

Generalizing Pooling Functions in CNNs: Mixed, Gated, and Tree

Chen-Yu Lee, Patrick Gallagher, and Zhuowen Tu

Abstract—In this paper, we seek to improve deep neural networks by generalizing the pooling operations that play a central role in the current architectures. We pursue a careful exploration of approaches to allow pooling to learn and to adapt to complex and variable patterns. The two primary directions lie in: (1) learning a pooling function via (two strategies of) combining of max and average pooling, and (2) learning a pooling function in the form of a tree-structured fusion of pooling filters that are themselves learned. In our experiments every generalized pooling operation we explore improves performance when used in place of average or max pooling. We experimentally demonstrate that the proposed pooling operations provide a boost in invariance properties relative to conventional pooling and set the state of the art on several widely adopted benchmark datasets. These benefits come with only a light increase in computational overhead during training (ranging from additional 5% to 15% in time complexity) and a very modest increase in the number of model parameters (e.g. additional 1, 9, and 27 parameters for mixed, gated, and 2-level tree pooling operators, respectively). To gain more insights about our proposed pooling methods, we also visualize the learned pooling masks and the embeddings of the internal feature responses for different pooling operations. Our proposed pooling operations are easy to implement and can be applied within various deep neural network architectures.

Index Terms—Convolutional Neural Networks, Deep Learning, Pooling Functions, Supervised Classification.



1 INTRODUCTION

The recent resurgence of neurally-inspired systems such as deep belief nets (DBN) [10], convolutional neural networks (CNNs) [20], and the sum-and-max infrastructure [36] has derived significant benefit from building more sophisticated network structures [37], [43] and from bringing learning to non-linear activations [6], [27]. The pooling operation has also played a central role, contributing to invariance to data variation and perturbation. However, pooling operations have been little revised beyond the current primary options of average, max, and stochastic pooling [3], [46]; this despite indications that e.g. choosing from more than just one type of pooling operation can benefit performance [35].

In this paper, we desire to bring learning and “responsiveness” (i.e., the characteristics of the region being pooled) into the pooling operation. Various approaches are possible, but here we pursue two in particular. In the first approach, we consider combining typical pooling operations (specifically, max pooling and average pooling); within this approach we further investigate two strategies by which to combine these operations. One of the strategies is “unresponsive”; for reasons discussed later, we call this strategy *mixed max-average pooling*. The

other strategy is “responsive”; we call this strategy *gated max-average pooling*, where the ability to be responsive is provided by a “gate” in analogy to the usage of gates elsewhere in deep learning.

Another natural generalization of pooling operations is to allow the pooling operations that are being combined to themselves be learned. Hence in the second approach, we learn to combine pooling filters that are themselves learned. Specifically, the learning is performed within a binary tree (with number of levels that is pre-specified rather than “grown” as in traditional decision trees) in which each leaf is associated with a learned pooling filter. As we consider internal nodes of the tree, each parent node is associated with an output value that is the mixture of the child node output values, until we finally reach the root node. The root node corresponds to the overall output produced by the tree. We refer to this strategy as *tree pooling*. Tree pooling is intended (1) to learn pooling filters directly from the data; (2) to learn how to combine leaf node pooling filters in a differentiable fashion; (3) to bring together these other characteristics within a hierarchical tree structure.

When the mixing of the node outputs is allowed to be “responsive”, the resulting tree pooling operation becomes an integrated method for learning pooling filters and combinations of those filters that are able to display a range of different behaviors depending on the characteristics of the region being pooled. In our experiments, we will see evidence that tree pooling is particularly useful at the lower network layers where the feature responses are denser whereas mixed max-average and gated max-average pooling are more advantageous at the higher, sparser network layers. We refer

-
- C.-Y. Lee was with the Department of Electrical and Computer Engineering, University of California, San Diego, La Jolla, CA 92093, U.S.A.
E-mail: chl260@ucsd.edu
 - P. Gallagher was with the Department of Cognitive Science, University of California, San Diego, La Jolla, CA 92093, U.S.A.
E-mail: patrick.w.gallagher@gmail.com
 - Z. Tu is with the Department of Cognitive Science, University of California, San Diego, La Jolla, CA 92093, U.S.A.
E-mail: ztu@ucsd.edu

to a combined strategy adopting max-average pooling at higher layers and tree pooling at lower layers as tree + max-average pooling.

