

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

S. S. Borodin¹, K. D. Yakovlev¹, and O. Y. Bakhteev^{1,2}

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

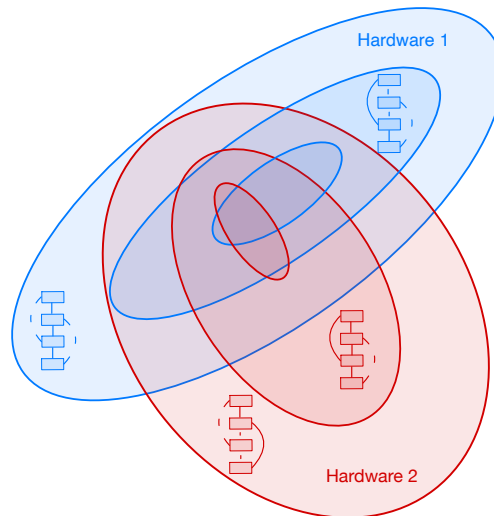
¹ MIPT, Russia

² Dorodnicyn Computing Center RAS, Russia

The paper investigates the problem of structural pruning with respect to target hardware properties. A method considers performance of a target platform to optimize both accuracy and latency on the platform. We use a hypernet to generate a pruned model for a desired trade-off between accuracy and latency. The hypernet is trained end-to-end with backpropagation through the main model. The model adapts 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 CIFAR-10 dataset with ResNet18 as a backbone-model using different hardware (e.g. ...) and compare the resulting architectures with

DOI: 10.21469/22233792

1 Introduction



work in progress...

Since introduction of residual networks [3] models were able to grow in size without vanishing/exploding gradients. As bigger models accomplish better accuracy, they become over-parameterized [7].

One of the techniques to reduce over-parameterization is pruning [2, 5] that can reduce the number of parameters by factors. But we can target other characteristics to reduce. In this

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

paper we target empirically measured latency. There are other time related properties, but such latency better describes real world requirements [6].

(literature review and state-of-the-art) ...

In this work we propose a method of using hypernet [1] to prune main model with respect to performance of target platform. ...

We call a main model the backbone-model, and a model that generates parameters for the backbone-model the hypernet. In this paper we consider a backbone-model as computational graph, where each edge is an atomic operation such as convolution, pooling, activation etc. Parameter of latency regularization referred to as trade-off. Hypernet is a model that for a given trade-off generates subset of edges i.e. pruned model.

Since extremely small latency trade-off provides empty graph with zero latency and random quality, we conduct a base experiment to find such limit for the parameter that it keeps computational graph connected. Once such limit is found we conduct main experiment. In the main experiment the hypernet is trained end-to-end with backpropagation through the main model.

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.

Pruned model is parameterised with $\gamma \in \{0, 1\}^{|E|}$ such that the output of edge e is multiplied with γ_e . This means that we compute the output of edge e iff $\gamma_e = 1$. Since we use hypernetwork to generate γ for regularization parameter $\lambda \in \mathbb{R}$ the generated parameters are denoted as $\gamma(\lambda)$. Thus we formulate optimization problem as:

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

where $\mathcal{L}_{\text{valid}}(\gamma) := \mathcal{L}(\hat{y}_\gamma(\mathbf{X}), Y)$ is a loss function for a pruned network \hat{y}_γ on validation dataset (\mathbf{X}, Y) , $T(e)$ is the time to execute corresponding atomic operation. Since optimization on discrete space is not differentiable we use Gumbel-Softmax approximation [4].

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. γ^* 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 Γ_{valid} .

33 3 Computational experiment

34 (plan) Brief description of basic and main experiments. Describe backbone-model with Gumbel-
35 Softmax relaxation.

36 (old) The experiment contains of two parts. In the first part we get pretrained model, fix λ
37 and then optimize γ . In the second part we train hypernetwork to find γ for arbitrary λ . We
38 use ResNet18 [3] as base model and dataset Cifar-10 to train and validate.

