

BlockDrop: Dynamic Inference Paths in Residual Networks

Zuxuan Wu^{1*}, Tushar Nagarajan^{2*}, Abhishek Kumar³, Steven Rennie⁴
Larry S. Davis¹, Kristen Grauman², Rogerio Feris³
¹ UMD, ² UT Austin, ³ IBM Research, ⁴ Fusemachines Inc.

Abstract

Very deep convolutional neural networks offer excellent recognition results, yet their computational expense limits their impact for many real-world applications. We introduce *BlockDrop*, an approach that learns to dynamically choose which layers of a deep network to execute during inference so as to best reduce total computation without degrading prediction accuracy. Exploiting the robustness of Residual Networks (ResNets) to layer dropping, our framework selects on-the-fly which residual blocks to evaluate for a given novel image. In particular, given a pretrained ResNet, we train a policy network in an associative reinforcement learning setting for the dual reward of utilizing a minimal number of blocks while preserving recognition accuracy. We conduct extensive experiments on CIFAR and ImageNet. The results provide strong quantitative and qualitative evidence that these learned policies not only accelerate inference but also encode meaningful visual information. Built upon a ResNet-101 model, our method achieves a speedup of 20% on average, going as high as 36% for some images, while maintaining the same 76.4% top-1 accuracy on ImageNet.

1. Introduction

Deep neural networks are now ubiquitous in computer vision owing to their recent successes in several important tasks. However, great strides in accuracy have been accompanied by increasingly complex and deep network architectures. This presents a problem for domains where fast inference is essential, particularly in delay-sensitive and real-time scenarios such as autonomous driving, robotic navigation, or user-interactive applications on mobile devices.

Most existing work pursues model compression techniques to speed up a deep network [19, 4, 25, 41, 36, 32, 16, 54, 31]. While significant speed-ups are possible, the approach yields a one-size-fits-all network that requires the same fixed set of features to be extracted for all novel im-

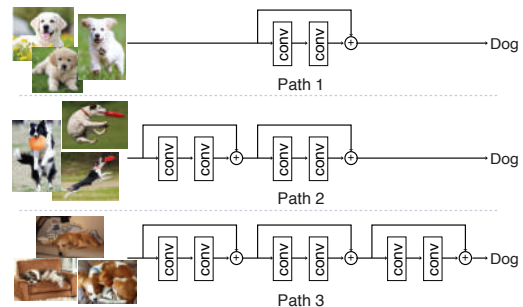


Figure 1: **A conceptual overview of our approach.** Rather than execute all blocks of a ResNet, our approach learns a policy to select the minimal configuration of blocks that is needed to correctly classify a given input image. The resulting instance-specific paths in the network not only reflect the image’s difficulty (easier samples use fewer blocks) but also encode meaningful visual information (patterns of blocks correspond to clusters of visual features).

ages, no matter their complexity. In contrast, an important feature of the human perception system is its ability to adaptively allocate time and scrutiny for visual recognition [49]. For example, a single glimpse is sufficient to recognize some objects and scenes, whereas more time and attention is required to clearly understand occluded or complicated ones [52].

In this spirit, we explore the problem of dynamically allocating computation across a deep network. In particular, we consider Residual Networks (ResNet) [18] both due to their strong track record for recognition tasks [18, 8, 17] as well as their tolerance to removal of layers [50]. ResNets are composed of *residual blocks*, consisting of two or more convolutional layers and skip-connections, which enable direct paths between any two residual blocks. These skip-connections make ResNets behave like ensembles of relatively shallow networks, and hence the removal of a certain residual block generally has only a modest impact on performance [50]. However, the preliminary study of block dropping in ResNets [50] applies a global, manually defined dropping scheme (the same blocks for all images), which leads to increased errors when more blocks are dropped.

We propose to learn optimal block dropping strategies

* Authors contributed equally

that simultaneously preserve both prediction accuracy and minimal block usage based on image-specific decisions. When a novel input is presented to the network trained for recognition, a *dynamic inference path* is followed, selectively choosing which blocks to compute for that instance. See Figure 1. The approach not only improves computational efficiency during inference (*i.e.*, for a similar prediction accuracy, being able to drop more residual blocks than a static global scheme), but also facilitates further insights into ResNets, *e.g.*, whether different blocks encode information about objects, whether the computation needed to classify depends on the difficulty level of the example.

To this end, we introduce *BlockDrop*, a reinforcement learning approach to derive instance-specific inference paths in ResNets. The main idea is to learn a model (referred to as the *policy network*) that, given a novel image, outputs the posterior probabilities of all the binary decisions for dropping or keeping each block in a pretrained ResNet. The policy network is trained using curriculum learning to maximize a reward that incentivizes the use of as few blocks as possible while preserving the prediction accuracy. In addition, the pretrained ResNet is further jointly finetuned with the policy network to produce feature transformations tailored for block dropping behavior. Our approach can be seen as an instantiation of associative reinforcement learning [46] where all the decisions are taken in a single step given the context (*i.e.*, the input instance)¹; this makes policy execution lightweight and scalable to very deep networks.

We conduct extensive experiments on CIFAR [27] and ImageNet [10]. BlockDrop achieves 93.6% and 73.7% accuracy using just 33% and 55% of blocks in a pretrained ResNet-110 on CIFAR-10 and CIFAR-100, respectively, outperforming state-of-the-art methods [14, 15, 12, 32] by clear margins. Furthermore, BlockDrop speeds up a ResNet-101 model on ImageNet by 20% while maintaining the same 76.4% top-1 accuracy². Qualitatively, we observe that the dropping policies learned with BlockDrop are correlated with the visual patterns in the images, *e.g.*, within the “orange” class, images containing a pile of oranges take an inference path that is different from that taken by the close-up images of oranges. Furthermore, BlockDrop policies for *easy* images with clearly visible objects utilize fewer residual blocks compared to the *difficult* images that contain other occluding or background objects. Note that although our analysis in this paper is focused on vanilla ResNets, our approach could also be applied to other recently proposed ResNet variants such as ResNeXt [55] or Multi-Residual Networks [1], as well as other tasks beyond image classification.

