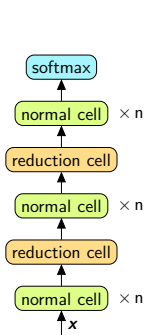
NASNet Search Space

Structure of a cell.

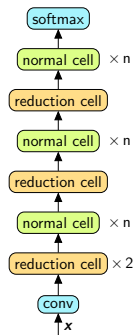


Transferring Architectures

Architectures from cell-based search spaces allow for easy transferability across different datasets.



(a) CIFAR-10



(b) ImageNet

Outline

1. Introduction

2. One-Shot Architecture Search

2.1 Overview

2.2 Shortcomings

- Memory Consumption

- DARTS Collapse

- DARTS Discretization

- Ranking Discrepancy

2.3 Other flavours of One-shot NAS

3. Zero-Shot NAS

4. Transfer Learning for NAS

Outline

1. Introduction

2. One-Shot Architecture Search

2.1 Overview

2.2 Shortcomings

- Memory Consumption

- DARTS Collapse

- DARTS Discretization

- Ranking Discrepancy

2.3 Other flavours of One-shot NAS

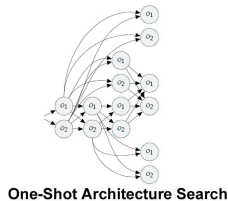
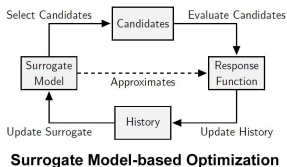
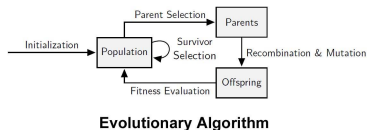
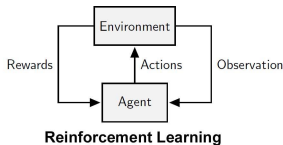
3. Zero-Shot NAS

4. Transfer Learning for NAS

NAS Optimizers

We distinguish several methods that maximize the response function:

- **Reinforcement learning:** learn to sample α that maximize f .
- **Evolutionary algorithms:** evolve α that maximize f .
- **Surrogate model-based optimization:** approximate f by \hat{f} and use it to maximize f .
- **One-shot architecture search:** learn one model and use it to max f .



One-Shot Architecture Search

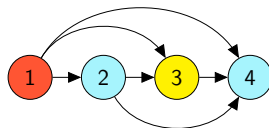
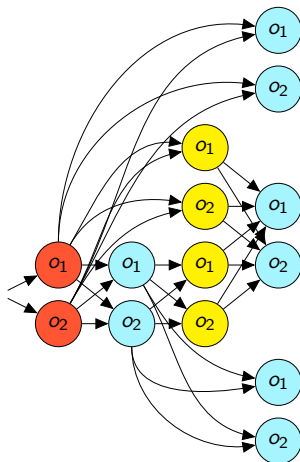
Until now,

- ▶ the candidate architecture is trained from scratch to obtain validation accuracy
- ▶ Previously trained candidate architectures' weights were not reused.

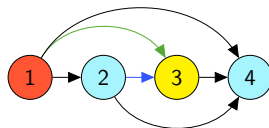
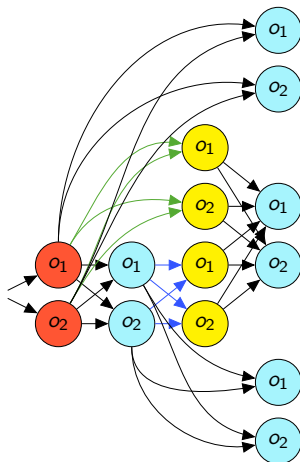
To overcome this, in one-shot architecture search the

- ▶ Entire search space is a directed acyclic graph - SuperNet
- ▶ Candidate architecture α is sampled from SuperNet
- ▶ The weights of all the operations are shared

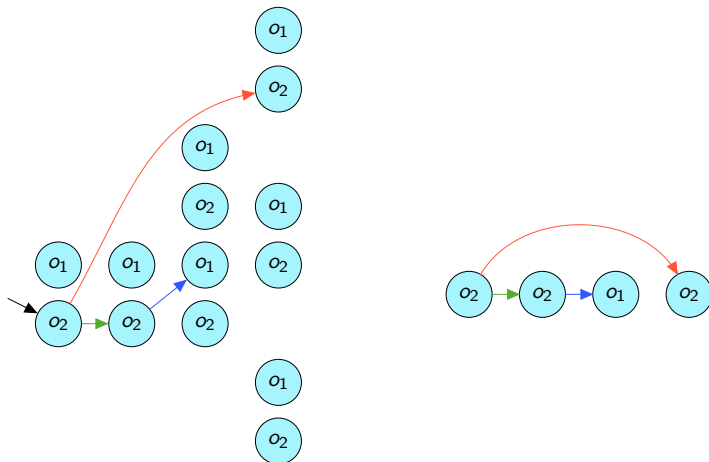
One-Shot Architecture Search



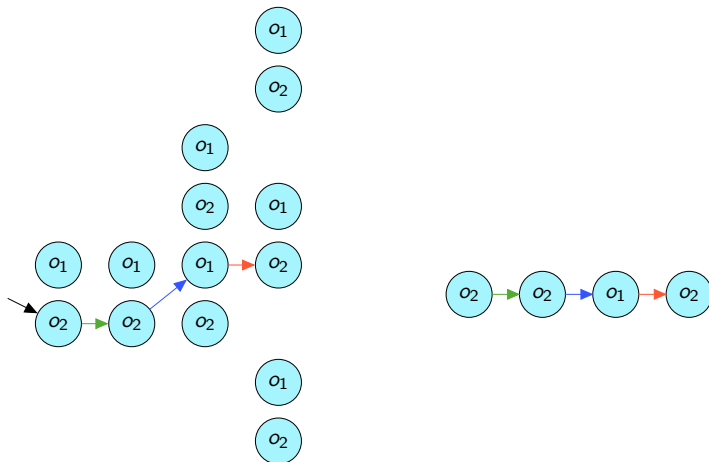
One-Shot Architecture Search



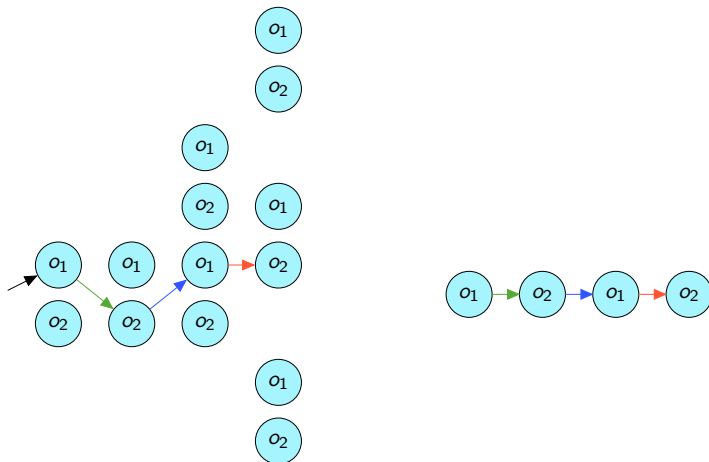
One-Shot Architecture Search



One-Shot Architecture Search



One-Shot Architecture Search



One-Shot Architecture Search

- ▶ Cross-entropy loss of α is computed on a minibatch of training data
- ▶ SuperNet parameters θ are updated using the gradients of the loss.
- ▶ Accelerated the search from 360 GPU days to 0.32 GPU days.
- ▶ Best architecture obtained by NAS is again trained from scratch

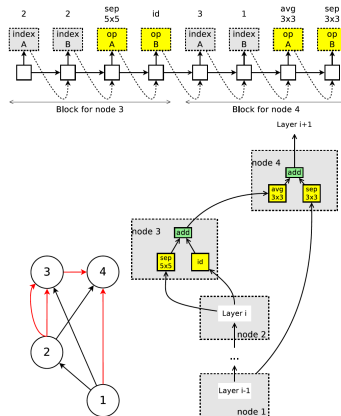
One-Shot Architecture Search

Sample Strategies

- ▶ Reinforcement learning (Pham et al.)
- ▶ Surrogate model-based optimization (Luo et al.)
- ▶ Learn a parameterized distribution (Casale et al.)
- ▶ Random sampling (Bender et al.)

Efficient Neural Architecture Search (ENAS)

Uses LSTM controller trained using RL to predict candidate network



Hieu Pham et al. "Efficient Neural Architecture Search via Parameter Sharing". In: *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*. 2018, pp. 4092–4101

Efficient Neural Architecture Search (ENAS)

Algorithm 1 ENAS

Input: Controller's policy parameters ω , SuperNet's parameters θ ,
for every iteration **do**
 Controller's policy samples candidate model α
 Compute cross-entropy loss $\nabla_{\theta} E_{\alpha}$ on m for a mini-batch of training data
 Fix ω and perform SGD on θ using $\nabla_{\theta} E_{\alpha}$
 Fix θ and update ω to maximize expected reward on validation data.

Differentiable Architecture Search (DARTS)

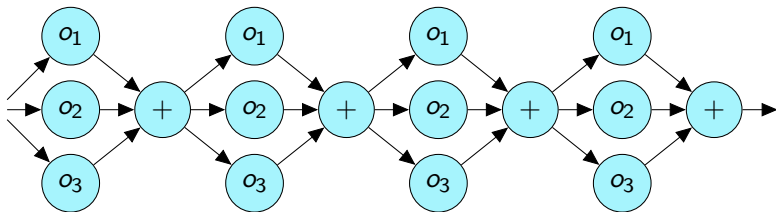
- ▶ In ENAS, choosing operations at every edge is a discrete decision
- ▶ DARTS Makes it continuous defining mixed operation

$$\bar{o}^{(i,j)}(x) = \sum_{o \in \mathcal{O}} \frac{\exp(\alpha_o^{(i,j)})}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{o'}^{(i,j)})} o(x) \quad (3)$$

- ▶ Architecture α is parameterized by β and network weights θ
- ▶ Magnitude of an operation: $\frac{\exp(\alpha_o^{(i,j)})}{\sum_{o' \in \mathcal{O}} \exp(\alpha_{o'}^{(i,j)})}$
- ▶ Derive discrete architecture by (1) $o^{(i,j)} = \operatorname{argmax}_{o \in \mathcal{O}} \alpha_o^{(i,j)}$ (2) choose top-k incoming edges

Hanxiao Liu, Karen Simonyan, and Yiming Yang. "DARTS: Differentiable Architecture Search". In: *Proceedings of the International Conference on Learning Representations, ICLR 2019, New Orleans, Louisiana, USA*. 2019

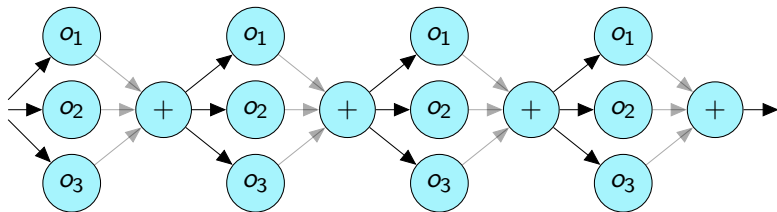
Differentiable Architecture Search



Relaxation of binary structural parameters α leads to differentiable loss:

$$\min_{\alpha(\beta) \in A} \mathcal{L} \left(\arg \min_{m_{\alpha(\beta), \theta} \in M_{\alpha(\beta)}} \mathcal{L}(m_{\alpha(\beta), \theta}, d_{\text{train}}) + \mathcal{R}(\theta), d_{\text{valid}} \right) \quad (4)$$