39 3.1 Basic experiment

40 Main purpose of the basic experiment is to approximate Λ . Λ is the max value for λ which
41 keeps better than random performance. It is important to obtain Λ , because we want to train
42 hypernetwork for all λ which does not lead to broken computational graph.

43 Algorithm 1 plots figure 1 which gives rough approximation for Λ . But for appropriate
44 $(\lambda_1, \dots, \lambda_n)$ estimation is acceptable, since hypernetwork fitting algorithm allows such imprec-
45 cision.

46 We expect the accuracy versus latency plot to contain three parts: plateau for small λ since
47 regularization does not impact enough, slope to random performance, random performance for
48 $\lambda > \Lambda$.

49 As you can see on figure 1 our hypothesis has been confirmed. We can see plateau for $\lambda \leq 2$,
50 slope for $2 \leq \lambda \leq 5$ and random performance for $\lambda \geq 5$. So it is safe to assume $\Lambda > 5$ for the
51 main experiment.

52 3.2 Main experiment

53 In the main experiment we want to train hypernetwork to predict γ for each $\lambda \in [0, \Lambda]$. Since
54 discrete search space of γ is finite, it is possible that obtained hypernetwork would result in
55 piecewise-constant function.

56 It is useful for convergence to initialize hypernetwork to predict such γ that keeps com-
57 putational graph not broken. So beforehand we find such vector of logits γ_{init} that provides
58 better than random performance on relaxed backbone-model. To find such vector we generate
59 $\gamma_{\text{init}} \sim \mathcal{N}(\text{shift}, 1)$ for increasing shift till the condition is met.

60 Once γ_{init} is found we optimize hypernetwork end-to-end for arbitrary λ as shown in 2. More
61 specific we assume that $\lambda \sim \mathcal{U}[0, \Lambda]$. We used piecewise-constant function with 10 intervals with

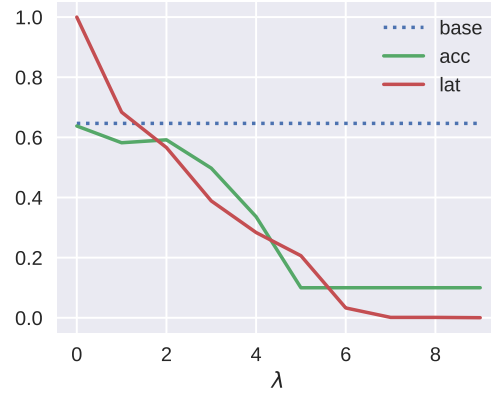


Figure 1 Accuracy and normalized latency plotted versus regularization parameter. The results were obtained on CIFAR-10 dataset with pruned ResNet18.

Algorithm 1 basic algorithm

```

1:  $(\mathbf{X}, Y) \leftarrow$  train dataset
2: load backbone-model
3: for  $\lambda = \lambda_1, \dots, \lambda_n$  do
4:   while not converged do
5:      $\hat{y}_\gamma \leftarrow$  backbone-model.prune( $\gamma$ )
6:      $\text{loss} \leftarrow \mathcal{L}(\hat{y}_\gamma(\mathbf{X}), Y) + \lambda \mathbf{T}^\top \gamma$ 
7:     Optimize  $\gamma$  for loss
8:   end while
9: end for
10: plot figure 1

```

Algorithm 2 hypernetwork train

```

1:  $(\mathbf{X}, Y) \leftarrow$  train dataset
2: load backbone-model
3: initialize hypernetwork with  $\gamma_{\text{init}}$   $\triangleright$  make hypernetwork predict  $\gamma_{\text{init}}$  for all  $\lambda$ 
4: while not converged do
5:   Sample  $\lambda \sim U[0, \Lambda]$ 
6:    $\gamma \leftarrow$  hypernetwork( $\lambda$ )
7:    $\hat{y}_\gamma \leftarrow$  backbone-model.prune( $\gamma$ )
8:    $\text{loss} \leftarrow \mathcal{L}(\hat{y}_\gamma(\mathbf{X}), Y) + \lambda \mathbf{T}^\top \gamma$ 
9:   Optimize hypernetwork for loss
10: end while

```

length = 1. On figure 2 we can see some of intervals have same accuracy and latency. Numerical values for merged intervals are shown in table 1.