¹It can also be seen as contextual bandits [29] although we do not operate in an online setting which has an objective of minimizing the *regret*.

²<https://goo.gl/EwHQcq>

2. Related Work

Layer Dropping in Residual Networks. Dropping layers in residual networks has been used as a regularization mechanism, similar to Dropout [44] or DropConnect [53], for *training* very deep networks (*e.g.*, over 1000 layers) with stochastic depth [22]. Unlike our method, residual layer dropping in stochastic depth networks happens only during the training stage, but at *test time* the layers remain fixed. Veit *et al.* [50] show that ResNets are resilient to layer dropping at test time, which motivates our approach; however, they do not provide a way to dynamically choose which layers could be removed from a network without sacrificing accuracy. More recently, Huang and Wang [23] propose a method for selecting a subset of residual blocks to be executed based on a sparsity constraint. In contrast to these approaches, we propose an *instance-specific* residual block removal scheme to speed up ResNets during inference.

Model Compression. The need to deploy top-performing deep neural network models on mobile devices motivates techniques that can effectively reduce the storage and computational costs of such networks, including knowledge distillation [19, 40, 4], low-rank factorization [25, 47, 41], filter pruning [30, 36, 32, 57], quantization [16, 54, 31], compression with structured matrices [6, 43], network binarization [38, 7, 33], and hashing [5]. Efficient network architectures such as *SqueezeNet* [24] and *MobileNet* [20] have also been explored for training compact deep nets. In contrast to this line of work where the same amount of computation is applied to all images, we focus on efficient inference by dynamically choosing a subset of blocks to be executed *conditioned on the input image*. More importantly, our method is *complementary* to these model compression techniques: the residual blocks that are kept for evaluation can be further pruned for even greater speed up.

Conditional Computation. Several *conditional computation* methods have been proposed to dynamically execute different modules of a network model on a per-example basis [3, 2]. Sparse activations in combination with gating functions are usually adopted to selectively turn on and off a subset of modules based on the input. These gating functions can be learned with reinforcement learning [2, 34, 11]. These models typically associate a reward with a series of decisions computed after each layer/path; the resulting policy execution overhead makes it expensive to scale them up to very deep models with hundreds or thousands of layers. In contrast, our policy network makes all routing decisions in a *single step*, resulting in lower overhead cost for the routing itself and thus larger computational savings. Reinforcement learning has also been applied for dynamic feature prioritization in images [26] and video [45, 56], actively deciding which frames or image regions to visit next. These techniques could be used in tandem with our approach.

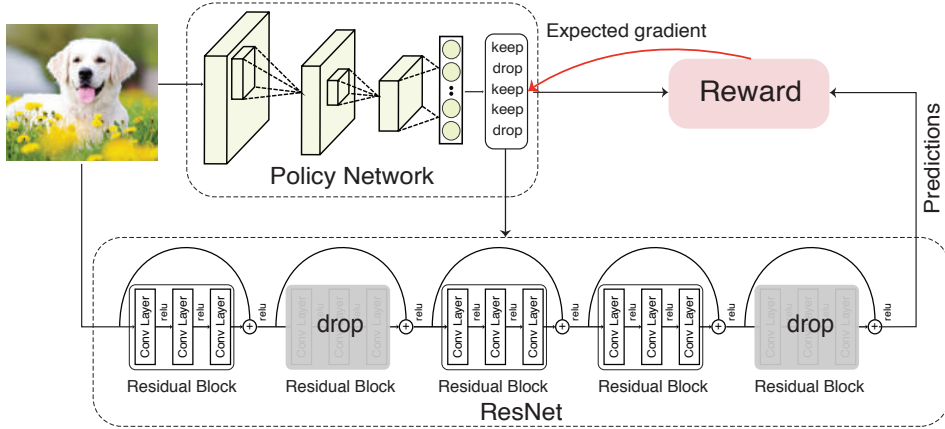


Figure 2: **Illustration of our proposed framework.** Given a new image, the policy network outputs dropping and keeping decisions for each block in a pretrained ResNet, which then makes a prediction by evaluating the active blocks only. Policy rewards account for both block usage and prediction accuracy. The policy network is trained to optimize the expected reward with a curriculum learning strategy, and then jointly finetuned with the ResNet.

Early Prediction. Our work relates more strongly to *early prediction models*, a class of conditional computation models that exit once a criterion (*e.g.*, sufficient confidence for classification) is satisfied at early layers. Cascade detectors [13, 51] are among the earliest methods that exploit this idea in computer vision, often relying on handcrafted control decisions learned separately from visual features. More recently, joint learning of features and early decisions has been studied for deep neural networks. Teerapittayanon *et al.* [48] propose *BranchyNet*, a network composed of branches at each layer to make early classification decisions. Similarly, Adaptive Computation Time (ACT) [15] augments an RNN with a halting unit whose activation determines the probability that computation should continue.

Figurnov *et al.* [14] further extend this idea to the spatial domain in ResNets by applying ACT to each spatial position of multiple image blocks. Like our work, their formulation identifies instance-specific ResNet configurations, but it only allows configurations that use early, contiguous blocks in each predefined segment of the ResNet. These early blocks usually encode low-level features in high-dimensional feature maps, and may lack the discriminative power required for the task. This issue can be mitigated by using images at different scales [35, 21], but at a higher computational cost. Instead, we allow *any* block to contribute to our network, allowing for a much higher variability in potential ResNet configurations and policies.

3. Approach

Given a test image, our goal is to find the best configuration of computational blocks in a pretrained ResNet model, such that a minimum number of blocks is used, without incurring a decrease in classification accuracy. Treating

the task of finding this configuration as a search problem quickly becomes intractable for deeper models as the number of potential configurations grows exponentially with the number of blocks. Learning a soft-attention mask over the blocks also presents problems, namely the difficulty of converting this mask into binary decisions which would require carefully handcrafted thresholds. In addition, such a thresholding operation is non-differentiable, making it non-trivial to directly adopt a supervised learning framework.

