

# Adaptation of the architecture of the deep learning model with performance control\*

*S. S. Borodin<sup>1</sup>, K. D. Yakovlev<sup>1</sup>, and O. Y. Bakhteev<sup>1,2</sup>*

{borodin.ss, iakovlev.kd, bakhteev}@phystech.edu

<sup>1</sup> MIPT, Russia

<sup>2</sup> Dorodnicyn Computing Center RAS, Russia

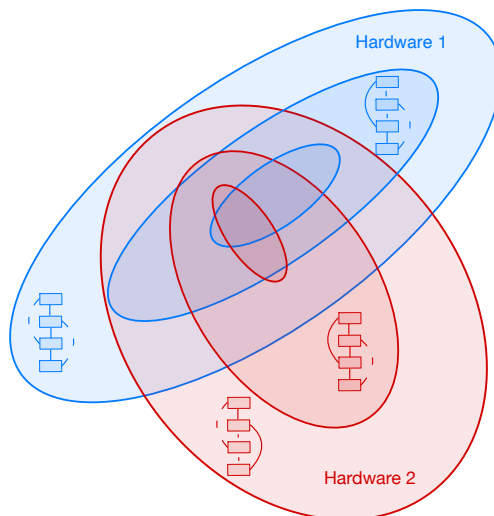
Will be rewritten since problem was corrected.

The paper investigates the problem of deep learning model selection. A method can consider different requirements to optimize performance on a wide range of hardware. The method allows to control trade-off between prediction quality and performance on the intermediate stage of model selection without additional model fitting. The method is based on a modification of a differential architecture search algorithm (DARTS) called DARTS-CC. The method can adapt to benefits and weaknesses of hardware, which is especially important for mobile devices with limited computation budget. To evaluate the performance of the proposed algorithm, we conduct experiments on the Fashion-MNIST and CIFAR-10 datasets using/simulating different hardware (e.g. ...) and compare the resulting architecture with architectures obtained by other neural architecture search methods.

DOI: 10.21469/22233792

## 1 Introduction

Will be rewritten since problem was corrected.



\*The research was supported by the Russian Foundation for Basic Research (grants 00-00-0000 and 00-00-00001).

Discovering architectures for different problems is a difficult task. If you want to find an architecture that considers hardware properties the task becomes inhumanly hard. Optimizing for mobile devices is crucial because of their limited computation budget. With a wide variety of devices it is important to optimize searching not for one, but many different capabilities and requirements.

There are different approaches to finding architectures: evolutionary, RL (reinforced learning) and gradient-based. According to [7] gradient-based show similar prediction quality to other approaches but mainly with lower search cost.

One of the most popular algorithm is DARTS [5]. Main idea of the algorithm is to represent NN as DAG. Where each connection can be of the “candidate operations (e.g., convolution, max pooling, zero)” with different possibilities.

FBNet [3] uses this algorithm to optimize latency. Model shows great results decreasing FLOPS and CPU latency almost without decreasing accuracy.

DARTS-CC [2] uses this algorithm to optimize model complexity with trade-off control. This modification of DARTS can produce different architectures in one fitting run with control to trade-off between quality and complexity.

We propose a method that will combine approaches of DARTS-CC [2] and FBNet [3] to obtain multiple architectures optimized for different requirements with different trade-offs between prediction quality and performance in one run.

## 2 Problem statement

Given a model with structure  $\Gamma = (V, E)$ , where  $E$  is a set of atomic operations such as convolution, pooling, activation etc. The task is to find a subset of edges for which the pruned model gives comparable quality with significant speed-up on a particular device.

Optimization problem:

$$\mathbb{E}_\lambda \mathbb{E}_\gamma \left[ \mathcal{L}_{\text{valid}}(\gamma(\lambda)) + \lambda \sum_{e \in E} T(e) \gamma(\lambda)_e \right] \rightarrow \min_{\gamma(\lambda)},$$

where  $\mathcal{L}_{\text{val}}(\gamma)$  is a loss function for a pruned network which is a shorthand of  $\mathcal{L}_{\text{valid}}(\hat{y}_\gamma(x), y)$ ,  $\lambda \in \mathbb{R}$  is regularization parameter,  $\gamma(\lambda) \in \{0, 1\}^{|E|}$ ,  $T(e)$  — time to execute corresponding atomic operation. Since optimization on discrete space is not differentiable we use Gumbel-Softmax approximation [2].