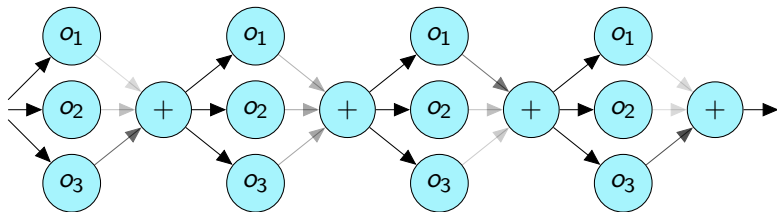
Differentiable Architecture Search



Relaxation of binary structural parameters α leads to differentiable loss:

$$\min_{\alpha(\beta) \in A} \mathcal{L} \left(\arg \min_{m_{\alpha(\beta)}, \theta \in M_{\alpha(\beta)}} \mathcal{L}(m_{\alpha(\beta)}, \theta, d_{\text{train}}) + \mathcal{R}(\theta), d_{\text{valid}} \right) \quad (4)$$

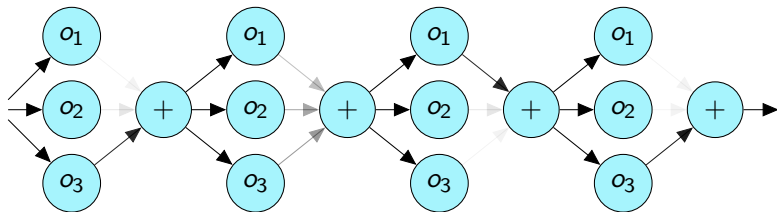
Differentiable Architecture Search



Relaxation of binary structural parameters α leads to differentiable loss:

$$\min_{\alpha(\beta) \in A} \mathcal{L} \left(\arg \min_{m_{\alpha(\beta)}, \theta \in M_{\alpha(\beta)}} \mathcal{L}(m_{\alpha(\beta)}, \theta, d_{\text{train}}) + \mathcal{R}(\theta), d_{\text{valid}} \right) \quad (4)$$

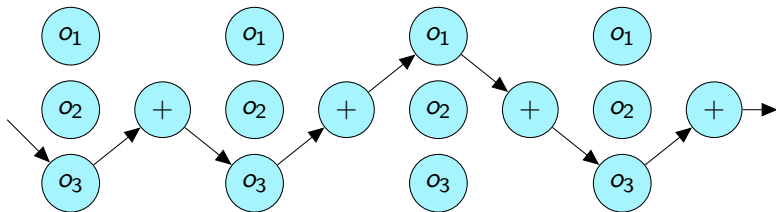
Differentiable Architecture Search



Relaxation of binary structural parameters α leads to differentiable loss:

$$\min_{\alpha(\beta) \in A} \mathcal{L} \left(\arg \min_{m_{\alpha(\beta)}, \theta \in M_{\alpha(\beta)}} \mathcal{L}(m_{\alpha(\beta)}, \theta, d_{\text{train}}) + \mathcal{R}(\theta), d_{\text{valid}} \right) \quad (4)$$

Differentiable Architecture Search



Relaxation of binary structural parameters α leads to differentiable loss:

$$\min_{\alpha(\beta) \in A} \mathcal{L} \left(\arg \min_{m_{\alpha(\beta), \theta} \in M_{\alpha(\beta)}} \mathcal{L}(m_{\alpha(\beta), \theta}, d_{\text{train}}) + \mathcal{R}(\theta), d_{\text{valid}} \right) \quad (4)$$

Bi-level optimization

$$\min_{\alpha} \mathcal{L}_{val}(\theta^*(\alpha), \alpha) \quad (5)$$

$$\text{s.t. } \theta^*(\alpha) = \operatorname{argmin}_{\theta} \mathcal{L}_{train}(\theta, \alpha) \quad (6)$$

Algorithm 2 DARTS – Differentiable Architecture Search

Input: A mixed operation $\bar{o}^{(i,j)}$

- 1: **while** not converged **do**
 - 2: Update architecture α by descending
 - 3: $\nabla_{\alpha} \mathcal{L}_{val}(\theta - \xi \nabla_{\theta} \mathcal{L}_{train}(\theta, \alpha), \alpha)$
 - 4: ($\xi = 0$ if using first-order approximation)
 - 5: Update weights θ by descending $\nabla_{\theta} \mathcal{L}_{train}(\theta, \alpha)$
 - 6: Derive the final architecture based on the learned α .
-

Outline

1. Introduction

2. One-Shot Architecture Search

2.1 Overview

2.2 Shortcomings

- Memory Consumption

- DARTS Collapse

- DARTS Discretization

- Ranking Discrepancy

2.3 Other flavours of One-shot NAS

3. Zero-Shot NAS

4. Transfer Learning for NAS

Pitfalls of DARTS

- ▶ All parameters need to be stored in memory.
- ▶ DARTS Collapse: Final architectures comprise of too many skip connections
- ▶ Discretization step results in architectures with higher validation loss

Outline

1. Introduction

2. One-Shot Architecture Search

2.1 Overview

2.2 Shortcomings

- Memory Consumption

- DARTS Collapse

- DARTS Discretization

- Ranking Discrepancy

2.3 Other flavours of One-shot NAS

3. Zero-Shot NAS

4. Transfer Learning for NAS

Memory consumption of DARTS

- ▶ In DARTS output of an edge is a weighted sum of all the operations
$$\sum_{i=1}^N \frac{\exp(\alpha_i)}{\sum_j \exp(\alpha_j)} o_i(x)$$
- ▶ It requires all possible combinations of the operations to be stored in memory
- ▶ The batch size used to train the SuperNet is small
- ▶ Searching on Imagenet takes several days.

ProxylessNAS

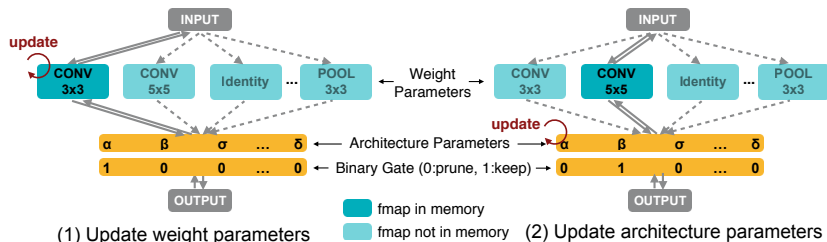
- ▶ Feature maps of N paths / operations are in memory
- ▶ In ProxylessNAS only one path / operation is active at a time.
Use binary gates for each edge:

$$g = \text{binarize}(p_1, \dots, p_N) = \begin{cases} [1, 0, \dots, 0] & \text{with probability } p_1, \\ \dots & \\ [0, 0, \dots, 1] & \text{with probability } p_N. \end{cases} \quad (7)$$

$$m_{\mathcal{O}}^{\text{Binary}}(x) = \sum_{i=1}^N g_i o_i(x) = \begin{cases} o_1(x) & \text{with probability } p_1 \\ \dots & \\ o_N(x) & \text{with probability } p_N. \end{cases} \quad (8)$$

Han Cai, Ligeng Zhu, and Song Han. “ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware”. In: *Proceedings of the International Conference on Learning Representations, ICLR 2019, New Orleans, Louisiana, USA*. 2019

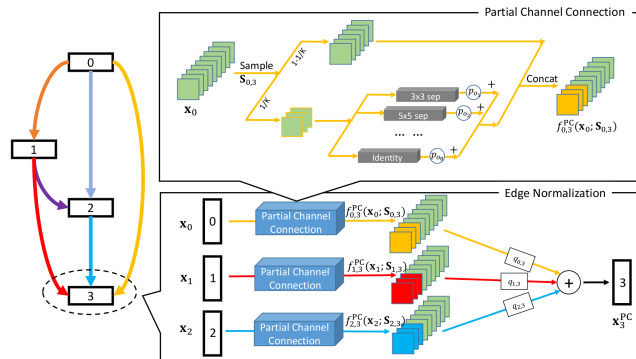
ProxylessNAS (cont.)



- Devise it as multiple binary selection tasks
- Requires only 2 paths in memory at any point.
- Able to search on ImageNet in 8.3 days

PC-DARTS

- Use channel mask $S_{i,j}$ to sample a $1/K$ channels each time
- K will determine accuracy vs search cost trade-off



Yuhui Xu et al. "PC-DARTS: Partial Channel Connections for Memory-Efficient Architecture Search". In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. 2020

PC-DARTS (cont.)

$$f_{i,j}^{\text{PC}}(\mathbf{x}_i; \mathbf{S}_{i,j}) = \sum_{o \in \mathcal{O}} \frac{\exp \{ \alpha_{i,j}^o \}}{\sum_{o' \in \mathcal{O}} \exp \{ \alpha_{i,j}^{o'} \}} \cdot o(\mathbf{S}_{i,j} * \mathbf{x}_i) + (1 - \mathbf{S}_{i,j}) * \mathbf{x}_i. \quad (9)$$

- ▶ To account for changing sampled channels, edge normalization is introduced
- ▶ All the edges contributing to the output of node j are assigned weights

$$\mathbf{x}_j^{\text{PC}} = \sum_{i < j} \frac{\exp \{ \gamma_{i,j} \}}{\sum_{i' < j} \exp \{ \gamma_{i',j} \}} \cdot f_{i,j}(\mathbf{x}_i). \quad (10)$$

- ▶ First train SuperNet for 15 epochs
- ▶ Increased batch size stabilizes the training
- ▶ The bias of choosing parameter-free operations is less pronounced

Outline

1. Introduction

2. One-Shot Architecture Search

2.1 Overview

2.2 Shortcomings

- Memory Consumption

- DARTS Collapse**

- DARTS Discretization

- Ranking Discrepancy

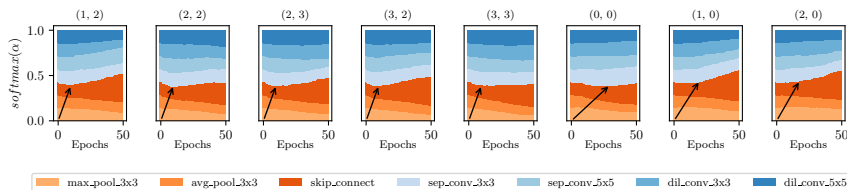
2.3 Other flavours of One-shot NAS

3. Zero-Shot NAS

4. Transfer Learning for NAS

DARTS Collapse

- ▶ Skip connections increase as search progresses
- ▶ Skip connections make it easier for the SuperNet to train although they do not boost the accuracy of the final discretized architecture

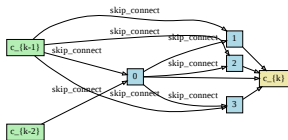


The softmax evolution where skip connections gradually become dominant. Image taken from Chu et al.

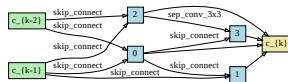
DARTS Collapse

- S1:** This search space uses a different set of only two operators per edge.
- S2:** Operated used are $\{3 \times 3 \text{ SepConv}, \text{SkipConnect}\}$.
- S3:** The set of candidate operations per edge is $\{3 \times 3 \text{ SepConv}, \text{SkipConnect}, \text{Zero}\}$.
- S4:** The set of candidate operations per edge is $\{3 \times 3 \text{ SepConv}, \text{Noise}\}$.