We therefore leverage *policy search methods* from reinforcement learning to derive the optimal block dropping schemes that encourage correct predictions with minimal block usage. To this end, we first revisit the architecture of ResNet in Sec. 3.1, and discuss why it is a good fit for block dropping. Then we introduce our policy network in Sec. 3.2, which learns to dynamically select inference paths conditioned on the input image. Finally, we present the training algorithm of our model in Sec. 3.3.

3.1. Pretrained Residual Networks

ResNets consist of multiple stacked *residual blocks* which are essentially regular convolutional layers that are bypassed by identity skip-connections. If we denote the input to the i -th residual block as y_i , and the function represented by its residual block as \mathcal{F}_i , the output of this residual block is given by: $y_{i+1} = \mathcal{F}_i(y_i) + y_i$, which is directly fed as input to the next residual block.

The presence of identity skip-connections induces direct paths between any two residual blocks, and hence top layers in the network are able to access information from bottom layers during a forward pass while gradients can be directly passed from higher layers to lower layers in the back-propagation phase. Veit *et al.* [50] demonstrated that removing (or dropping) a residual block at test time (*i.e.*, having

$y_{i+1} = y_i$) does not lead to a significant accuracy drop. This behavior is due to the fact that ResNets can be viewed as an ensemble of many paths—as opposed to single-path models like AlexNet [28] and VGGNet [42]—and so information can be preserved even with the deletion of paths.

The results in [50] suggest that different blocks *do not share strong dependencies*. However, the study also shows classification errors do increase when more blocks are removed from the model during inference. We contend this is the result of their adopting a global dropping strategy for all images. We posit the best dropping schemes, which lead to correct predictions with the minimal number of blocks, must be instance-specific.

3.2. Policy Network for Dynamic Inference Paths

The *configurations* in the context of ResNets represent decisions to keep/drop each block, where each decision to drop a block corresponds to removing a subset of paths from the network. We refer to these decisions as our *dropping strategy*. To derive the optimal dropping strategy given an input instance, we develop a *policy network* to output a binary *policy vector*, representing the actions to keep or drop a block in a pretrained ResNet. During training, a reward is given considering both block usage and prediction accuracy, which is generated by running the ResNet with only active blocks in the policy vector. See Figure 2 for an overview.

Unlike standard reinforcement learning, we train the policy to predict *all actions at once*. This is essentially a single-step Markov Decision Process (MDP) given the input state and can also be viewed as contextual bandit [29] or associative reinforcement learning [46]. We examine the positive impact of this design choice on scalability in Sec. 4.2.

Formally, given an image \mathbf{x} and a pretrained ResNet with K residual blocks, we define a policy of block-dropping behavior as a K -dimensional Bernoulli distribution:

$$\pi_{\mathbf{W}}(\mathbf{u}|\mathbf{x}) = \prod_{k=1}^K s_k^{\mathbf{u}_k} (1 - s_k)^{1 - \mathbf{u}_k} \quad (1)$$

$$\mathbf{s} = f_{pn}(\mathbf{x}; \mathbf{W}), \quad (2)$$

where f_{pn} denotes the *policy network* parameterized by weights \mathbf{W} and \mathbf{s} is the output of the network after the $\sigma(x) = \frac{1}{1 + e^{-x}}$ function. We choose the architecture of f_{pn} (details below in Sec. 4) such that the cost of running it is negligible compared to ResNet, *i.e.*, so that policy execution overhead remains low. The k -th entry of the vector, $s^k \in [0, 1]$, represents the likelihood of its corresponding residual block in the original ResNet being dropped. An action $\mathbf{u} \in \{0, 1\}^K$ is selected based on \mathbf{s} . Here, $\mathbf{u}^k = 0$ and $\mathbf{u}^k = 1$ indicate dropping and keeping the k -th residual block, respectively.

Only the blocks that are not dropped according to \mathbf{u} will be evaluated in the forward pass. To encourage both correct

predictions as well as minimal block usage, we associate the actions taken with the following reward function:

$$R(\mathbf{u}) = \begin{cases} 1 - (\frac{|\mathbf{u}|_0}{K})^2 & \text{if correct} \\ -\gamma & \text{otherwise.} \end{cases} \quad (3)$$

Here, $(\frac{|\mathbf{u}|_0}{K})^2$ measures the percentage of blocks utilized; when a correct prediction is produced, we incentivize block dropping by giving a larger reward to a policy that uses fewer blocks. We penalize incorrect predictions with γ , which controls the trade-off between efficiency (block usage) and accuracy (*i.e.*, a larger value leads to more correct, but less efficient policies). We use this parameter to vary the *operating point* of our model, allowing different models to be trained depending on the target budget constraint. Finally, to learn the optimal parameters of the policy network, we maximize the following expected reward:

$$J = \mathbb{E}_{\mathbf{u} \sim \pi_{\mathbf{W}}} [R(\mathbf{u})]. \quad (4)$$

In summary, our model works as follows: f_{pn} is used to decide which blocks of the ResNet to keep conditioned on the input image, a prediction is generated by running a forward pass with the ResNet using *only* these blocks, and a reward is observed based on correctness and efficiency.