**Theorem 1.** Let  $\Gamma_{\text{valid}} \subset \{0, 1\}^D$  is a set for which the computational graph is not broken. Let also loss function be in the following form

$$\mathcal{L}_{\text{valid}}(\gamma) = \frac{1}{n} \sum_{i=1}^n \ell(f(\gamma), y_i)$$

where  $f(\gamma) \in \Delta$ . Let also

$$\inf_{p \in \Delta} \frac{1}{n} \sum_{i=1}^n \ell(p, y_i) > \sup_{\gamma \in \Gamma_{\text{valid}}} \mathcal{L}_{\text{valid}}(\gamma).$$

Then the following statement is true. Suppose that  $\gamma^* \notin \Gamma_{\text{valid}}$  is a solution of the following problem

$$\gamma^* = \arg \min_{\gamma} \mathcal{L}_{\text{valid}}(\gamma) + \lambda \mathbf{T}^\top \gamma.$$

Then

$$\lambda \geq \frac{\inf_{p \in \Delta} \frac{1}{n} \sum_{i=1}^n \ell(p, y_i) - \sup_{\gamma \in \Gamma_{\text{valid}}} \mathcal{L}_{\text{valid}}(\gamma)}{\mathbf{T}^\top \mathbf{1}} =: \Lambda > 0.$$

**Proof.**  $\gamma^*$  the solution of the problem. Hence,  $\forall \gamma' \in \Gamma_{\text{valid}}$

$$\mathcal{L}_{\text{valid}}(\gamma^*) + \lambda \mathbf{T}^\top \gamma^* \leq \mathcal{L}_{\text{valid}}(\gamma') + \lambda \mathbf{T}^\top \gamma' \leq \sup_{\gamma \in \Gamma_{\text{valid}}} \mathcal{L}_{\text{valid}}(\gamma) + \lambda \mathbf{T}^\top \mathbf{1}.$$

Since  $\forall i \ \mathbf{T}_i \geq 0$  the following statement is true

$$\mathcal{L}_{\text{valid}}(\gamma^*) + \lambda \mathbf{T}^\top \gamma^* = \inf_{\gamma \in \{0,1\}^D} \frac{1}{n} \sum_{i=1}^n \ell(f(\gamma), y_i) + \lambda \mathbf{T}^\top \gamma \geq \inf_{p \in \Delta} \frac{1}{n} \sum_{i=1}^n \ell(p, y_i).$$

Thus,

$$\sup_{\gamma \in \Gamma_{\text{valid}}} \mathcal{L}_{\text{valid}}(\gamma) + \lambda \mathbf{T}^\top \mathbf{1} \geq \inf_{p \in \Delta} \frac{1}{n} \sum_{i=1}^n \ell(p, y_i).$$

31

32 **Corollary 1.** For  $\lambda < \Lambda$  the solution lies in  $\Gamma_{\text{valid}}$ . ■

---

### Algorithm 1 hypernet fitting

---

```

1: load pre-trained model
2: initialize hypernet for the model
3: while hypernet( $\lambda = 0$ ) has random accuracy do
4:   Sample  $\gamma_{\text{sampled}} \sim \mathcal{N}(\text{bias}, 0)$ 
5:    $\gamma \leftarrow \text{hypernet}(\lambda)$ 
6:   Optimize hypernet for loss  $\|\gamma - \gamma_{\text{sampled}}\|_2^2$ 
7: end while
8: while not converged do
9:   Sample  $\lambda \sim U[0, \Lambda]$ 
10:   $\gamma \leftarrow \text{hypernet}(\lambda)$ 
11:  Optimize hypernet for loss  $\mathcal{L}_{\text{valid}}(\gamma) + \lambda \mathbf{T}^\top \gamma$ 
12: end while

```

---

## 3 Computational experiment

The experiment contains of two parts. In the first part we get pretrained model, fix  $\lambda$  and then optimize  $\gamma$ . In the second part we train hypernetwork to find  $\gamma$  for arbitrary  $\lambda$ . We will use ResNet18 [1] as base model and dataset Cifar-10 to train and validate.

Main purpose of the first part is to approximate  $\Lambda$ .  $\Lambda$  is the max value for  $\lambda$  which keeps positive metrics. But before all of it we measure execution time of each module. Then for each  $\lambda = \lambda_1, \lambda_2, \dots, \lambda_n$  we optimize  $\gamma$  for loss with time regularization. Finally, we plot accuracy versus  $\lambda$ . It is expected that plot will have plateau for small  $\lambda$  and then will decrease to worst quality for big  $\lambda$ .

## References

- [1] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.

- 46 [2] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax.  
47 *arXiv preprint arXiv:1611.01144*, 2016.

48 *Received January 01, 2017*