3.3 Ablation study/AUC comparison

work in progress...

Compare the method in terms of random pruning. First, sample may random pruned structures. Second, plot the curve speedup/accuracy and compare AUC of the proposed method (w/o hypernetwork) and AUC of random architectures.



Figure 2 Accuracy and normalized latency plotted versus regularization parameter. The results were obtained on CIFAR-10 dataset with pruned ResNet18. This plot shows performance of pruned with hypernetwork backbone-model.

Table 1 Comparison of accuracy and acceleration of pruned with hypernetwork backbone-model for different λ . Here “accuracy diff”-column shows difference between non-pruned model and pruned model accuracy, “acceleration”-column is calculated as $1/\text{latency}$.

λ	accuracy (%)	accuracy diff. (%)	acceleration
base	64	0	1
[0; 2]	63.42	0.58	1.397
[2; 4]	60.33	3.66	1.585
[4; 5]	50.94	13.06	1.798
[5; 8]	52.44	11.56	2.567

To present value of our method we obtain pruned models by our method and method ... (WIP). Then we measure latency and accuracy of each model and plot accuracy versus latency. To conduct numerical comparison we calculate area under the curve (AUC) for both methods. The results are written in table 2

Table 2 Comparison of AUC for different methods

method	AUC
hypernet (ours)	0.342
random	0.065

(For now) As you can see on 3 our method provides trade-off between accuracy and latency of pruned networks, while random provides broken computational graph.

3.4 Ablation study/Comparison of discrete and relaxed model

Since optimization of original task on discrete space is not differentiable, we relax optimization. This approach changes optimization problem. In this experiment we ensure that solution for the relaxed problem is an approximation of a solution for the original problem.

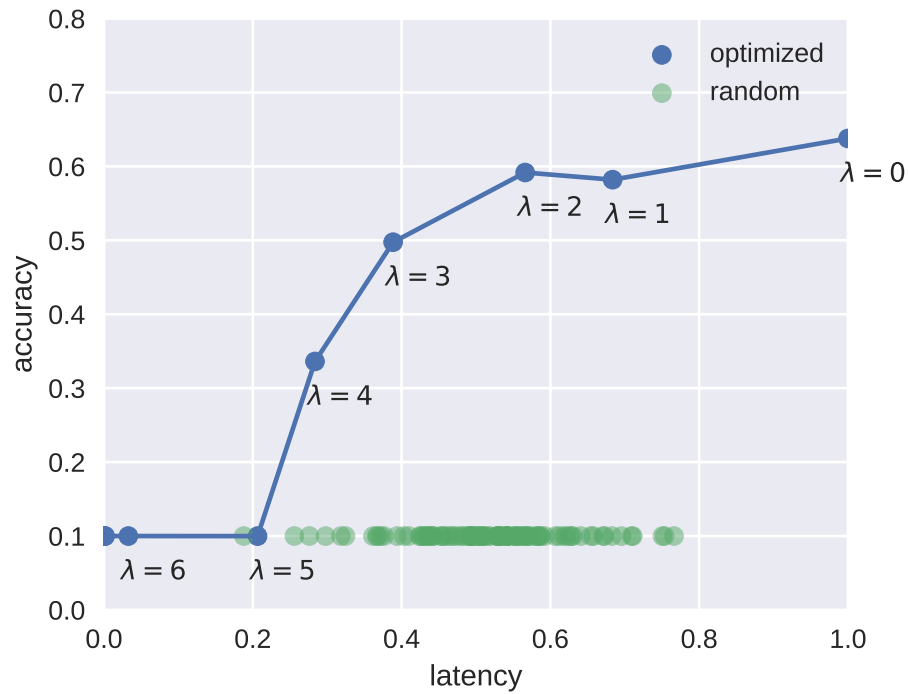


Figure 3 Comparing accuracy-latency dependence of models obtained by our method and random subarchitectures.