3.3. Training the BlockDrop Policy

Expected gradient. To maximize Eqn. 4, we utilize policy gradient [46], one of the seminal policy search methods [9], to compute the gradients of J . In contrast to typical reinforcement learning methods where policies are sampled from a multinomial distribution [46], our policies are generated from a K -dimensional Bernoulli distribution. With $\mathbf{u}_k \in \{0, 1\}$, the gradients can be derived similarly as:

$$\begin{aligned} \nabla_{\mathbf{W}} J &= \mathbb{E}[R(\mathbf{u}) \nabla_{\mathbf{W}} \log \pi_{\mathbf{W}}(\mathbf{u}|\mathbf{x})] \\ &= \mathbb{E}[R(\mathbf{u}) \nabla_{\mathbf{W}} \log \prod_{k=1}^K s_k^{\mathbf{u}_k} (1 - s_k)^{1 - \mathbf{u}_k}] \\ &= \mathbb{E}[R(\mathbf{u}) \nabla_{\mathbf{W}} \sum_{k=1}^K \log[s_k^{\mathbf{u}_k} + (1 - s_k)(1 - \mathbf{u}_k)]], \end{aligned} \quad (5)$$

where again \mathbf{W} denotes the parameters of the policy network. We approximate the expected gradient in Eqn. 5 with Monte-Carlo sampling using all samples in a mini-batch. These gradient estimates are unbiased, but exhibit high variance [46]. To reduce variance, we utilize a self-critical baseline $R(\tilde{\mathbf{u}})$ as in [39], and rewrite Eqn. 5 as:

$$\nabla_{\mathbf{W}} J = \mathbb{E}[A \nabla_{\mathbf{W}} \sum_{k=1}^K \log[s_k^{\mathbf{u}_k} + (1 - s_k)(1 - \mathbf{u}_k)]], \quad (6)$$

where $A = R(\mathbf{u}) - R(\tilde{\mathbf{u}})$ and $\tilde{\mathbf{u}}$ is defined as the maximally probable configuration under the current policy, \mathbf{s} : *i.e.*, $\mathbf{u}_i = 1$ if $s_i > 0.5$, and $\mathbf{u}_i = 0$ otherwise [39].

We further encourage exploration by introducing a parameter α to bound the distribution \mathbf{s} and prevent it from saturating, by creating a modified distribution \mathbf{s}' :

$$\mathbf{s}' = \alpha \cdot \mathbf{s} + (1 - \alpha) \cdot (\mathbf{1} - \mathbf{s}).$$

This bounds the distribution in the range $1 - \alpha \leq s'_i \leq \alpha$, from which we then sample the policy vector.

Curriculum learning. Policy gradient methods are typically extremely sensitive to their initialization. Indeed, we found that starting from a randomly initialized policy and optimizing for both accuracy and block usage is not effective, due to the extremely large dimension of the search space, which scales exponentially with the total number of blocks (there are 2^K possible on/off configurations of the blocks). Note that in contrast with applications such as image captioning where ground-truth action sequences (captions) can be used to train an initial policy [39], here no such “expert examples” are available, other than the standard single execution path that executes all blocks.

Therefore, to efficiently search for good action sequences, we take inspiration from the idea of curriculum learning [3]. During epoch t , for $1 \leq t < K$, we keep the first $K - t$ blocks on, and learn a policy only for the last t blocks. As t increases, the activity of more blocks are optimized, until finally all blocks are included (*i.e.*, when $t \geq K$). Using this approach, the activation of each block is first optimized according to unmodified input features in order to assess the utility of the block, and then is gradually exposed to increasingly different feature inputs as t increases and the policy for the last t blocks is jointly trained. This procedure is efficient, and it is effective at identifying and removing blocks that are redundant for the input data instance being considered. It is similar in spirit to [37, 39] that gradually exposes sequences when training with REINFORCE for text generation.

Joint finetuning. After curriculum learning, our policy network is able to identify which residual blocks in the original ResNet to drop for a given input image. Though the policy network is trained to preserve accuracy as much as possible, removing blocks from the pre-trained ResNet will inevitably result in a mismatch between training and testing conditions. We therefore jointly finetune the ResNet with the policy network, so that it can adapt itself to the learned block dropping behavior. The principle of our joint training procedure is similar to that of stochastic depth [22], with the exception that the drop rates are not fixed, but are instead controlled by the policy network. Alg. 1 presents the complete training procedure for our framework.

Algorithm 1 The pseudo-code for training our network.

Input: An input image \mathbf{x} and its label

- 1: Initialize the weights of policy network \mathbf{W} randomly
- 2: Set epochs for curriculum learning and joint finetuning to M^{cl} and M^{ft} , respectively; and set α
- 3: **for** $t \leftarrow 1$ to M^{cl} **do**
- 4: $\mathbf{s} \leftarrow f_{pn}(\mathbf{x}; \mathbf{W})$
- 5: $\mathbf{s} \leftarrow \alpha \cdot \mathbf{s} + (1 - \alpha) \cdot (\mathbf{1} - \mathbf{s})$
- 6: **if** $t < K$ **then**
- 7: set $\mathbf{s}_{[1:K-t]} = \mathbf{1}$ ▷ curriculum training
- 8: **end if**
- 9: $\mathbf{u} \sim \text{Bernoulli}(\mathbf{s})$
- 10: Execute the ResNet according to \mathbf{u}
- 11: Evaluate reward $R(\mathbf{u})$ with Eqn. 3
- 12: Back-propagate gradients computed with Eqn. 6
- 13: **end for**
- 14: **for** $t \leftarrow 1$ to M^{ft} **do**
- 15: Jointly finetune ResNet and policy network
- 16: **end for**

4. Experiment

4.1. Experimental Setup

Datasets and evaluation metrics. We evaluate our method on three benchmarks: CIFAR-10, CIFAR-100 [27], and IMAGENET (ILSVRC2012) [10]. The CIFAR datasets consist of 60,000 32×32 colored images, with 50,000 images for training and 10,000 for testing. They are labeled for 10 and 100 classes for CIFAR-10 and CIFAR-100, respectively. Performance is measured by classification accuracy. ImageNet contains 1.2M training images labeled for 1,000 categories. We test on the validation set of 50,000 images and report top-1 accuracy.