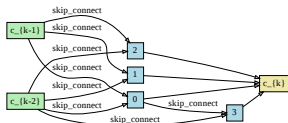
DARTS Collapse



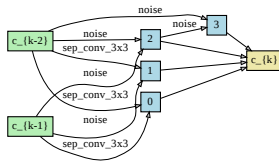
(a) Space 1



(b) Space 2



(c) Space 3



(d) Space 4

The normal cells standard DARTS finds on spaces S1-S4. Image taken from Zela et al.

DARTS Collapse

- ▶ Dropout after skip connection as suggested by P-DARTS
- ▶ Explicitly limit the number of skip connections: DARTS+
- ▶ Experiment done by FairDARTS

Methods	Cifar10-Acc
Random (M=2)	97.01 \pm 0.24
Random (M=2, MultAdds \geq 500M)	97.14 \pm 0.28
DARTS without skip-connection	96.88 \pm 0.18
DARTS (First Order) + Gaussian (cosine decay)	97.12 \pm 0.23
DARTS (First Order)	97.00 \pm 0.14

Operation choice no longer mutually exclusive

Apply a *sigmoid activation* (σ) for each $\alpha_{o_{i,j}}$, so that each operation can be switched on or off independently without being suppressed.

$$\bar{o}_{i,j}(x) = \sum_{o \in \mathcal{O}} \sigma(\alpha_{o_{i,j}}) o(x). \quad (11)$$

For the sigmoid of architectural weights to tend towards 0 or 1, additional loss is used

$$L_{0-1} = -\frac{1}{N} \sum_i^N (\sigma(\alpha_i) - 0.5)^2 \quad (12)$$

$$L_{total} = \mathcal{L}_{val}(w^*(\alpha), \alpha) + w_{0-1} L_{0-1}. \quad (13)$$

The final architecture is discretized by using a threshold ($\sigma_{threshold}$) instead of argmax

Xiangxiang Chu et al. "Fair DARTS: Eliminating Unfair Advantages in Differentiable Architecture Search". In: *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XV*. 2020, pp. 465–480

Outline

1. Introduction

2. One-Shot Architecture Search

2.1 Overview

2.2 Shortcomings

- Memory Consumption

- DARTS Collapse

- DARTS Discretization**

- Ranking Discrepancy

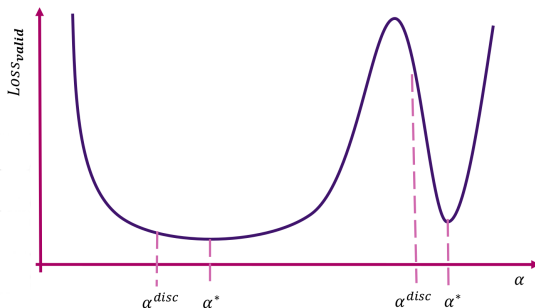
2.3 Other flavours of One-shot NAS

3. Zero-Shot NAS

4. Transfer Learning for NAS

DARTS Discretization

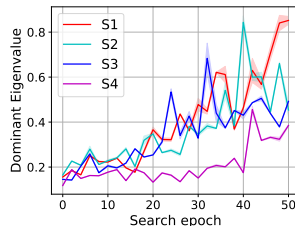
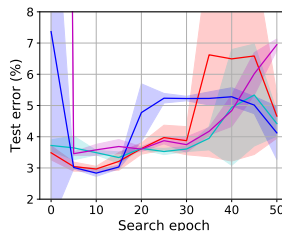
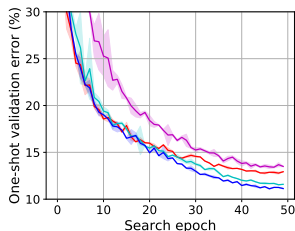
- ▶ DARTS found a sharp local minima
- ▶ Validation loss increased on discretization
- ▶ Need to find smoother local minima



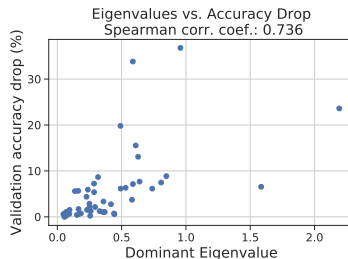
Arber Zela et al. "Understanding and Robustifying Differentiable Architecture Search". In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. 2020

DARTS Discretization

- RobustDARTS studied the relationship between the eigenvalues of the Hessian matrix of validation loss $\nabla_{\alpha(\beta)}^2 L_{valid}$ and the generalization error.



Accuracy Drop During Discretization Step



- ▶ Early stop if $\lambda_{max}^{-\alpha}(i-k) / \lambda_{max}^{-\alpha}(i) < 0.75$
- ▶ Increase l2 regularization of the network weights
- ▶ Apply cutout augmentation along with scheduled drop path

Anneal and Prune

- ▶ Avoid discretization by gradually removing operations from the mixed operation
- ▶ Anneal each operation to make its strength:

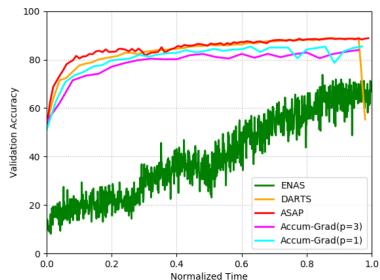
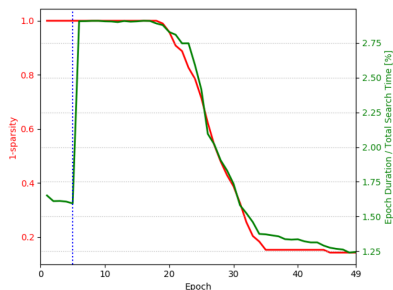
$$\Phi_o(\alpha^{(i,j)}; T) = \frac{\exp(\frac{\alpha_o^{(i,j)}}{T})}{\sum_{o' \in \mathcal{O}} \exp(\frac{\alpha_{o'}^{(i,j)}}{T})} \quad (14)$$

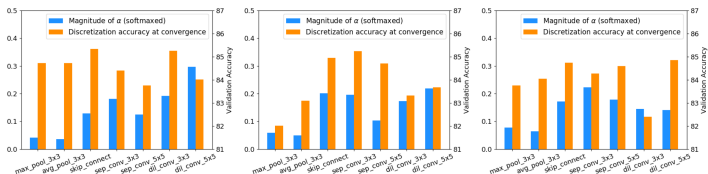
- ▶ Train SuperNet for some grace cycles τ
- ▶ For every iteration during bi-level optimization:
 - ▶ Prune operation if $\Phi_o(\alpha^{(i,j)}; T) < \text{threshold}$
 - ▶ Update threshold and T

Asaf Noy et al. "ASAP: Architecture Search, Anneal and Prune". In: *The 23rd International Conference on Artificial Intelligence and Statistics, AISTATS 2020, 26-28 August 2020, Online [Palermo, Sicily, Italy]*. 2020, pp. 493–503

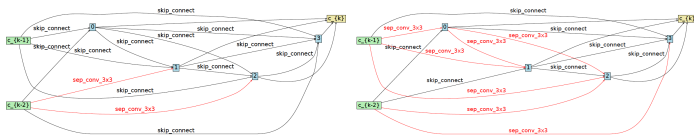
Anneal and Prune

- ▶ As SuperNet size reduces, search speed increases
- ▶ Searches on CIFAR-10 in 4.8 hours



Efficacy of α 

(a) Magnitude and strength based selection from trained supernet



(a) Magnitude

(b) Strength

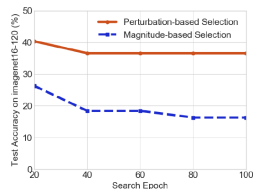
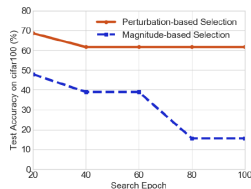
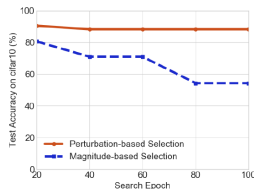
(b) Best operations chosen in Zela et al.'s S2 search space

Ruochen Wang et al. "Rethinking Architecture Selection in Differentiable NAS". In: *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021

Perturbation based architecture selection

Given a trained supernet,

- ▶ Randomly select an edge e
- ▶ For e , select the operation based on maximum drop in validation accuracy
- ▶ Fine tune the Supernet for few epochs



Uniform α

Space	DARTS	DARTS+PT	DARTS+PT (fix α)
DARTS Space	3.00	2.61	2.87
NAS-Bench-201	45.7	11.89	6.20

Dataset	Space	DARTS	DARTS+PT	DARTS+PT (fix α)
C10	S1	3.84	3.50	2.86
	S2	2.75	2.79	2.59
	S3	3.34	2.49	2.52
	S4	7.20	2.64	2.58
C100	S1	29.46	24.48	24.40
	S2	26.05	23.16	23.30
	S3	28.90	22.03	21.94
	S4	22.85	20.80	20.66
SVHN	S1	4.58	2.62	2.39
	S2	3.53	2.53	2.32
	S3	3.41	2.42	2.32
	S4	3.05	2.42	2.39

Outline

1. Introduction

2. One-Shot Architecture Search

2.1 Overview

2.2 Shortcomings

- Memory Consumption

- DARTS Collapse

- DARTS Discretization

- Ranking Discrepancy**

2.3 Other flavours of One-shot NAS

3. Zero-Shot NAS

4. Transfer Learning for NAS

Ranking Discrepancy of Weight-Sharing

Experiments are performed on reduced search space of NASBench-101 with 3 operations and 7 nodes.

Accuracy of NAS algorithms on 10 different searches.

NAS algo	Mean Acc.	Best Acc.	Best Rank	p(> random)
DARTS	92.21 \pm 0.61	93.02	57079	0.24
NAO	92.59 \pm 0.59	93.33	19552	0.62
ENAS	91.83 \pm 0.42	92.54	96939	0.07
NAO w/o WS	93.08 \pm 0.71	94.22	3543	0.92
ENAS w/o WS	93.54 \pm 0.45	94.22	4610	0.90
Best Arch	90.93 \pm 5.84	95.06		

Kaicheng Yu et al. "Evaluating The Search Phase of Neural Architecture Search".
In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. 2020

Ranking Discrepancy of Weight-Sharing

- Correlation between the architecture rankings found with and without weight-sharing for 200 architectures.

