

# BPNAS: Bayesian Progressive Neural Architecture Search

Hyunwoong Chang<sup>1,\*</sup>, Anirban Samaddar<sup>2,\*</sup>, Sandeep Madireddy<sup>2,†</sup>

<sup>1</sup> Texas A&M University, <sup>2</sup> Argonne National Laboratory,

\* Equal contribution, † Corresponding author



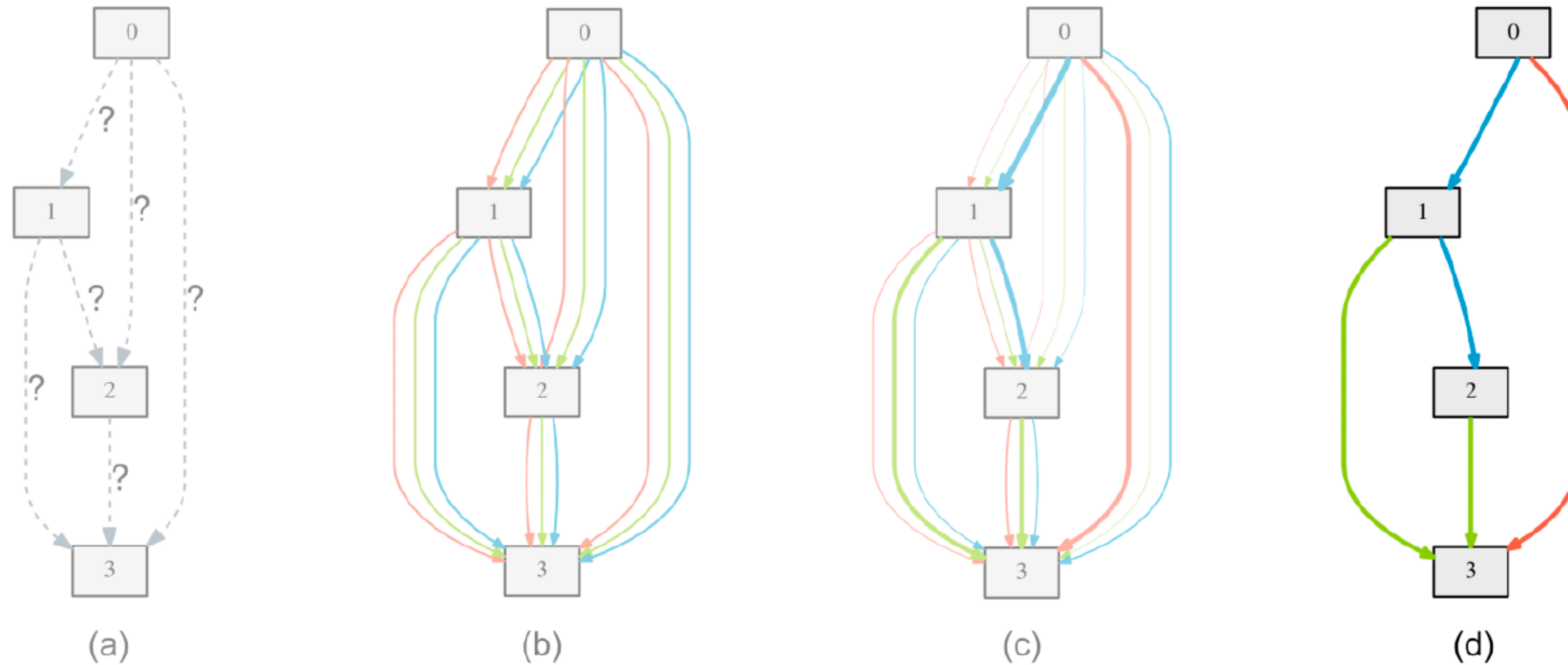
**Anirban Samaddar**



# Introduction

- Objective of NAS is to automate designing of neural networks
- Traditional NAS methods use evolutionary algorithms and reinforcement learning (heavy computation burden)
- Differentiable NAS<sup>1</sup> is an attractive alternative that incorporates differentiability in the search process hence increase efficiency

# Differentiable NAS



A cell is a directed acyclic graph consisting  $N$  ( $=4$  here) nodes with unknown operations (convolution, skip-connect etc.)