Pretrained ResNet. For CIFAR-10 and CIFAR-100, we experiment with two ResNet models that achieve promising results. In particular, ResNet-32 and ResNet-110 start with a convolutional layer followed by 15 and 54 residual blocks, respectively. These residual blocks, each of which contains two convolutional layers, are evenly distributed into 3 segments with down-sampling layers in between. Finally, a fully-connected layer with 10/100 neurons is applied. See [18] for details. For ImageNet, we adopt ResNet-101 with a total of 33 residual blocks, organized into four segments (*i.e.*, [3, 4, 20, 3]). Here, each residual block contains three convolutional layers based on the *bottleneck* design [18] for computational efficiency. These models are pretrained to match state-of-the-art performance on the corresponding datasets when run without our policy network.

Policy network architecture. For our policy network, we use ResNets with a fraction of the depth of the base model. For CIFAR, we use a ResNet with 3 blocks (equivalently ResNet-8), while for ImageNet, we use a ResNet with 4 blocks (equivalently ResNet-10). In addition, we downsam-

		CIFAR-10				CIFAR-100			
		Acc	K	Acc (<i>ft</i>)	K (<i>ft</i>)	Acc	K	Acc (<i>ft</i>)	K (<i>ft</i>)
ResNet-32	FirstK	16.6	10	84.3	7	23.3	13	66.5	14
	RandomK	20.5	10	88.9	7	38.3	13	67.6	14
	DistributeK	23.4	10	90.2	7	31.9	13	66.7	14
	Ours	88.6	9.4	91.3	6.9	58.3	12.4	68.7	13.1
	Full ResNet	92.3	15	92.3	15	69.3	15	69.3	15
ResNet-110	FirstK	13.3	21	71.3	17	63.5	50	57.9	31
	RandomK	14.5	21	90.1	17	66.3	50	68.4	31
	DistributeK	13.0	21	92.7	17	49.6	50	69.9	31
	Ours	75.4	20.1	93.6	16.9	72.1	49.1	73.7	30.2
	Full ResNet	93.2	54	93.2	54	72.2	54	72.2	54

Table 1: **Accuracy and block usage with our policies vs. heuristic baselines**, with and without jointly finetuning (*ft*) for all methods. For fair comparisons, K is selected based on the average block usage of our method, and this can be different before and after finetuning. Note that the average value of K for our method is reported here for brevity. It is determined dynamically per image, and can be as low as 3 (out of 54) in ResNet-110 on CIFAR-10.

ple images to 112×112 as the input of the policy network for ImageNet experiments. The computation required for the policy network is 4.8% and 3.0% of the total ResNet computation for the CIFAR (ResNet-110) and ImageNet (ResNet-101) models respectively, making policy computations negligible (it takes about 0.5 ms per image on average for ImageNet). While a recurrent model (*e.g.*, LSTM) could also serve as the policy network, we found a CNN to be more efficient with similar performance.

Implementations details. We adopt PyTorch for implementation and utilize ADAM as the optimizer. We set α to 0.8, learning rate to $1e-4$, and use a batch size of 2048 during curriculum learning. For joint finetuning, we adjust the batch size to 256 and 320 on CIFAR and ImageNet, respectively, and adjust the learning rate to $1e-5$ for ImageNet. Our code is available at <https://goo.gl/NqyNeN>.

4.2. Quantitative Results

Learned policies vs. heuristics. We compare our block dropping strategy to the following alternative methods: (1) FIRSTK, which keeps only the first K residual blocks active; (2) RANDOMK, which keeps K randomly selected residual blocks active; (3) DISTRIBUTEK, which evenly distributes K blocks across all segments. For all baselines, we choose K to match the average number of blocks used by BlockDrop, rounding up as needed. DistributeK allows us to see if feature combinations of different blocks learned by BlockDrop are better than features learned from the restricted set of early blocks of each segment. This setting resembles the allowable feature combinations from early stopping models applied to ResNets.

The results in Table 1 highlight the advantage of our instance-specific policy. On CIFAR-10, the learned policies give an accuracy of 88.6% and 75.4% using an average of 9.4 and 20.1 blocks from the original ResNet-32

and ResNet-110 respectively, outperforming the baselines by a large margin. Furthermore, the instance-specific nature of our method allows us to capture the inherent variance in the computational requirements of our dataset. We notice a wide distribution in block usage depending on the image. With ResNet-110, nearly 15% of the images use fewer than 10 blocks, with some images using as few as 3 blocks. This variance cannot be captured by any static policies. Similar trends are observed on CIFAR-100. This confirms that dropping residual blocks with policies computed in a learned manner is indeed significantly better than heuristic dropping behaviors. The fact that RandomK performs better than FirstK is interesting, suggesting the value of having residual blocks at different segments to learn feature representations at different scales.

Impact of joint finetuning. Next we analyze the impact of joint finetuning (*cf.* Sec. 3.3) for both our approach and the baselines, denoted *ft* in Table 1.

Joint finetuning further significantly improves classification accuracy using fewer (or almost the same) number of blocks. In particular, on CIFAR-10, it offers absolute performance gains of 2.7% and 18.2% using 2.5 and 3.2 fewer blocks with ResNet-32 and ResNet-110 respectively compared with curriculum training alone. Similarly, on CIFAR-100, joint finetuning improves accuracies and brings down block usage with ResNet-110. For ResNet-32, we observe 0.7 more blocks on average are used after finetuning, which might be due to the challenging nature of CIFAR-100 requiring more blocks to make correct predictions. Comparing ResNet-110 with ResNet-32, we observe that the computational speed-ups are more dramatic for deeper ResNets owing to the fact that there are more blocks with potentially diverse features to select from. When built upon ResNet-110, our method outperforms the pretrained model by 0.4% and 1.5% (absolute) using 31% and 55.9% of the original

blocks on CIFAR-10 and CIFAR-100, respectively. Additionally, we observe that some images use as few as 5 blocks for inference. These results confirm that joint finetuning can indeed assist the ResNet to adapt to the removal of blocks by refining its feature representations while maintaining its capacity for instance-specific variation.