#Nodes	Kendall's τ
4	0.441
5	0.314
6	0.214
7	0.195

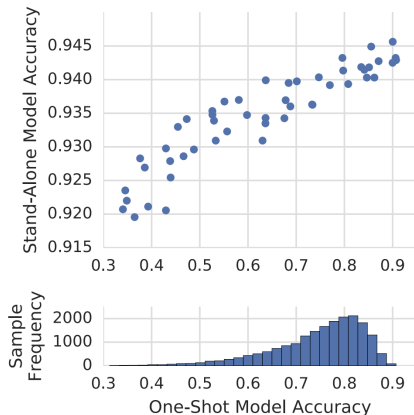
Effective Training of One-Shot Architectures

- ▶ Operations in one-shot model are subjected to co-adaptation
- ▶ Removing operations deteriorates performance
- ▶ Add dropout for every operation
- ▶ Use a variant of batch-normalization
- ▶ Apply L2 normalization only for the selected paths

Gabriel Bender et al. "Understanding and Simplifying One-Shot Architecture Search". In: *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*. 2018, pp. 549–558

Effective Training of One-Shot Architectures

- ▶ Sample 2000 architectures
- ▶ Train from scratch for 28 epochs



Outline

1. Introduction

2. One-Shot Architecture Search

2.1 Overview

2.2 Shortcomings

- Memory Consumption

- DARTS Collapse

- DARTS Discretization

- Ranking Discrepancy

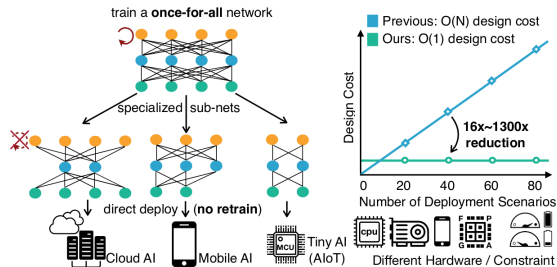
2.3 Other flavours of One-shot NAS

3. Zero-Shot NAS

4. Transfer Learning for NAS

Once-for-All

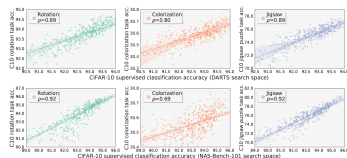
- ▶ A single network is trained to support versatile architectural configurations including depth, width, kernel size, and resolution.
- ▶ Training is difficult since weights will interfere with each other.
- ▶ Progressive shrinking: train the largest network and then fine-tune the network to support smaller sub-networks.



Han Cai et al. "Once-for-All: Train One Network and Specialize it for Efficient Deployment". In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. 2020

Unsupervised NAS

- During the search phase, use pretext tasks: rotation, colorization and solving jigsaw puzzle
- Train the best architecture on labeled data



method	search dataset & task	top-1 acc.	top-5 acc.	FLOPs (M)	params (M)
NAS-DARTS	CIFAR-10 Sup. Cls	73.3	91.3	574	4.7
NAS-P-DARTS	CIFAR-10 Sup. Cls	75.6	92.6	557	4.9
NAS-PC-DARTS	CIFAR-10 Sup. Cls	74.9	92.2	586	5.3
NAS-PC-DARTS	IN1K Sup. Cls	75.8	92.7	597	5.3
NAS-DARTS	CIFAR-10 Sup. Cls	74.9 \pm 0.08	92.2 \pm 0.06	538	4.7
NAS-DARTS	IN1K Sup. Cls	76.3 \pm 0.06	92.9 \pm 0.04	590	5.3
UnNAS-DARTS	IN1K Rot	75.8 \pm 0.18	92.7 \pm 0.16	558	5.1
UnNAS-DARTS	IN1K Color	75.7 \pm 0.12	92.6 \pm 0.16	547	4.9
UnNAS-DARTS	IN1K Jigsaw	75.9 \pm 0.15	92.8 \pm 0.10	567	5.2
NAS-DARTS	IN22K Sup. Cls	75.9 \pm 0.09	92.7 \pm 0.08	585	5.2
UnNAS-DARTS	IN22K Rot	75.7 \pm 0.23	92.7 \pm 0.11	549	5.0
UnNAS-DARTS	IN22K Color	75.9 \pm 0.21	92.8 \pm 0.10	547	5.0
UnNAS-DARTS	IN22K Jigsaw	75.9 \pm 0.31	92.8 \pm 0.13	559	5.1
NAS-DARTS	Cityscapes Seg	75.8 \pm 0.13	92.6 \pm 0.17	566	5.1
UnNAS-DARTS	Cityscapes Rot	75.9 \pm 0.19	92.7 \pm 0.07	554	5.1
UnNAS-DARTS	Cityscapes Color	75.2 \pm 0.15	92.4 \pm 0.05	594	5.1
UnNAS-DARTS	Cityscapes Jigsaw	75.5 \pm 0.06	92.6 \pm 0.08	566	5.0

Chenxi Liu et al. “Are Labels Necessary for Neural Architecture Search?” In: *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part IV*. ed. by Andrea Vedaldi et al. Vol. 12349. Lecture Notes in Computer Science. Springer, 2020, pp. 798–813

Conclusion

- ▶ NASNet Search Space
- ▶ One-shot model flavour of optimizers
 - ▶ ENAS and DARTS
 - ▶ Drawbacks of DARTS
 - ▶ Problem ranking weight-shared models
 - ▶ Once for all network and UNNAS

Outline

1. Introduction
2. One-Shot Architecture Search
3. Zero-Shot NAS
4. Transfer Learning for NAS
5. Conclusions

Definition and Motivation for Zero-Shot NAS

Definition

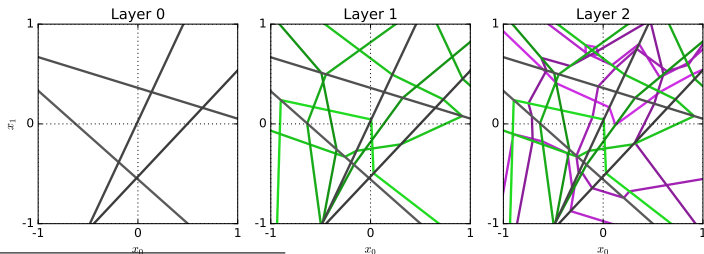
- ▶ No training.
- ▶ Depends on intrinsic properties of the architecture or data.
- ▶ Oftentimes a training-free score is computed which can be maximized instead of the validation score.

Motivation

- ▶ In some applications, even efficient NAS is too expensive, e.g. hardware-aware NAS.

Neural Network Expressivity

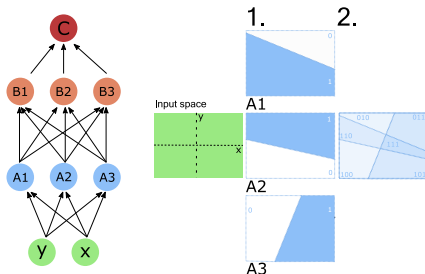
- ▶ Neural networks with ReLU activations represent piecewise linear functions.
- ▶ The input space is divided into convex polytopes in which the network behaves linearly.
- ▶ The number of those linear regions is a measure of the expressivity of the network.



Maithra Raghu et al. "On the Expressive Power of Deep Neural Networks". In: *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*. 2017, pp. 2847–2854

Neural Network Expressivity

- ▶ Each linear region has an activation pattern $a \in \{0, 1\}^{N_a}$ where N_a is the number of activations.
- ▶ a_i is 1 for a region if the neuron's activation is larger than 0.



Joe Mellor et al. "Neural Architecture Search without Training". In: *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*. 2021, pp. 7588–7598

NAS without Training

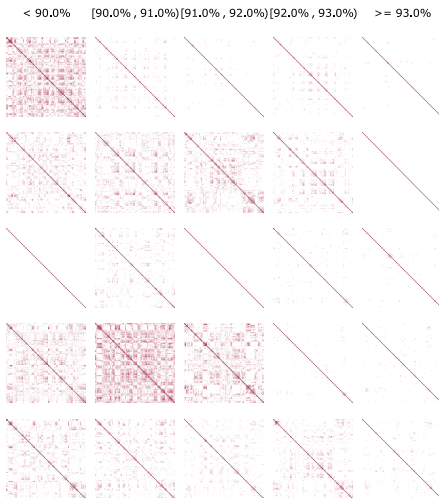
Assumption: An architecture is better if it maps the training data to different linear regions.

NASWOT Score for an architecture:

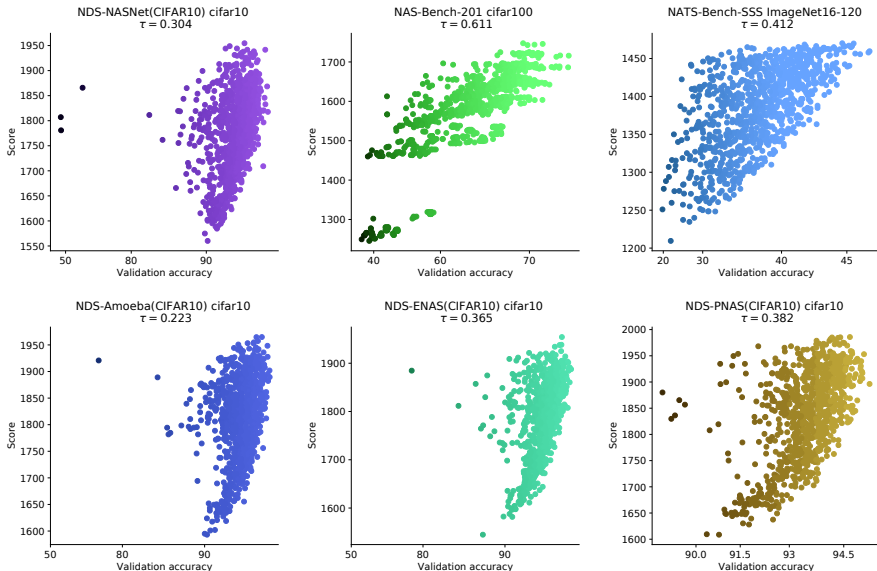
- ▶ Estimate the activation patterns a_1, \dots, a_n for n training points.
- ▶ Construct the matrix K with $K_{i,j} = N_a - \|a_i - a_j\|_1$.
- ▶ Estimate the score: $\text{NASWOT} = |K|$

Joe Mellor et al. “Neural Architecture Search without Training”. In: *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*. 2021, pp. 7588–7598

NASWOT Matrix



NASWOT Score Correlation



Zen Score

- ▶ Counting the number of linear regions is difficult and it does not account for the complexity of the regions themselves.
- ▶ The Zen-score measures the Gaussian complexity of the network instead:

$$\text{ZEN} = \mathbb{E}_{\mathbf{x}, \theta} \|\nabla_{\mathbf{x}} f(\mathbf{x}|\theta)\|_F \quad (15)$$

- ▶ Simplification: In practice overflow and batch normalization need a special treatment.
- ▶ Limitation: CNNs can only convolutional layers during scoring.

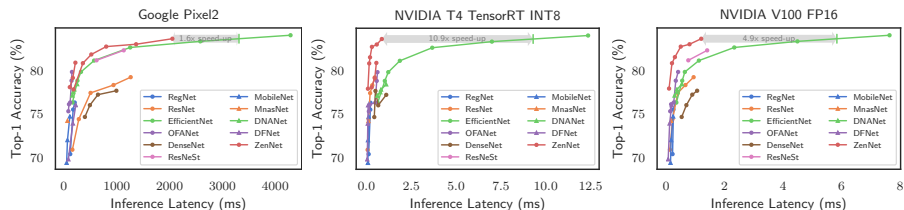
Ming Lin et al. “Zen-NAS: A Zero-Shot NAS for High-Performance Image Recognition”. In: *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. 2021, pp. 337–346

Zero-Shot Score Comparison

Proxy	CIFAR-10	CIFAR-100
Zen-Score	96.2%	80.1%
FLOPs	93.1%	64.7%
grad	92.8%	65.4%
synflow	95.1%	75.9%
TE-Score	96.1%	77.2%
NASWOT	96.0%	77.5%
Random	93.5±0.7%	71.1±3.1%

Proxy	Model	N	Time	Speed-Up
TE-Score	ResNet-18	16	0.34	1/28x
	ResNet-50	16	0.77	1/20x
NASWOT	ResNet-18	16	0.040	1/3.3x
	ResNet-50	16	0.059	1/1.6x
Zen-Score	ResNet-18	16	0.012	1.0
	ResNet-50	16	0.037	1.0

Zero-Shot Hardware-Aware NAS

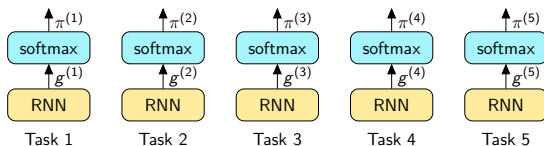


Outline

1. Introduction
2. One-Shot Architecture Search
3. Zero-Shot NAS
4. Transfer Learning for NAS
 - 4.1 Transfer NAS
 - 4.2 Few-Shot NAS
 - 4.3 Learning Curve Ranking
5. Conclusions

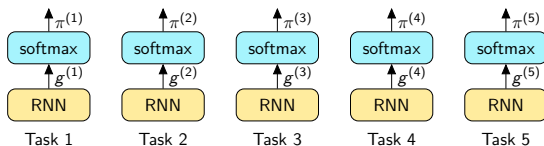
Motivation for Transfer Learning

- ▶ Standard NAS methods solve every problem independently.
- ▶ No knowledge is shared between different optimization problem.
- ▶ Every search starts from scratch again.