Initially we obtain pruned relaxed models for $\lambda_1, \dots, \lambda_n$. For each pruned model we discretize γ , thus obtaining discrete models. Finally we measure accuracy for each model and plot the results: accuracy of relaxed models versus accuracy of discrete models for each λ_k . Complete equivalence between tasks should result in identity curve. As we can see on figure 4 actual curve is close to identity curve, but it does not follow identity curvature.

May be plot multiple runs. May be plot latency in the same way. Convert to pdf.

References

- [1] David Ha, Andrew Dai, and Quoc V. Le. Hypernetworks, 2016.
- [2] Song Han, Jeff Pool, John Tran, and William Dally. Learning both weights and connections for efficient neural network. *Advances in neural information processing systems*, 28, 2015.
- [3] 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.
- [4] Eric Jang, Shixiang Gu, and Ben Poole. Categorical reparameterization with gumbel-softmax. *arXiv preprint arXiv:1611.01144*, 2016.
- [5] Yann LeCun, John Denker, and Sara Solla. Optimal brain damage. *Advances in neural information processing systems*, 2, 1989.
- [6] Mingxing Tan, Bo Chen, Ruoming Pang, Vijay Vasudevan, Mark Sandler, Andrew Howard, and Quoc V Le. Mnasnet: Platform-aware neural architecture search for mobile. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 2820–2828, 2019.
- [7] Chiyuan Zhang, Samy Bengio, Moritz Hardt, Benjamin Recht, and Oriol Vinyals. Understanding deep learning (still) requires rethinking generalization. *Communications of the ACM*, 64(3):107–115, 2021.

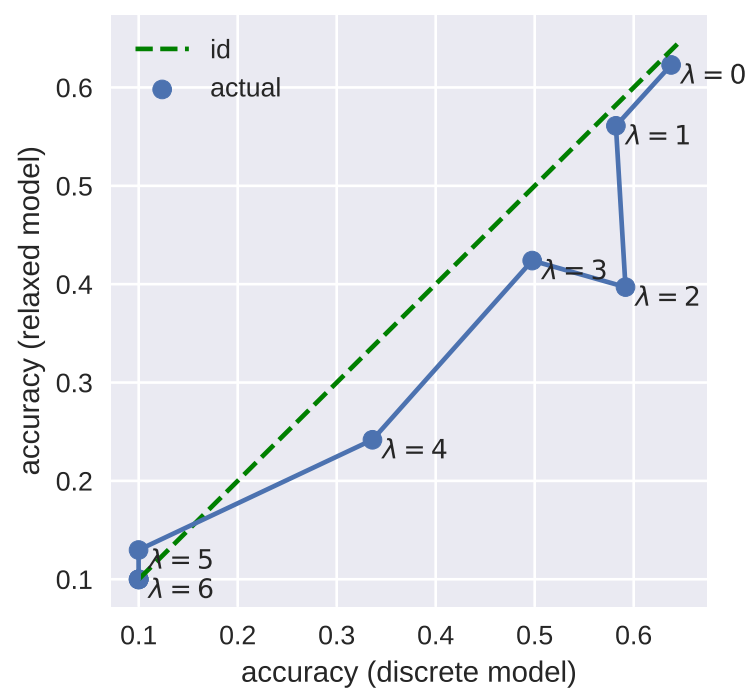


Figure 4 Comparing accuracy of relaxed and discrete models for different λ .

Received January 01, 2017