BlockDrop vs. state-of-the-art methods. We next compare BlockDrop to several techniques from the literature. We vary γ , which controls our algorithm’s trade-off between block usage and accuracy, to get a range of models with varying computational requirements. We compute the average FLOPs utilized to classify each image in the test set; FLOPs are a hardware independent metric, allowing for fair comparisons across models.³

We compare to the following state-of-the-art methods⁴: (1) ACT and (2) SACT [14], (3) PFEC [32], (4) LCCL [12]. ACT and SACT learn a halting score at the end of each block, and exit the model when a high-confidence is obtained. PFEC and LCCL reduce the parameters of convolutional layers by either pruning or sparsity constraints, which is complementary to our method. Other model compression methods cited earlier do not report results on larger ResNet models, and hence are not available to compare here.

Figure 3 (a) presents the results on CIFAR. We observe that our best model offers 0.4% performance gain in accuracy (93.6% vs. 93.2%) using 65% fewer FLOPs on average (1.73×10^8 vs. 5.08×10^8) over the original ResNet-110 model. The performance gains might result from the regularization effect of dropping blocks when finetuning the network as in [22]. Compared to ACT and SACT, our method only requires 50% of the FLOPs to achieve the same level of precision (>93.0%). BlockDrop also exhibits a much higher variance in its FLOPs over other methods. Compared to SACT, this variance is 3 times larger, allowing some samples to achieve a speedup as high as 85% with correct predictions. Further, BlockDrop also outperforms PFEC [32] and LCCL [12], which are complementary compression techniques and can be utilized together with our framework to speed up convolution operations.

Figure 3 (b) presents the results for ImageNet. Compared with the original ResNet-101 model, BlockDrop again achieves slightly better results (76.8% vs. 76.4%) with 6% speed up (1.47×10^{10} vs. 1.56×10^{10} FLOPs). BlockDrop performs on par with the full ResNet with a 20% speed up (1.25×10^{10} vs. 1.56×10^{10} FLOPs) when we relax γ slightly. This 20% acceleration without degradation in accuracy is quite promising. For example, in a high-precision

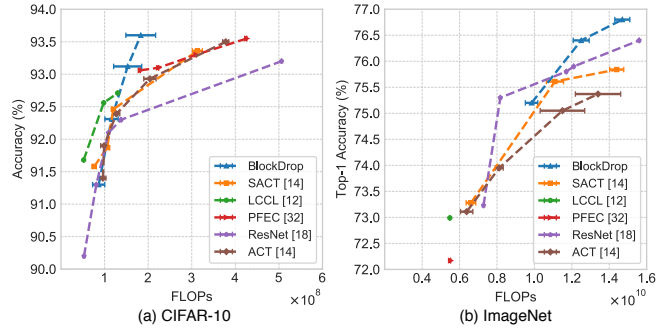


Figure 3: **FLOPs vs. accuracy on CIFAR-10 and ImageNet.** Results compared to several state-of-the-art methods. Error bars denote the standard deviation across images.

		Time (ms)	Speed-up
ResNet-32	Full ResNet	7.71	–
	Ours-single	6.56	14.9%
	Ours-seq	9.92	-28.7%
ResNet-110	Full ResNet	24.1	–
	Ours-single	10.9	52.3%
	Ours-seq	29.1	-20.7%

Table 2: **Impact of our single-step policy inference on efficiency for CIFAR-10.** See text for details.

image recognition service accepting 1 billion daily API calls, such a speedup would save around 1000 hours of computation on a single P6000 GPU (0.024 seconds/image).

Efficiency advantage of single-step policy. The single-step design of our policy network—where the full dynamic inference path is computed without revisiting intermediate outputs of the network—has important efficiency advantages. In short, it permits lower policy execution overhead. To examine the impact empirically, we devised a variant of BlockDrop that uses traditional RL policy learning to instead make sequential decisions (see Supp. for details). We select models of both variants that attain equivalent accuracy, with the same number of blocks. To ensure fair comparison, we run all three models on the same single NVIDIA P6000 GPU while disabling other processes.

Table 2 shows the results for CIFAR-10. We report the time per test image and the speed-up over the original ResNet run in entirety with no block dropping. This result confirms the efficiency advantage of our single-step design: to reach the same accuracy, we need much less overhead (e.g., less than 60% of the time required by the sequential variant). In fact, the sequential variant takes even *longer* to run than the original full ResNet models, yielding a negative speed-up. These results reaffirm our choice to compute all actions in one shot rather than compute them sequentially. They also stress the importance of accounting for any overhead a deep net speed-up scheme incurs to make its speed-up decisions.

³Note that we consider the multiply-accumulate operation as a two step process yielding two floating point operations and we only compute FLOPs for convolutional layers and linear layers as they account for most of the computation for inference.

⁴For ACT and SACT on CIFAR, we train models with the authors’ code. For the rest, we compare to numbers in the respective papers.



Figure 4: **Policies learned for four ImageNet classes, *volcano*, *orange*, *hamster* and *castle*.** These policies correspond to a set of active paths in the ResNet, which seem to cater to different “states” of images of the particular class. For *volcano*, these include features like smoke, lava, *etc.*, while for *orange* they include whether it is sliced/whole, quantity.

4.3. Qualitative Results

Finally, we provide qualitative results based on our learned policies. We investigate the visual patterns encoded in these learned policies and then analyze the relation between block usage and instance difficulty.

Visual patterns in policies. Intuitively, related images can be recognized by their similar characteristics (*e.g.*, low-level clues like texture and color). Here, we analyze similarity in terms of the *policies they utilize* by sampling dominant policies for each class and visualizing samples from them. Figure 4 shows samples utilizing three different policies for four classes. It can be clearly seen that images under the same policy are similar, and different policies encode different styles, although they all correspond to the same semantic concept. For example, the first inference path for the “orange” class caters to images containing a pile of oranges, and close up views of oranges activate the second inference path, while images containing slices of oranges are routed through the third inference path. These results indicate that different paths encode meaningful semantic visual patterns, based on the input images. While this happens in standard ResNets as well, all images necessarily utilize all the paths, and disentangling this information is not possible.