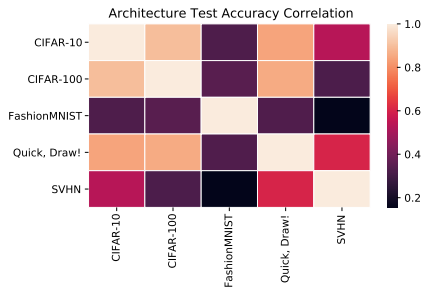
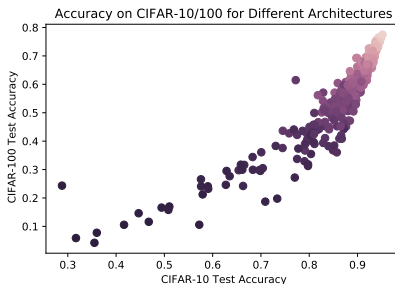
Motivation for Transfer Learning

- ▶ Standard NAS methods solve every problem independently.
- ▶ No knowledge is shared between different optimization problem.
- ▶ Every search starts from scratch again.



- ▶ **Can you reuse the knowledge of source tasks 1 to n for a new target task $n + 1$?**

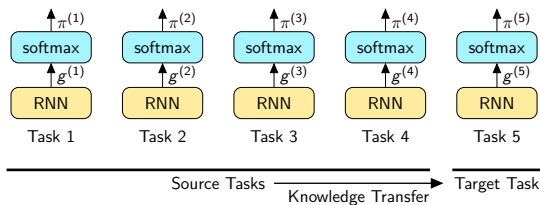
Metadata for Architecture Selection



- ▶ Strong architecture test accuracy correlation across tasks
- ▶ A good architecture on one task is very likely a good candidate for another

Motivation for Transfer Learning

How can we use metadata to improve Neural Architecture Search?



Agenda

- ▶ We cover various ways of using the *metadata* (knowledge of source tasks) to improve NAS on the target task.

Agenda

- ▶ We cover various ways of using the *metadata* (knowledge of source tasks) to improve NAS on the target task.
- ▶ Transfer Neural Architecture Search
 - ▶ Methods that incorporate transfer learning methods directly into NAS methods.

Agenda

- ▶ We cover various ways of using the *metadata* (knowledge of source tasks) to improve NAS on the target task.
- ▶ Transfer Neural Architecture Search
 - ▶ Methods that incorporate transfer learning methods directly into NAS methods.
- ▶ Few-Shot Learning
 - ▶ Methods that combine NAS with meta-learning.

Agenda

- ▶ We cover various ways of using the *metadata* (knowledge of source tasks) to improve NAS on the target task.
- ▶ Transfer Neural Architecture Search
 - ▶ Methods that incorporate transfer learning methods directly into NAS methods.
- ▶ Few-Shot Learning
 - ▶ Methods that combine NAS with meta-learning.
- ▶ Learning Curve Prediction
 - ▶ Methods that accelerate NAS methods by using early stopping methods that use transfer learning.

Outline

1. Introduction
2. One-Shot Architecture Search
3. Zero-Shot NAS
4. Transfer Learning for NAS
 - 4.1 Transfer NAS
 - 4.2 Few-Shot NAS
 - 4.3 Learning Curve Ranking
5. Conclusions

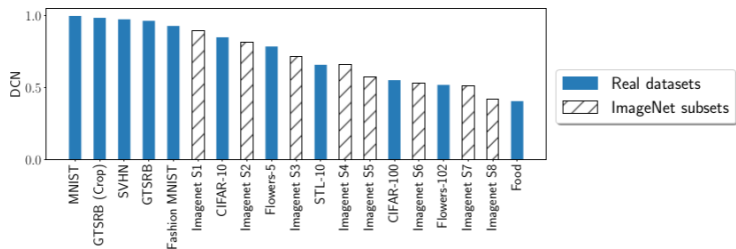
Trainless Accuracy Predictor Architecture Search

- ▶ TAPAS is a zero-shot transfer Neural Architecture Search algorithm
- ▶ The best architecture is searched using an evolutionary algorithm
- ▶ Instead of training and evaluating each architecture, a surrogate model is used
- ▶ This surrogate model is trained on metadata from similar datasets

Roxana Istrate et al. "TAPAS: Train-less Accuracy Predictor for Architecture Search". In: *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence, (AAAI-19), Honolulu, Hawaii, USA*. 2019

Dataset Similarity

- ▶ A dataset is defined by its difficulty (DCN)
- ▶ The DCN is defined by the validation accuracy obtained by a fixed architecture (landmarker)
- ▶ Assumption: datasets are similar iff their DCN is similar



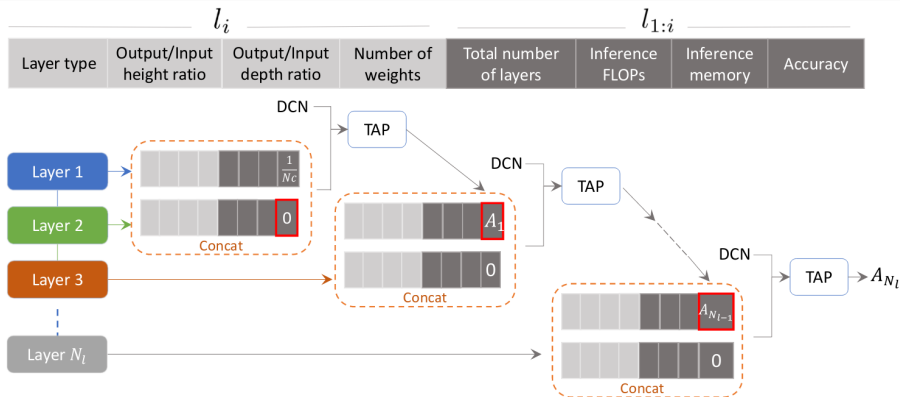
Metadata (LDE)

- ▶ 11 publicly available datasets and 8 datasets generated from ImageNet.
- ▶ 800 architectures are trained per datasets.
- ▶ Architectures are trained incrementally, adding one layer at a time.

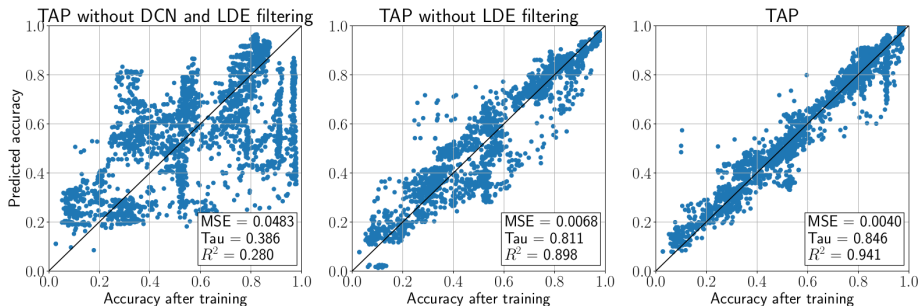
Surrogate Model

- ▶ Surrogate model (TAP) is designed using two stacked LSTMs
- ▶ Datasets with DCN similar to given dataset are selected.
- ▶ TAP is trained on the corresponding metadata.
- ▶ Architecture Encoding and DCN are inputs

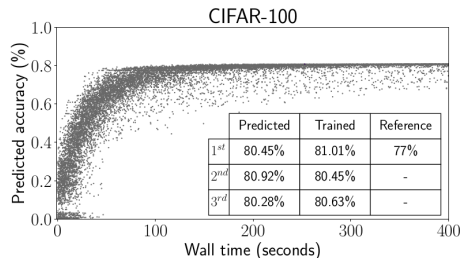
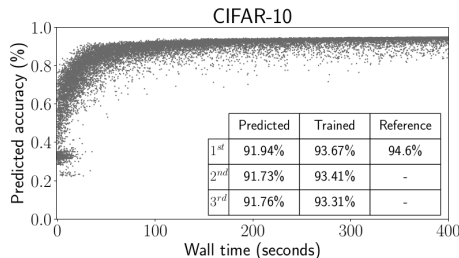
Encoding and Architecture Search



TAP Predictions



TAPAS Simulated Search



- ▶ TAPAS simulates large-scale evolution of image classifiers algorithm¹
- ▶ The algorithm took 250 hours

¹Esteban Real et al. "Large-Scale Evolution of Image Classifiers". In: *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*. 2017, pp. 2902–2911

Rapid Neural Architecture Search

- ▶ Very similar to the idea of TAPAS:
- ▶ Conditioned on the raw data, a graph generator samples promising candidates.
- ▶ A surrogate predicts the performance of an architecture given the dataset.
- ▶ Surrogate and graph generator are pretrained on related data.

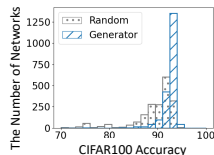
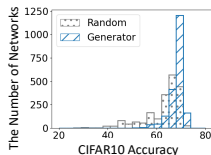
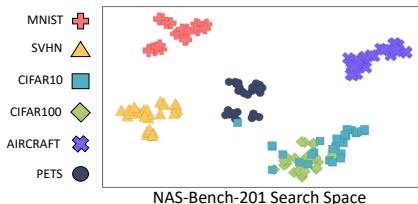
Hayeon Lee, Eunyoung Hyung, and Sung Ju Hwang. "Rapid Neural Architecture Search by Learning to Generate Graphs from Datasets". In: *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. 2021

RapidNAS - Architecture Sampling

For each dataset D_i we have architecture, validation score tuples (a_j, y_j) .

Architecture Sampling

- ▶ Jointly learn a dataset encoder and architecture sampler by maximizing the ELBO.
 - ▶ Dataset encoder: maps a dataset D_i to an embedding z_i .
 - ▶ Architecture sampler: samples architectures conditioned on z_i
- ▶ Training data limited to good performing architectures.
- ▶ The objective is to sample architectures which are good architectures for D_i .



RapidNAS - Surrogate Model

For each dataset D_i we have architecture, validation score tuples (a_j, y_j) .