We make the following observations: First, in an experiment to test invariance (shown in Figure 4) the specific configurations of our proposed gated max-average pooling and tree pooling investigated display better invariance properties than conventional pooling operations across a wide range of transformation types and amounts. Second, our results (detailed in Tables 1, 2, 4, and 5) are obtained using only a modest number of additional parameters — for example, using only 45 additional parameters, we improve AlexNet performance on ImageNet by a 6% relative increase (top-5, single-view). Third, we find that the tree + max-average pooling configuration gives the best overall range of performance (Table 1); we interpret this to indeed indicate that the initial tree pooling layer here is well-suited to the denser low-layer feature while the following gated max-average pooling layer(s) is/are better suited to the sparser high-layer feature maps.

We pursue experimental validation and find that in the architectures we investigate, replacing standard pooling operations with any of our proposed generalized pooling methods boosts performance on each of the standard benchmark datasets, as well as on the larger and more complex ImageNet dataset [34]. We attain state-of-the-art results on MNIST, CIFAR10 (with and without data augmentation), and SVHN. Our proposed pooling operations can be used as drop-in replacements for standard pooling operations in various current architectures and can be used in tandem with other performance-boosting approaches such as learning activation functions, training with data augmentation, or modifying other aspects of network architecture — we confirm improvements when used in a deeply-supervised nets (DSN) style architecture, as well as in AlexNet and GoogLeNet. Our proposed pooling operations are also simple to implement, computationally undemanding (ranging from 5% to 15% additional overhead in timing experiments), differentiable, and use only a modest number of additional parameters. Since our proposed methods serve as drop-in replacements for standard pooling operations, they can be used in tandem with other performance boosting techniques or within alternative architectures [19], [26], [37], [43], [45].

2 RELATED WORK

In the current deep learning literature, popular pooling functions include max, average, and stochastic pooling [2], [3], [46]. A recent effort using more complex pooling operations, spatial pyramid pooling [9], is mainly designed to deal with images of varying size, rather than delving into different pooling functions or incorporating learning. Learning pooling functions is analogous to receptive field learning [5], [8], [11], [16]. However methods like [16] lead to a more difficult learning procedure

that in turn leads to a less competitive result, e.g. an error rate of 16.89% on unaugmented CIFAR10.

Since our tree pooling approach involves a tree structure in its learning, we observe an analogy to “logic-type” approaches such as decision trees [31] or “logical operators” [28]. Such approaches have played a central role in artificial intelligence for applications that require “discrete” reasoning, and are often intuitively appealing. Unfortunately, despite the appeal of such logic-type approaches, there is a disconnect between the functioning of decision trees and the functioning of CNNs — the output of a standard decision tree is non-continuous with respect to its input (and thus nondifferentiable). This means that a standard decision tree is not able to be used in CNNs, whose learning process is performed by back propagation using gradients of differentiable functions. Part of what allows us to pursue our approaches is that we ensure the resulting pooling operation is differentiable and thus usable within network backpropagation.

A recent work, referred to as auto-encoder trees [14], also pays attention to a differentiable use of tree structures in deep learning but is distinct from our method as it focuses on learning encoding and decoding methods (rather than pooling methods) using a “soft” decision tree for a generative model. In the supervised setting, [4] incorporates multilayer perceptrons within decision trees, but simply uses trained perceptrons as splitting nodes in a decision forest; not only does this result in training processes that are separate (and thus more difficult to train than an integrated training process), this training process does not involve the learning of any pooling filters. The work in [17] investigates routing decisions by using a sigmoid function in fully connected layers and achieves a performance boost with 30 trees. Techniques of learning decision functions are also presented in [12], [41].

Since we explore pooling operations in which we both learn pooling filters and also learn how to combine those filters, and since the filter combinations can differ based on the characteristics of the region being pooled, our proposed methods begin with the ability to present a much richer range of responses than conventional pooling methods both during learning and during actual forward operation. From this initial potential, we move to investigate whether the promise of these methods is reflected in actual test performance. Our focus on empirical verification comes both because theory (with a few notable exceptions) has yet to make significant inroads in deep learning, and because the primary recent theoretical work on pooling operations [3] is grounded in assumptions that are not applicable in our setting.

Evaluation of the proposed method after publication: After the acceptance of the conference version of our work [23], the proposed generalized pooling operations have been validated in [29] where the effectiveness of the proposed method has been illustrated in a comprehensive study on the ImageNet benchmark [34].