Instance difficulty. Instance difficulty is well understood in the context of prediction confidence, where easy and difficult examples are classified with high and low probabilities, respectively. Inspired by the above analysis that revealed interesting correlations between the inference policies and the visual patterns in the images, we try to characterize instance difficulty in terms of block usage. We hypothesize that simple examples (*e.g.* images with clear objects, without occlusions) require fewer computations to be correctly recognized. To qualitatively analyze the correlations between instance difficulty and block usage, we utilize learned policies that lead to high-confidence predictions for each class.

Figure 5 illustrates samples from ImageNet. The top row contains images that are correctly classified with the least number of blocks, while samples in the bottom row utilize the most blocks. We see that samples using fewer blocks are indeed easier to identify since they contain single frontal-view objects positioned in the center, while several objects,



Figure 5: **Samples from ImageNet classes.** Easy and hard samples from *goldfish*, *artichoke*, *spacecraft* and *bridge* to illustrate how block usage translates to instance difficulty.

occlusion, or cluttered background occur in samples that require more blocks. This confirms our hypothesis that block usage is a function of instance difficulty. We stress that this “sorting” into easy or hard cases falls out automatically; it is learned by BlockDrop.

5. Conclusion

We presented BlockDrop, an approach for faster inference in ResNets by selectively choosing residual blocks to evaluate in a learned and optimized manner conditioned on inputs. In particular, we trained a policy network to predict blocks to drop in a pretrained ResNet while trying to retain the prediction accuracy. The ResNet is further jointly finetuned to produce smooth feature representations tailored for block dropping behavior. We conducted extensive experiments on CIFAR and ImageNet, observing considerable gains over existing methods in terms of the efficiency-accuracy trade-off. Further, we also observe that the policies learned encode semantic information in the images.

Acknowledgments: Rogerio Feris is supported by IARPA via DOI/IBC contract number D17PC00341. The U.S. Government is authorized to reproduce and distribute reprints for Governmental purposes notwithstanding any copyright annotation thereon. Disclaimer: The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies or endorsements, either expressed or implied, of IARPA, DOI/IBC, or the U.S. Government. Kristen Grauman is supported in part by an IBM Faculty Award and IBM Open Collaboration Award.

References

- [1] M. Abdi and S. Nahavandi. Multi-residual networks. *arXiv preprint arXiv:1609.05672*, 2016. 2
- [2] E. Bengio, P.-L. Bacon, J. Pineau, and D. Precup. Conditional computation in neural networks for faster models. In *ICML Workshop on Abstraction in Reinforcement Learning*, 2016. 2
- [3] Y. Bengio. Deep learning of representations: Looking forward. In *SLCP*, 2013. 2, 5
- [4] G. Chen, M. Chandraker, T. Han, W. Choi, and X. Yu. Learning efficient object detection models with knowledge distillation. In *NIPS*, 2017. 1, 2
- [5] W. Chen, J. Wilson, S. Tyree, K. Weinberger, and Y. Chen. Compressing neural networks with the hashing trick. In *ICML*, 2015. 2
- [6] Y. Cheng, F. Yu, R. Feris, S. Kumar, A. Choudhary, and S. F. Chang. An exploration of parameter redundancy in deep networks with circulant projections. In *ICCV*, 2015. 2
- [7] M. Courbariaux, I. Hubara, D. Soudry, R. El-Yaniv, and Y. Bengio. Binarized neural networks: Training deep neural networks with weights and activations constrained to ± 1 or ± 1 . *arXiv preprint arXiv:1602.02830*, 2016. 2
- [8] J. Dai, Y. Li, K. He, and J. Sun. R-fcn: Object detection via region-based fully convolutional networks. In *NIPS*, 2016. 1
- [9] M. P. Deisenroth, G. Neumann, J. Peters, et al. A survey on policy search for robotics. *Foundations and Trends® in Robotics*, 2013. 4
- [10] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei. Imagenet: A large-scale hierarchical image database. In *CVPR*, 2009. 2, 5
- [11] L. Denoyer and P. Gallinari. Deep sequential neural network. *arXiv preprint arXiv:1410.0510*, 2014. 2
- [12] X. Dong, J. Huang, Y. Yang, and S. Yan. More is less: A more complicated network with less inference complexity. In *CVPR*, 2017. 2, 7
- [13] P. Felzenszwalb, R. Girshick, and D. McAllester. Cascade object detection with deformable part models. In *CVPR*, 2010. 3
- [14] M. Figurnov, M. D. Collins, Y. Zhu, L. Zhang, J. Huang, D. Vetrov, and R. Salakhutdinov. Spatially adaptive computation time for residual networks. In *CVPR*, 2017. 2, 3, 7
- [15] A. Graves. Adaptive computation time for recurrent neural networks. *arXiv preprint arXiv:1603.08983*, 2016. 2, 3
- [16] S. Han, H. Mao, and W. J. Dally. Deep compression: Compressing deep neural network with pruning, trained quantization and Huffman coding. In *ICLR*, 2016. 1, 2
- [17] K. He, G. Gkioxari, P. Dollár, and R. Girshick. Mask r-cnn. In *ICCV*, 2017. 1
- [18] K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In *CVPR*, 2016. 1, 5
- [19] G. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. In *arXiv preprint arXiv:1503.02531*, 2015. 1, 2
- [20] A. G. Howard, M. Zhu, B. Chen, D. Kalenichenko, W. Wang, T. Weyand, M. Andreetto, and H. Adam. Mobilenets: Efficient convolutional neural networks for mobile vision applications. In *CVPR*, 2017. 2
- [21] G. Huang, D. Chen, T. Li, F. Wu, L. van der Maaten, and K. Q. Weinberger. Multi-scale dense convolutional networks for efficient prediction. *arXiv preprint arXiv:1703.09844*, 2017. 3
- [22] G. Huang, Y. Sun, Z. Liu, D. Sedra, and K. Q. Weinberger. Deep networks with stochastic depth. In *ECCV*, 2016. 2, 5, 7
- [23] Z. Huang and N. Wang. Data-driven sparse structure selection for deep neural networks. *arXiv preprint arXiv:1707.01213*, 2017. 2
- [24] F. N. Iandola, S. Han, M. W. Moskewicz, K. Ashraf, W. J. Dally, and K. Keutzer. Squeezenet: Alexnet-level accuracy with 50x fewer parameters and <0.5 mb model size. In *arXiv:1602.07360*, 2016. 2
- [25] Y. Ioannou, D. Robertson, J. Shotton, R. Cipolla, and A. Criminisi. Training cnns with low-rank filters for efficient image classification. In *ICLR*, 2016. 1, 2
- [26] S. Karayev, M. Fritz, and T. Darrell. Anytime recognition of objects and scenes. In *CVPR*, 2014. 2
- [27] A. Krizhevsky and G. Hinton. Learning multiple layers of features from tiny images. 2009. 2, 5
- [28] A. Krizhevsky, I. Sutskever, and G. E. Hinton. Imagenet classification with deep convolutional neural networks. In *NIPS*, 2012. 4
- [29] J. Langford and T. Zhang. The epoch-greedy algorithm for multi-armed bandits with side information. In *NIPS*, 2008. 2, 4
- [30] Y. LeCun, J. Denker, and S. Solla. Optimal brain damage. In *NIPS*, 1989. 2
- [31] H. Li, S. De, Z. Xu, C. Studer, H. Samet, and T. Goldstein. Training quantized nets: A deeper understanding. In *NIPS*, 2017. 1, 2
- [32] H. Li, A. Kadav, I. Durdanovic, H. Samet, and H. P. Graf. Pruning filters for efficient convnets. In *ICLR*, 2017. 1, 2, 7
- [33] Z. Li, X. Wang, X. Lv, and T. Yang. Sep-nets: Small and effective pattern networks. *arXiv preprint arXiv:1706.03912*, 2017. 2
- [34] L. Liu and J. Deng. Dynamic deep neural networks: Optimizing accuracy-efficiency trade-offs by selective execution. *arXiv preprint arXiv:1701.00299*, 2017. 2
- [35] M. McGill and P. Perona. Deciding how to decide: Dynamic routing in artificial neural networks. In *ICML*, 2017. 3
- [36] A. Polyak and L. Wolf. Channel-level acceleration of deep face representations. *IEEE Access*, 3:2163–2175, 2015. 1, 2
- [37] M. Ranzato, S. Chopra, M. Auli, and W. Zaremba. Sequence level training with recurrent neural networks. In *ICLR*, 2016. 5
- [38] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi. Xnet: Imagenet classification using binary convolutional neural networks. In *ECCV*, 2016. 2
- [39] S. J. Rennie, E. Marcheret, Y. Mroueh, J. Ross, and V. Goel. Self-critical sequence training for image captioning. In *CVPR*, 2017. 4, 5
- [40] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio. Fitnets: Hints for thin deep nets. In *arXiv preprint arXiv:1412.6550*, 2014. 2