Surrogate Model

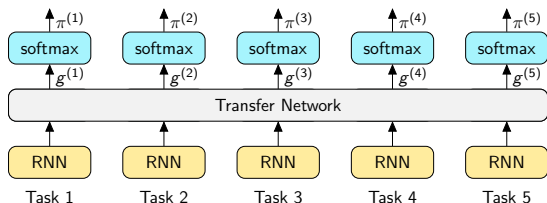
- ▶ Learn a dataset encoder, graph encoder and validation score predictor by minimizing the squared error.
- ▶ Predictor can be applied across datasets since it takes the raw data as input.

RapidNAS - Results

Target Dataset	NAS Method	Params (M)	Search Time (GPU Sec)	Speed Up	Search Cost (\$)	Accuracy (%)
CIFAR-10	REA	-	$0.02+T$	-	-	$93.92_{\pm 0.30}$
	BOHB	-	$3.59+T$	-	-	$93.61_{\pm 0.52}$
	RSPS	-	10200	$147\times$	4.13	$84.07_{\pm 3.61}$
	SETN	-	30200	$437\times$	12.25	$87.64_{\pm 0.00}$
	GDAS	-	25077	$363\times$	10.17	$93.61_{\pm 0.09}$
	PC-DARTS	1.17	10395	$150\times$	4.21	$93.66_{\pm 0.17}$
	DrNAS	1.53	21760	$315\times$	8.82	$94.36_{\pm 0.00}$
	RapidNAS	1.11	69	$1\times$	0.028	$94.37_{\pm 0.03}$
	Baseline	0.86	0			$94.02_{\pm 0.06}$
CIFAR-100	REA	-	$0.02+T$	-	-	$71.84_{\pm 0.99}$
	BOHB	-	$3.59+T$	-	-	$70.85_{\pm 1.28}$
	RSPS	-	18841	$196\times$	7.64	$52.31_{\pm 5.77}$
	SETN	-	58808	$612\times$	23.85	$59.09_{\pm 0.24}$
	GDAS	-	51580	$537\times$	20.91	$70.70_{\pm 0.30}$
	PC-DARTS	0.26	19951	$207\times$	8.09	$66.64_{\pm 2.34}$
	DrNAS	1.20	34529	$359\times$	14.00	$73.51_{\pm 0.00}$
	RapidNAS	1.07	96	$1\times$	0.039	$73.51_{\pm 0.00}$
	Baseline	0.86	0			$73.00_{\pm 0.09}$

XferNAS

- ▶ Warmstart:
 - ▶ Learn an initial policy which does better than random.
- ▶ Minimally invasive:
 - ▶ Easy integration.
 - ▶ Converge to the original NAS optimizer's behavior.
- ▶ Solution: share weights across tasks.



Martin Wistuba. "XferNAS: Transfer Neural Architecture Search". In: *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2020, Ghent, Belgium, September 14-18, 2020*

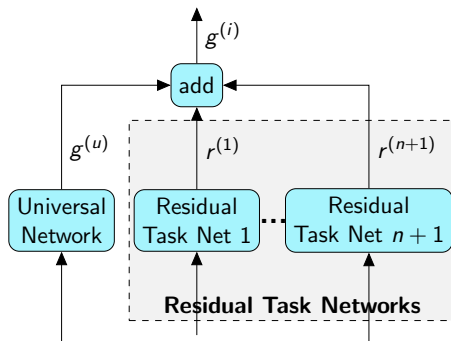
Transfer Network

Core idea is to separate the task-dependent function $g^{(i)}$ into

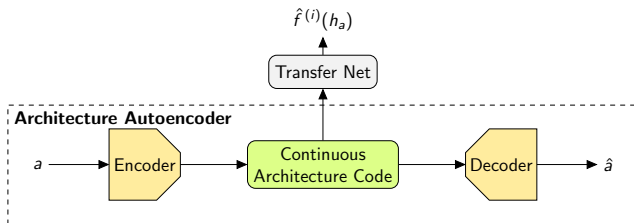
- ▶ a universal function $g^{(u)}$ (warmstart initialization) and
- ▶ a task-dependent residual $r^{(i)}$.

Thus,

$$g^{(i)} = g^{(u)} + r^{(i)} . \quad (16)$$



XferNAS

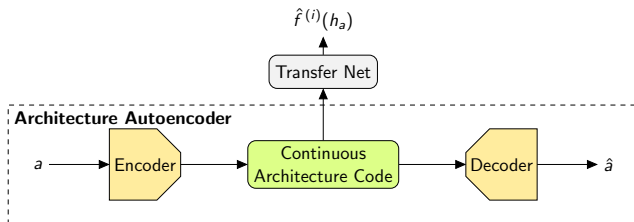
**Example:** Integration into NAO.

- ▶ Auto-encoder with surrogate model \hat{f} that predicts the accuracy of an architecture based on its code h .

$$L = \alpha L_{\text{pred}} + (1 - \alpha) L_{\text{rec}} \quad (17)$$

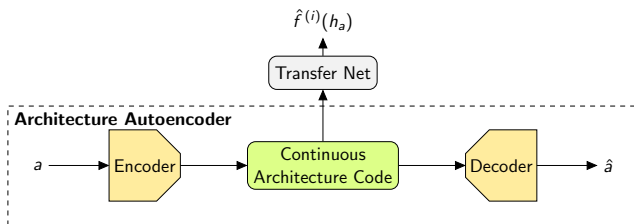
- ▶ L_{pred} : error when predicting accuracy.
- ▶ L_{rec} : auto-encoder reconstruction loss.

XferNAS - Search



1. Solve $h_a^* = \arg \max_{h_a} \hat{f}^{(i)}(h_a)$.
2. Estimate $a^* = \text{Decoder}(h_a^*)$.
3. Evaluate $f^{(i)}(a^*)$.
4. Update the prediction model.
5. Go to 1.

XferNAS vs. NAO



Advantages of XferNAS over NAO:

- ▶ Auto-encoder is trained at the beginning of the search.
- ▶ Knowledge is leverage to warmstart the search.

Results on CIFAR-10

Model	F	#op	Err	#pms	M	GPU Days
NASNet-A	32	13	3.41	3.3M	20000	2000
AmoebaNet-B	36	19	3.37	2.8M	27000	3150
AmoebaNet-B (c/o)	128	19	2.13	34.9M	27000	3150
PNAS	48	8	3.41	3.2M	1280	225
NAONet	36	11	3.18	10.6M	1000	200
NAONet (c/o)	128	11	2.11	128M	1000	200
TAPAS	/	/	6.33	2.7M	1	0
RapidNAS	/	/	5.63	1.1M	30	≈ 1.5
Best on CIFAR-100	32	19	4.14	6.1M	200	/
XferNASNet	32	19	3.37	4.5M	33	6
XferNASNet (c/o)	32	19	2.70	4.5M	33	6
XferNASNet	64	19	3.11	17.5M	33	6
XferNASNet (c/o)	64	19	2.19	17.5M	33	6
XferNASNet (c/o)	128	19	1.99	69.5M	33	6

Results on CIFAR-10

Model	F	#op	Err	#pms	M	GPU Days
NASNet-A	32	13	3.41	3.3M	20000	2000
AmoebaNet-B	36	19	3.37	2.8M	27000	3150
AmoebaNet-B (c/o)	128	19	2.13	34.9M	27000	3150
PNAS	48	8	3.41	3.2M	1280	225
NAONet	36	11	3.18	10.6M	1000	200
NAONet (c/o)	128	11	2.11	128M	1000	200
TAPAS	/	/	6.33	2.7M	1	0
RapidNAS	/	/	5.63	1.1M	30	≈ 1.5
Best on CIFAR-100	32	19	4.14	6.1M	200	/
XferNASNet	32	19	3.37	4.5M	33	6
XferNASNet (c/o)	32	19	2.70	4.5M	33	6
XferNASNet	64	19	3.11	17.5M	33	6
XferNASNet (c/o)	64	19	2.19	17.5M	33	6
XferNASNet (c/o)	128	19	1.99	69.5M	33	6

Results on CIFAR-10

Model	F	#op	Err	#pms	M	GPU Days
NASNet-A	32	13	3.41	3.3M	20000	2000
AmoebaNet-B	36	19	3.37	2.8M	27000	3150
AmoebaNet-B (c/o)	128	19	2.13	34.9M	27000	3150
PNAS	48	8	3.41	3.2M	1280	225
NAONet	36	11	3.18	10.6M	1000	200
NAONet (c/o)	128	11	2.11	128M	1000	200
TAPAS	/	/	6.33	2.7M	1	0
RapidNAS	/	/	5.63	1.1M	30	≈ 1.5
Best on CIFAR-100	32	19	4.14	6.1M	200	/
XferNASNet	32	19	3.37	4.5M	33	6
XferNASNet (c/o)	32	19	2.70	4.5M	33	6
XferNASNet	64	19	3.11	17.5M	33	6
XferNASNet (c/o)	64	19	2.19	17.5M	33	6
XferNASNet (c/o)	128	19	1.99	69.5M	33	6

Results on CIFAR-10

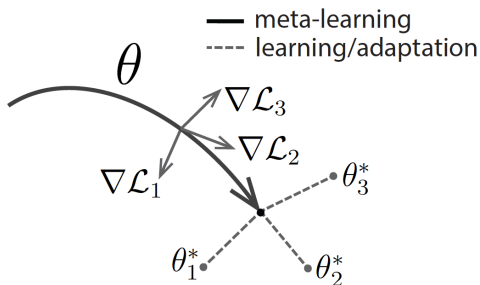
Model	F	#op	Err	#pms	M	GPU Days
NASNet-A	32	13	3.41	3.3M	20000	2000
AmoebaNet-B	36	19	3.37	2.8M	27000	3150
AmoebaNet-B (c/o)	128	19	2.13	34.9M	27000	3150
PNAS	48	8	3.41	3.2M	1280	225
NAONet	36	11	3.18	10.6M	1000	200
NAONet (c/o)	128	11	2.11	128M	1000	200
TAPAS	/	/	6.33	2.7M	1	0
RapidNAS	/	/	5.63	1.1M	30	≈ 1.5
Best on CIFAR-100	32	19	4.14	6.1M	200	/
XferNASNet	32	19	3.37	4.5M	33	6
XferNASNet (c/o)	32	19	2.70	4.5M	33	6
XferNASNet	64	19	3.11	17.5M	33	6
XferNASNet (c/o)	64	19	2.19	17.5M	33	6
XferNASNet (c/o)	128	19	1.99	69.5M	33	6

Outline

1. Introduction
2. One-Shot Architecture Search
3. Zero-Shot NAS
4. Transfer Learning for NAS
 - 4.1 Transfer NAS
 - 4.2 Few-Shot NAS
 - 4.3 Learning Curve Ranking
5. Conclusions

Model-Agnostic Meta-Learning

- ▶ Model learns from all the tasks
- ▶ Learn a representation that requires only few steps to the optimal representation for each task
- ▶ Performs well for few-shot learning problems



Chelsea Finn, Pieter Abbeel, and Sergey Levine. “Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks”. In: *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*. 2017, pp. 1126–1135

MAML Algorithm

Algorithm 3 Model-Agnostic Meta-Learning

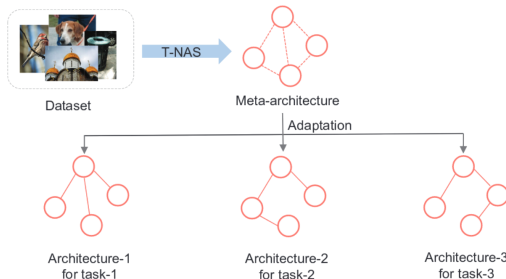
Input: $p(\mathcal{T})$: distribution over tasks