A super-network is assumed where output of each node is an weighted average of  $K$  operations

The weights are learned as a part of the training

An architecture is selected by choosing the operators with the maximum weights

Fig.1: Schematic of differentiable NAS<sup>1</sup>

# Differentiable NAS

- The differentiable NAS objective function:

$$\min_{\alpha} \mathcal{L}_{\text{val}}(\alpha, W^*) \text{ s.t. } W^* = \arg \min_W \mathcal{L}_{\text{train}}(\alpha, W)$$

$\alpha, W$  are the weights of the operators and the weights (and biases) of the network

- The loss is differentiable w.r.t the parameters and therefore they can be learned with existing gradient based tools

# A Bayesian Framework

- We propose a unified Bayesian framework for the architecture search —

$$\pi(\alpha, W, \Gamma) \propto \pi(\alpha, W | \Gamma) \pi(\Gamma)$$

$$\pi(\alpha, W | \Gamma) \propto e^{-\mathcal{L}(\alpha(\Gamma), W(\Gamma) | \Gamma)}$$

$$\pi(\Gamma) \propto \prod_{m=1}^{n_{edge}} \pi(\gamma_m) = \prod_{m=1}^{n_{edge}} e^{-c|\gamma_m|} \mathbf{1}_{\{|\gamma_m| > 0\}}(\gamma_m)$$

$n_{edge}$  is the number of edges, and  $\Gamma$ , a binary matrix, is introduced to select from the space of architectures

- The prior for  $\Gamma$  is set to ensure that it enforces a mode where every edge has one operator ( $|\gamma_m| = 1$ )

# BPNAS algorithm

- Our goal is to learn—  $(\hat{\alpha}^{\text{MAP}}, \hat{W}^{\text{MAP}}, \hat{\Gamma}^{\text{MAP}}) = \arg \max_{\alpha, W, \Gamma} \pi(\alpha, W, \Gamma)$
- Difficult to find this MAP estimate from the joint distribution of discrete and continuous random variable
- We devise an algorithm to prune out less important operators based on the weight matrix  $\alpha$

# BPNAS algorithm

---

## Bayesian Progressive architecture search (BPNAS):

---

**Input:** An initial super-network ( $\Gamma = \mathbf{1}$ ) with architecture parameters  $\alpha$ , weights  $W$ , posterior distribution  $\pi(W, \alpha, \Gamma)$ , and a threshold  $\delta$

**While**  $\|\gamma_m\|_0 > 1$  for some  $k$ :

Update  $W, \alpha$  by stochastic gradient descent

If  $\|\theta(\alpha_m) - \frac{\mathbf{1}_{\|\gamma_m\|_0}}{\|\gamma_m\|_0}\|_2 \geq \delta$  holds:

Select  $(l, m) = \operatorname{argmin}_{\{(l', m') : \Gamma_{l'm'} = 1\}} \alpha_{l'm'}$  and set  $\Gamma_{l'm'} = 0$

Reset  $\alpha$

**Output:** An architecture  $\Gamma$  with a single operation for each edge.

---

- We chose  $\theta(\alpha_m) \sim \text{Dirichlet}(\alpha_m)$

# Architecture Ensemble

- Our Bayesian framework enables efficient method to sample architectures from the posterior
- Upon convergence to an architecture  $\hat{\Gamma}$ , we restart the algorithm  $N$  times by randomly sampling an edge  $s$  and set  $\gamma_s = \mathbf{1}$
- We perform this (in parallel) to build an ensemble of architectures