3 GENERALIZING POOLING OPERATIONS

A typical convolutional neural network is structured as a series of convolutional layers and pooling layers. Each convolutional layer is intended to produce representations (in the form of activation values) that reflect aspects of local spatial structures and to consider multiple channels when doing so. More specifically, a convolution layer computes “feature response maps” that involve multiple channels within some localized spatial region. On the other hand, a pooling layer is restricted to act within just one channel at a time, “condensing” the activation values in each spatially-local region in the currently considered channel. An early reference related to pooling operations (although not explicitly using the term “pooling”) can be found in [11]. In modern visual recognition systems, pooling operations play a role in producing “downstream” representations that are more robust to the effects of variations in data while still preserving important motifs. The specific choices of average pooling [20], [21] and max pooling [32] have been widely used in many CNN-like architectures; [3] includes a theoretical analysis (albeit one based on assumptions that do not hold here).

Our goal is to bring learning and “responsiveness” into the pooling operation. We focus on two approaches in particular. In the first approach, we begin with the (conventional, non-learned) pooling operations of max pooling and average pooling and learn to combine them. Within this approach, we further consider two strategies by which to combine these fixed pooling operations. One of these strategies is “unresponsive” to the characteristics of the region being pooled; the learning process in this strategy will result in an effective pooling operation that is some specific, unchanging “mixture” of max and average. To emphasize this unchanging mixture, we refer to this strategy as *mixed max-average pooling*.

The other strategy is “responsive” to the characteristics of the region being pooled; the learning process in this strategy results in a “gating mask”. This learned gating mask is then used to determine a “responsive” mix of max pooling and average pooling; specifically, the value of the inner product between the gating mask and the current region being pooled is fed through a sigmoid, the output of which is used as the mixing proportion between max and average. To emphasize the role of the gating mask in determining the “responsive” mixing proportion, we refer to this strategy as *gated max-average pooling*.

Both the mixed strategy and the gated strategy involve combinations of fixed pooling operations; a complementary generalization to these strategies is to learn the pooling operations themselves. From this, we are in turn led to consider learning pooling operations and also learning to combine those pooling operations. Since these combinations can be considered within the context of a binary tree structure, we refer to this approach as *tree pooling*. We pursue further details in the following

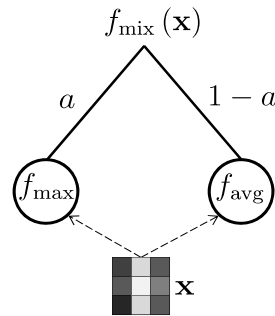


Fig. 1: Illustration of proposed “mixed” max-average pooling operations. \mathbf{x} is referred to as an input and α is the parameter balancing the importance of the max pooling and the average pooling operations.

sections.

3.1 Combining max and average pooling functions

3.1.1 “Mixed” max-average pooling

The conventional pooling operation is fixed to be either a simple average $f_{\text{ave}}(\mathbf{x}) = \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i$ or a maximum operation $f_{\text{max}}(\mathbf{x}) = \max_i \mathbf{x}_i$, where the vector \mathbf{x} contains the activation values from a local pooling region of N pixels (typical pooling region dimensions are 2×2 or 3×3) in an image or a channel.

At present, max pooling is often used as the default in CNNs. We touch on the relative performance of max pooling and, e.g., average pooling as part of a collection of exploratory experiments to test the invariance properties of pooling functions under common image transformations (including rotation, translation, and scaling); see Figure 4. The results indicate that, on the evaluation dataset, there are regimes in which either max pooling or average pooling demonstrates better performance than the other (although we observe that both of these choices are outperformed by our proposed pooling operations). In the light of observation that neither max pooling nor average pooling dominates the other, a first natural generalization is the strategy we call “mixed” max-average pooling, in which we learn specific mixing proportion parameters from the data. When learning such mixing proportion parameters one has several options (listed in order of increasing number of parameters): learning one mixing proportion parameter (a) per net, (b) per layer, (c) per layer/region being pooled (but used for all channels across that region), (d) per layer/channel (but used for all regions in each channel) (e) per layer/region/channel combination.