Input: β, γ : step size hyperparameters

- 1: randomly initialize θ
 - 2: **while** not done **do**
 - 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
 - 4: **for all** \mathcal{T}_i **do**
 - 5: $\theta'_i = \theta - \beta \nabla_{\theta} \mathcal{L}_{\mathcal{T}_i}(f_{\theta})$
 - 6: $\theta \leftarrow \theta - \gamma \nabla_{\theta} \sum_{\mathcal{T}_i \sim p(\mathcal{T})} \mathcal{L}_{\mathcal{T}_i}(f_{\theta'_i})$
-

T-NAS

- ▶ Objective: Given multiple tasks, learn a meta-architecture
- ▶ Using bilevel optimization combined with MAML to estimate α and network weights w
- ▶ Finetune both parameters for each new task



Dongze Lian et al. “Towards Fast Adaptation of Neural Architectures with Meta Learning”. In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. 2020

T-NAS

Algorithm 4 T-NAS

Input: $p(\mathcal{T})$: distribution over tasks

- 1: randomly initialize α and w
 - 2: **while** not done **do**
 - 3: Sample batch of tasks $\mathcal{T}_i \sim p(\mathcal{T})$
 - 4: **for all** \mathcal{T}_i **do**
 - 5: Alternately update α' and w'
 - 6: Update α and w
-

Architecture Evaluation

Methods	Arch.	#Param.	1-shot	5-shot
Matching nets (Vinyals et al., 2016)	4CONV	32.9K	$43.44 \pm 0.77\%$	$55.31 \pm 0.73\%$
ProtoNets (Snell et al., 2017)	4CONV	32.9K	$49.42 \pm 0.78\%$	$68.20 \pm 0.66\%$
Meta-LSTM (Ravi & Larochelle, 2017)	4CONV	32.9K	$43.56 \pm 0.84\%$	$60.60 \pm 0.71\%$
Bilevel (Franceschi et al., 2018)	4CONV	32.9K	$50.54 \pm 0.85\%$	$64.53 \pm 0.68\%$
CompareNets (Sung et al., 2018)	4CONV	32.9K	$50.44 \pm 0.82\%$	$65.32 \pm 0.70\%$
LLAMA (Grant et al., 2018)	4CONV	32.9K	$49.40 \pm 1.83\%$	-
MAML (Finn et al., 2017)	4CONV	32.9K	$48.70 \pm 1.84\%$	$63.11 \pm 0.92\%$
MAML (first-order) (Finn et al., 2017)	4CONV	32.9K	$48.07 \pm 1.75\%$	$63.15 \pm 0.91\%$
MAML++ (Antoniou et al., 2019)	4CONV	32.9K	$52.15 \pm 0.26\%$	$68.32 \pm 0.44\%$
Auto-Meta (small) (Kim et al., 2018)	Cell	28/28 K	$49.58 \pm 0.20\%$	$65.09 \pm 0.24\%$
Auto-Meta (large) (Kim et al., 2018)	Cell	98.7/94.0 K	$51.16 \pm 0.17\%$	$69.18 \pm 0.14\%$
BASE (Softmax) (Shaw et al., 2018)	Cell	1200K	-	$65.40 \pm 0.74\%$
BASE (Gumbel-Softmax) (Shaw et al., 2018)	Cell	1200K	-	$66.20 \pm 0.70\%$
Auto-MAML (ours)	Cell	23.2/26.1 K	$51.23 \pm 1.76\%$	$64.10 \pm 1.12\%$
T-NAS (ours)	Cell	24.3/26.5 K*	$52.84 \pm 1.41\%$	$67.88 \pm 0.92\%$
T-NAS++ (ours)	Cell	24.3/26.5 K*	$54.11 \pm 1.35\%$	$69.59 \pm 0.85\%$

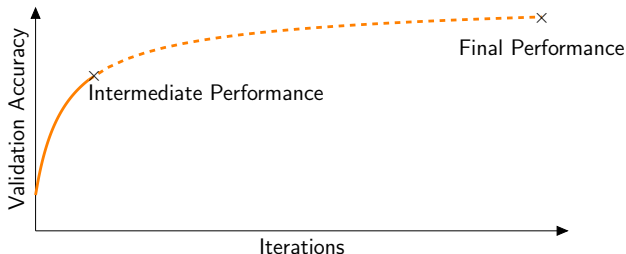


Outline

1. Introduction
2. One-Shot Architecture Search
3. Zero-Shot NAS
4. Transfer Learning for NAS
 - 4.1 Transfer NAS
 - 4.2 Few-Shot NAS
 - 4.3 Learning Curve Ranking
5. Conclusions

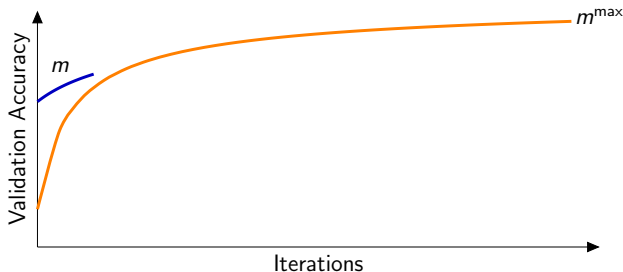
Motivation

- ▶ Hyperparameter and neural architecture optimization are computationally expensive.
- ▶ Human experts decrease this effort by monitoring the model's learning curve and terminate options early that are unlikely to improve over the currently best solutions.
- ▶ With the rise of AutoML, a system that is able to perform this automatically is desired.



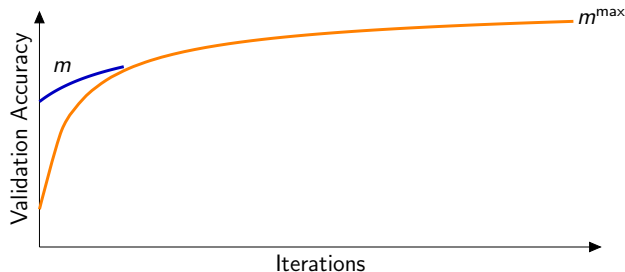
Use Simple Statistics

1. Use median/mean or last value in learning curve to make decision.



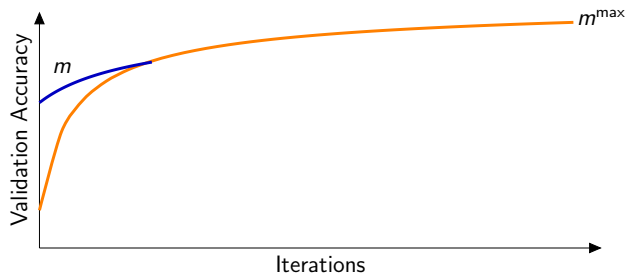
Use Simple Statistics

1. Use median/mean or last value in learning curve to make decision.



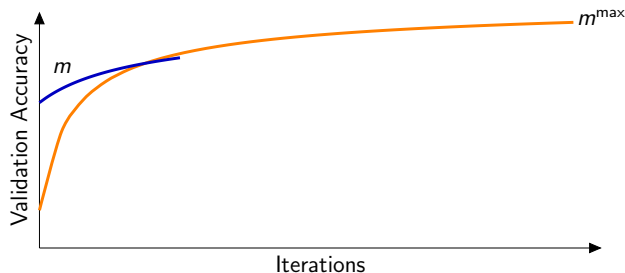
Use Simple Statistics

1. Use median/mean or last value in learning curve to make decision.



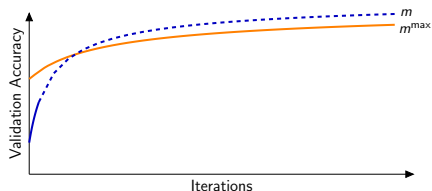
Use Simple Statistics

1. Use median/mean or last value in learning curve to make decision.



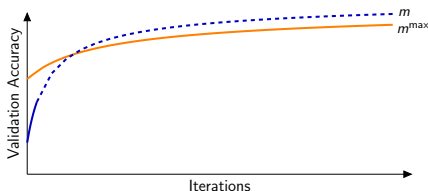
Simple Statistics - Problems

- Late bloomers will not be considered.

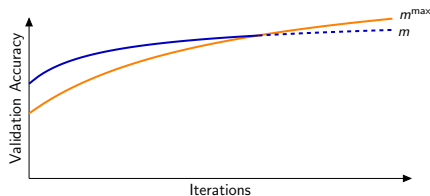


Simple Statistics - Problems

- Late bloomers will not be considered.



- Quick learners will be considered unnecessarily long.





Learning Curves Prediction

- ▶ Given a partial learning curve, predict the final performance.
- ▶ Use this prediction to estimate $p(m > m^{\max})$.
- ▶ Terminate all runs with $p(m > m^{\max}) \leq \delta$

Learning Curves Ranking

- ▶ Proposing to predict $p(m > m^{\max})$ directly.
- ▶ Defining the probability that m_i is better than m_j as

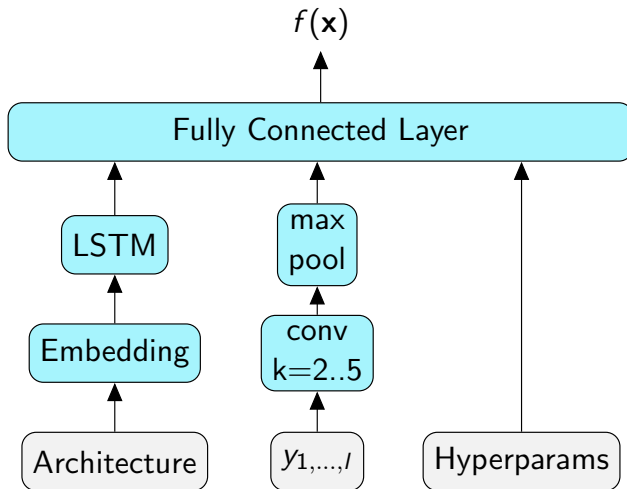
$$p(m_i > m_j) = \hat{p}_{i,j} = \frac{e^{f(\mathbf{x}_i) - f(\mathbf{x}_j)}}{1 + e^{f(\mathbf{x}_i) - f(\mathbf{x}_j)}}. \quad (18)$$

- ▶ Minimize the cross-entropy loss

$$\sum_{i,j} -p_{i,j} \log \hat{p}_{i,j} - (1 - p_{i,j}) \log(1 - \hat{p}_{i,j}) \quad (19)$$

Martin Wistuba and Tejaswini Pedapati. "Learning to Rank Learning Curves". In: *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 12-18 June 2020, Vienna, Austria*. 2020

Modelling f





Learning Curves Ranking with Transfer Learning

- ▶ Learning requires data which is not available.



Learning Curves Ranking with Transfer Learning

- ▶ Learning requires data which is not available.
- ▶ Solution 1: Do not learn, only consider given partial learning curve.



Learning Curves Ranking with Transfer Learning

- ▶ Learning requires data which is not available.
- ▶ Solution 1: Do not learn, only consider given partial learning curve.
- ▶ Solution 2: First collect sufficient learning curves and then train your model.