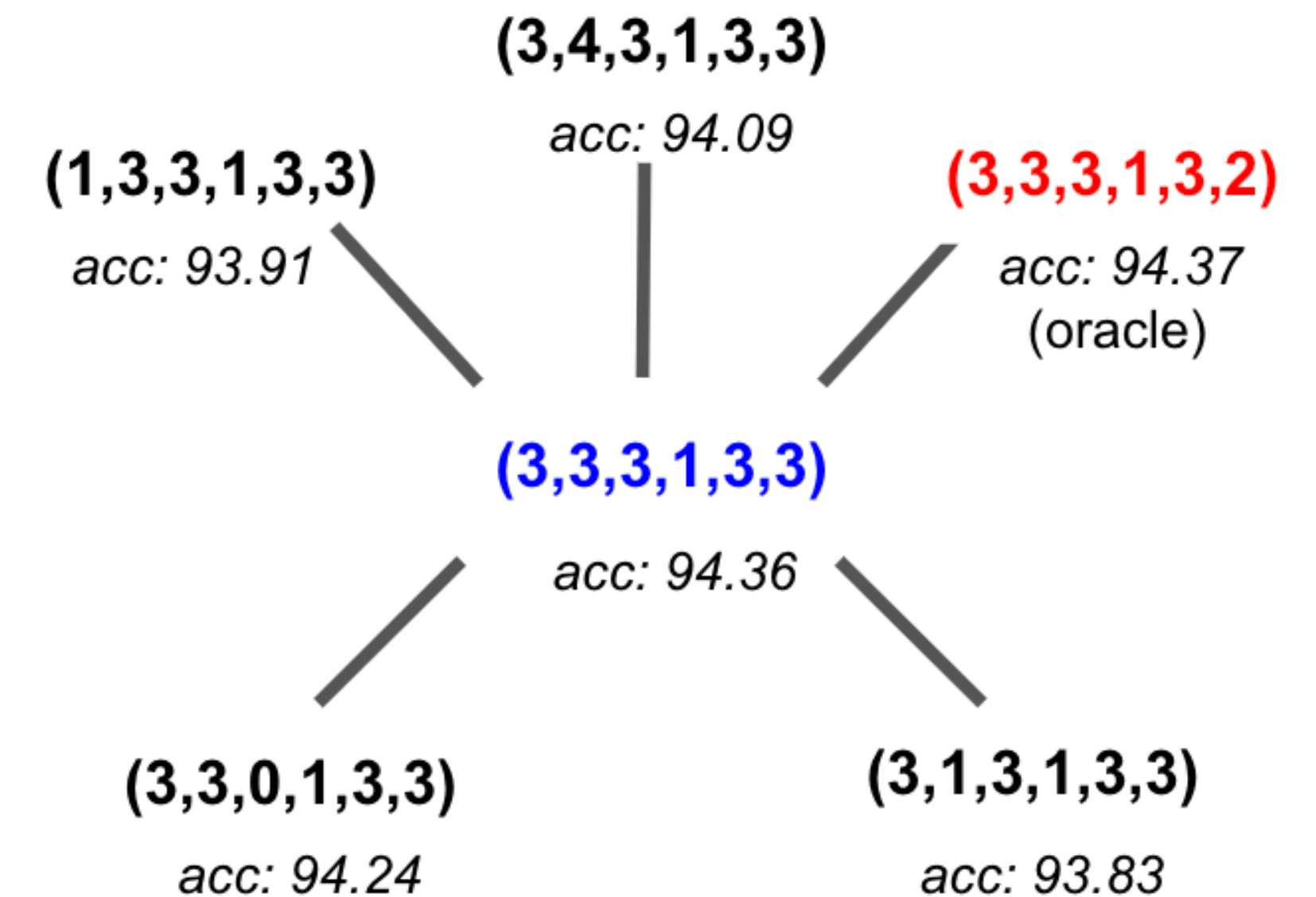


Fig.2: An ensemble in NAS-Bench-2012



# Results

## Architecture search

Methods	CIFAR-10		CIFAR-100		ImageNet	
	Test	Epochs	Test	Epochs	Test	Epochs
<b>ResNet</b>	93.97	100	70.86	100	43.63	100
<b>DrNAS<sup>3</sup></b>	<b>94.36 ± 0.0</b>	100	71.00 ± 1.3	100	<b>46.34 ± 0.0</b>	100
<b>BPNAS</b>	94.18 ± 0.3	<b>63</b>	<b>73.40 ± 0.2</b>	<b>50</b>	<b>46.34 ± 0.0</b>	<b>26</b>
<b>Optimal</b>	94.37	-	73.51	-	47.31	-

- On CIFAR-10 and ImageNet dataset, BPNAS performs similarly to the state-of-the-art while outperforming on CIFAR-100
- Since BPNAS used pruning it converges to the final architecture faster (~70% less epochs on ImageNet)

# Results

## Architecture Ensemble

Methods	CIFAR-10	CIFAR-100	ImageNet
<b>NES-RS<sup>4</sup></b>	94.17 $\pm$ 0.3	74.42 $\pm$ 0.8	45.66 $\pm$ 1.7
<b>NESBS<sup>5</sup></b>	94.08 $\pm$ 0.1	75.00 $\pm$ 0.2	47.32 $\pm$ 0.4
<b>BPNAS</b>	<b>95.10 <math>\pm</math> 0.3</b>	<b>77.47 <math>\pm</math> 1.0</b>	<b>50.27 <math>\pm</math> 0.5</b>

- On all datasets, BPNAS ensemble outperforms all the state-of-the-art architecture ensemble algorithms

Thank you!

# References

1. Liu, H., K. Simonyan, and Y. Yang (2018). “Darts: Differentiable architecture search”. In: arXiv preprint arXiv:1806.09055.
2. Dong, X. and Y. Yang (2020). “Nas-bench-201: Extending the scope of reproducible neural architecture search”. In: arXiv preprint arXiv:2001.00326.
3. Chen, X., R. Wang, M. Cheng, X. Tang, and C.-J. Hsieh (2020). “Drnas: Dirichlet neural architecture search”. In: arXiv preprint arXiv:2006.10355.
4. Zaidi, S., A. Zela, T. Elsken, C. C. Holmes, F. Hutter, and Y. Teh (2021). “Neural ensemble search for uncertainty estimation and dataset shift”. In: Advances in Neural Information Processing Systems 34, pp. 7898–7911.
5. Shu, Y., Y. Chen, Z. Dai, and B. K. H. Low (2022). “Neural ensemble search via Bayesian sampling”. In: Uncertainty in Artificial Intelligence. PMLR, pp. 1803–1812.
6. Zoph, B., V. Vasudevan, J. Shlens, and Q. V. Le (2018). “Learning transferable architectures for scalable image recognition”. In: Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 8697–8710.