- [41] T. N. Sainath, B. Kingsbury, V. Sindhvani, E. Arisoy, and B. Ramabhadran. Low-rank matrix factorization for deep neural network training with high-dimensional output targets. In *ICASSP*, 2013. 1, 2
- [42] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 4
- [43] V. Sindhvani, T. Sainath, and S. Kumar. Structured transforms for small-footprint deep learning. In *NIPS*, 2015. 2
- [44] N. Srivastava, G. E. Hinton, A. Krizhevsky, I. Sutskever, and R. Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *JMLR*, 2014. 2
- [45] Y.-C. Su and K. Grauman. Leaving some stones unturned: dynamic feature prioritization for activity detection in streaming video. In *ECCV*, 2016. 2
- [46] R. S. Sutton and A. G. Barto. *Reinforcement learning: An introduction*. MIT press Cambridge, 1998. 2, 4
- [47] C. Tai, T. Xiao, X. Wang, et al. Convolutional neural networks with low-rank regularization. In *ICLR*, 2016. 2
- [48] S. Teerapittayanon, B. McDanel, and H. Kung. Branchynet: Fast inference via early exiting from deep neural networks. In *ICPR*, 2016. 3
- [49] A. Ude. Integrating visual perception and manipulation for autonomous learning of object representations. *Can developmental robotics yield human-like cognitive abilities?*, 2012. 1
- [50] A. Veit, M. J. Wilber, and S. Belongie. Residual networks behave like ensembles of relatively shallow networks. In *NIPS*, 2016. 1, 2, 3, 4
- [51] P. Viola and M. J. Jones. Robust real-time face detection. *IJCV*, 2004. 3
- [52] D. B. Walther, B. Chai, E. Caddigan, D. M. Beck, and L. Fei-Fei. Simple line drawings suffice for functional mri decoding of natural scene categories. *PNAS*, 2011. 1
- [53] L. Wan, M. Zeiler, S. Zhang, Y. L. Cun, and R. Fergus. Regularization of neural networks using dropconnect. In *ICML*, 2013. 2
- [54] J. Wu, C. Leng, Y. Wang, Q. Hu, and J. Cheng. Quantized convolutional neural networks for mobile devices. In *CVPR*, 2016. 1, 2
- [55] S. Xie, R. Girshick, P. Dollár, Z. Tu, and K. He. Aggregated residual transformations for deep neural networks. In *CVPR*, 2017. 2
- [56] S. Yeung, O. Russakovsky, G. Mori, and L. Fei-Fei. End-to-end learning of action detection from frame glimpses in videos. In *CVPR*, 2016. 2
- [57] R. Yu, A. Li, C.-F. Chen, J.-H. Lai, V. I. Morariu, X. Han, M. Gao, C.-Y. Lin, and L. S. Davis. Nisp: Pruning networks using neuron importance score propagation. In *CVPR*, 2018. 2