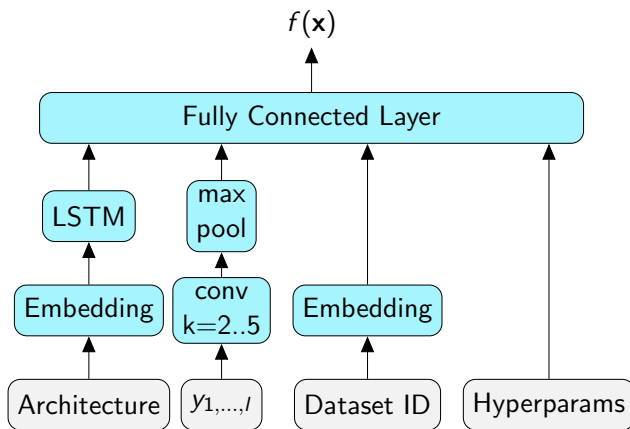


Learning Curves Ranking with Transfer Learning

- ▶ Learning requires data which is not available.
- ▶ Solution 1: Do not learn, only consider given partial learning curve.
- ▶ Solution 2: First collect sufficient learning curves and then train your model.
- ▶ Proposal: Use transfer learning to reduce this problem.

Considering Transfer Learning in our Modelling

To account for transfer learning, an embedding per dataset is added.

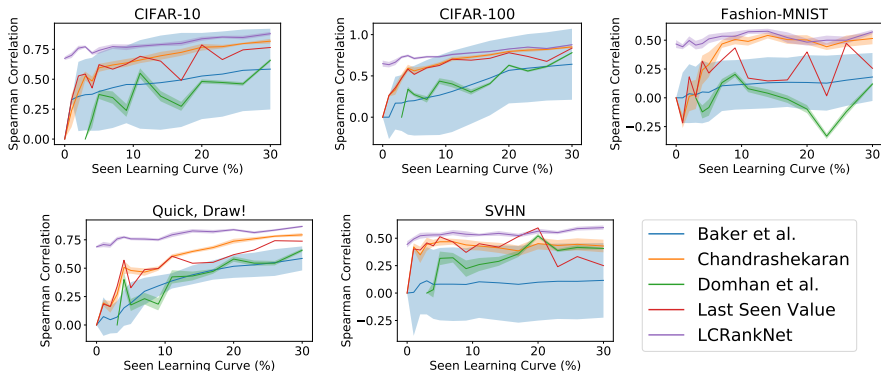




Setup

- ▶ Experiments are conducted on five different datasets: CIFAR-10, CIFAR-100, Fashion-MNIST, Quickdraw, and SVHN.
- ▶ To create the meta-knowledge, 200 architectures per dataset are chosen at random from the NASNet search space (i.e. 1000 unique architectures) and train it for 100 epochs.
- ▶ Experiments are conducted in a leave-one-dataset-out cross-validation.

Ranking Performance



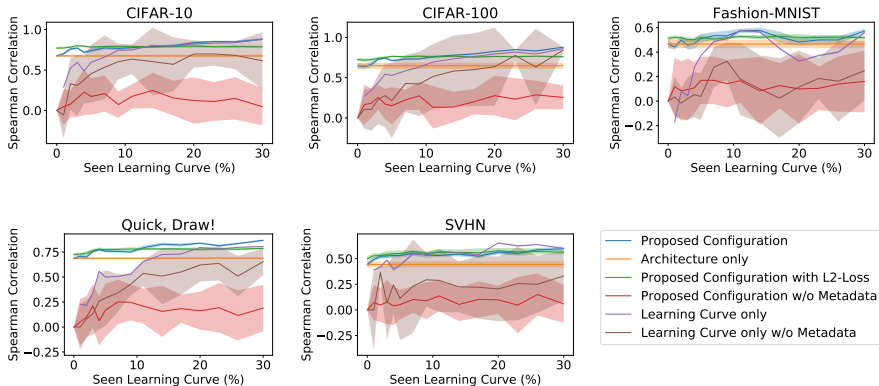
Random Search

Random Neural Architecture Search with Early Stopping.

- Regret: Difference of best solution and best solution without early stopping.
- Time: GPU time in hours.

Method	CIFAR-10		CIFAR-100		Fashion		Quickdraw		SVHN	
	REGR.	TIME	REGR.	TIME	REGR.	TIME	REGR.	TIME	REGR.	TIME
NO EARLY TERMINATION	0.00	1023	0.00	1021	0.00	1218	0.00	1045	0.00	1485
DOMHAN ET AL.	0.56	346	0.82	326	0.00	460	0.44	331	0.28	471
HYPERBAND	0.22	106	0.78	102	0.32	132	0.54	109	0.00	156
BAKER ET AL.	0.00	89	0.00	77	0.00	129	0.00	107	0.00	241
SUCCESSIVE HALVING	0.62	62	0.00	54	0.18	70	0.40	60	0.28	88
CHANDRA-SHEKARAN	0.62	30	0.00	35	0.28	41	0.30	82	0.06	164
LCRANKNET	0.22	20	0.00	11	0.10	19	0.00	28	0.10	74

Component Analysis





Conclusions

- ▶ Zero-shot may be an efficient way to search for neural architectures.
- ▶ Transfer learning for Neural Architecture Search has been explored in various ways.
 - ▶ By means of special neural architecture search methods,
 - ▶ meta-learning,
 - ▶ and early termination techniques for incremental model training.
- ▶ All imply that it can be used to significantly decrease the computational effort for NAS.
- ▶ Yet, it is a relatively unexplored research topic.

Outline

1. Introduction
2. One-Shot Architecture Search
3. Zero-Shot NAS
4. Transfer Learning for NAS
5. Conclusions

Final Conclusions

- ▶ A common search space for NAS.
- ▶ Various efficient optimizers based on parameter sharing and differentiable architecture search.
- ▶ Discussion of several problems being faced with these very methods.
- ▶ A deep dive on zero-shot and transfer learning for NAS.

Thank you for your attention.

Survey Paper: <https://arxiv.org/abs/1905.01392>

References I

- [1] Gabriel Bender et al. “Understanding and Simplifying One-Shot Architecture Search”. In: *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*. 2018, pp. 549–558.
- [2] Simone Bianco et al. “Benchmark Analysis of Representative Deep Neural Network Architectures”. In: *IEEE Access* 6 (2018), pp. 64270–64277.
- [3] Han Cai, Ligeng Zhu, and Song Han. “ProxylessNAS: Direct Neural Architecture Search on Target Task and Hardware”. In: *Proceedings of the International Conference on Learning Representations, ICLR 2019, New Orleans, Louisiana, USA*. 2019.

References II

- [4] Han Cai et al. “Once-for-All: Train One Network and Specialize it for Efficient Deployment”. In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. 2020.
- [5] Francesco Paolo Casale, Jonathan Gordon, and Nicolo Fusi. “Probabilistic Neural Architecture Search”. In: *CoRR* abs/1902.05116 (2019).
- [6] Xiangxiang Chu et al. “Fair DARTS: Eliminating Unfair Advantages in Differentiable Architecture Search”. In: *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part XV*. 2020, pp. 465–480.

References III

- [7] Chelsea Finn, Pieter Abbeel, and Sergey Levine. “Model-Agnostic Meta-Learning for Fast Adaptation of Deep Networks”. In: *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*. 2017, pp. 1126–1135.
- [8] Roxana Istrate et al. “TAPAS: Train-less Accuracy Predictor for Architecture Search”. In: *Proceedings of the Thirty-Third AAAI Conference on Artificial Intelligence, (AAAI-19), Honolulu, Hawaii, USA*. 2019.
- [9] Hayeon Lee, Eunyoung Hyung, and Sung Ju Hwang. “Rapid Neural Architecture Search by Learning to Generate Graphs from Datasets”. In: *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. 2021.

References IV

- [10] Dongze Lian et al. “Towards Fast Adaptation of Neural Architectures with Meta Learning”. In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. 2020.
- [11] Ming Lin et al. “Zen-NAS: A Zero-Shot NAS for High-Performance Image Recognition”. In: *2021 IEEE/CVF International Conference on Computer Vision, ICCV 2021, Montreal, QC, Canada, October 10-17, 2021*. 2021, pp. 337–346.
- [12] Chenxi Liu et al. “Are Labels Necessary for Neural Architecture Search?” In: *Computer Vision - ECCV 2020 - 16th European Conference, Glasgow, UK, August 23-28, 2020, Proceedings, Part IV*. Ed. by Andrea Vedaldi et al. Vol. 12349. Lecture Notes in Computer Science. Springer, 2020, pp. 798–813.

References V

- [13] Hanxiao Liu, Karen Simonyan, and Yiming Yang. “DARTS: Differentiable Architecture Search”. In: *Proceedings of the International Conference on Learning Representations, ICLR 2019, New Orleans, Louisiana, USA*. 2019.
- [14] Renqian Luo et al. “Neural Architecture Optimization”. In: *Advances in Neural Information Processing Systems 31: Annual Conference on Neural Information Processing Systems 2018, NeurIPS 2018, 3-8 December 2018, Montréal, Canada*. 2018, pp. 7827–7838.
- [15] Joe Mellor et al. “Neural Architecture Search without Training”. In: *Proceedings of the 38th International Conference on Machine Learning, ICML 2021, 18-24 July 2021, Virtual Event*. 2021, pp. 7588–7598.

References VI

- [16] Asaf Noy et al. “ASAP: Architecture Search, Anneal and Prune”. In: *The 23rd International Conference on Artificial Intelligence and Statistics, AISTATS 2020, 26-28 August 2020, Online [Palermo, Sicily, Italy]*. 2020, pp. 493–503.
- [17] Hieu Pham et al. “Efficient Neural Architecture Search via Parameter Sharing”. In: *Proceedings of the 35th International Conference on Machine Learning, ICML 2018, Stockholmsmässan, Stockholm, Sweden, July 10-15, 2018*. 2018, pp. 4092–4101.
- [18] Maithra Raghu et al. “On the Expressive Power of Deep Neural Networks”. In: *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*. 2017, pp. 2847–2854.

References VII

- [19] Esteban Real et al. “Large-Scale Evolution of Image Classifiers”. In: *Proceedings of the 34th International Conference on Machine Learning, ICML 2017, Sydney, NSW, Australia, 6-11 August 2017*. 2017, pp. 2902–2911.
- [20] Christian Szegedy et al. “Inception-v4, Inception-ResNet and the Impact of Residual Connections on Learning”. In: *Proceedings of the Thirty-First AAAI Conference on Artificial Intelligence, February 4-9, 2017, San Francisco, California, USA*. 2017, pp. 4278–4284.
- [21] Ruochen Wang et al. “Rethinking Architecture Selection in Differentiable NAS”. In: *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3-7, 2021*. OpenReview.net, 2021.

References VIII

- [22] Martin Wistuba. “XferNAS: Transfer Neural Architecture Search”. In: *Machine Learning and Knowledge Discovery in Databases - European Conference, ECML PKDD 2020, Ghent, Belgium, September 14-18, 2020*.
- [23] Martin Wistuba and Tejaswini Pedapati. “Learning to Rank Learning Curves”. In: *Proceedings of the 37th International Conference on Machine Learning, ICML 2020, 12-18 June 2020, Vienna, Austria. 2020*.
- [24] Yuhui Xu et al. “PC-DARTS: Partial Channel Connections for Memory-Efficient Architecture Search”. In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. 2020*.

References IX

- [25] Kaicheng Yu et al. "Evaluating The Search Phase of Neural Architecture Search". In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. 2020.
- [26] Arber Zela et al. "Understanding and Robustifying Differentiable Architecture Search". In: *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. 2020.