The form for each “mixed” pooling operation (written here for the “one per layer” option; the expression for other options differs only in the subscript of the mixing proportion a) is:

$$f_{\text{mix}}(\mathbf{x}) = a_{\ell} \cdot f_{\text{max}}(\mathbf{x}) + (1 - a_{\ell}) \cdot f_{\text{avg}}(\mathbf{x}), \quad (1)$$

where $a_{\ell} \in [0, 1]$ is a scalar mixing proportion specifying the specific combination of max and average; the

subscript ℓ is used to indicate that this equation is for the “one per layer” option. Figure 1 gives an illustration of the proposed “mixed” pooling operation. Once the output loss function E is defined, we can automatically learn each mixing proportion a (where we now suppress any subscript specifying which of the options we choose). Vanilla backpropagation for this learning is given by

$$\frac{\partial E}{\partial a} = \frac{\partial E}{\partial f_{\text{mix}}(\mathbf{x})} \frac{\partial f_{\text{mix}}(\mathbf{x})}{\partial a} = \delta (\max_i \mathbf{x}_i - \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i), \quad (2)$$

where $\delta = \partial E / \partial f_{\text{mix}}(\mathbf{x})$ is the error backpropagated from the following layer. Since pooling operations are typically placed in the midst of a deep neural network, we also need to compute the error signal to be propagated back to the previous layer:

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{x}_i} &= \frac{\partial E}{\partial f_{\text{mix}}(\mathbf{x}_i)} \frac{\partial f_{\text{mix}}(\mathbf{x}_i)}{\partial \mathbf{x}_i} \\ &= \delta \left[a \cdot \mathbf{1}[\mathbf{x}_i = \max_j \mathbf{x}_j] + (1 - a) \cdot \frac{1}{N} \right], \end{aligned} \quad (3)$$

where $\mathbf{1}[\cdot]$ denotes the 0/1 indicator function. In the experiment section, we report results for the “one parameter per pooling layer” option; the network for this experiment has 2 pooling layers and so has 2 more parameters than a network using standard pooling operations. We found that even this simple option yielded a surprisingly large performance boost. We also obtain results for a simple 50/50 mix of max and average, as well as for the option with the largest number of parameters: one parameter for each combination of layer/channel/region, or $pc \times ph \times pw$ parameters for each “mixed” pooling layer using this option (where pc is the number of channels being pooled by the pooling layer, and the number of spatial regions being pooled in each channel is $ph \times pw$). We observe that the increase in the number of parameters is not met with a corresponding boost in performance, and so we pursue the “one per layer” option.

3.1.2 “Gated” max-average pooling

In the previous section we considered a strategy that we referred to as “mixed” max-average pooling; in that strategy we learned a mixing proportion to be used in combining max pooling and average pooling. As mentioned earlier, once learned, each mixing proportion a remains fixed — it is “nonresponsive” insofar as it remains the same no matter what characteristics are present in the region being pooled. We now consider a “responsive” strategy that we call “gated” max-average pooling. In this strategy, rather than directly learning a mixing proportion that will be fixed after learning, we instead learn a “gating mask” (with spatial dimensions matching that of the regions being pooled). The scalar result of the inner product between the gating mask and the region being pooled is fed through a sigmoid to produce the value that we use as the mixing proportion. This strategy means that the actual mixing proportion

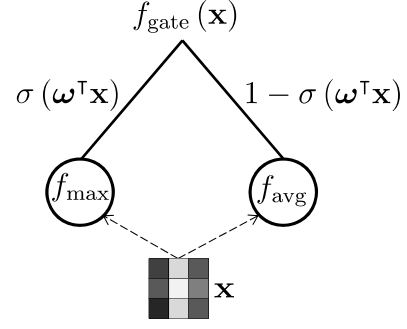


Fig. 2: Illustration of proposed “gated” max-average pooling operations. \mathbf{x} is referred to as an input and ω denotes the gating function balancing the importance of the max pooling and the average pooling operations.

can vary during use depending on characteristics present in the region being pooled. To be more specific, suppose we use \mathbf{x} to denote the values in the region being pooled and ω to denote the values in a “gating mask”. The “responsive” mixing proportion is then given by $\sigma(\omega^T \mathbf{x})$, where $\sigma(\omega^T \mathbf{x}) = 1 / (1 + \exp\{-\omega^T \mathbf{x}\}) \in [0, 1]$ is a sigmoid function.

Analogous to the strategy of learning a mixing proportion parameter, when learning gating masks one has several options (listed in order of increasing number of parameters): learning one gating mask (a) per net, (b) per layer, (c) per layer/region being pooled (but used for all channels across that region), (d) per layer/channel (but used for all regions in each channel) (e) per layer/region/channel combination. We suppress the subscript denoting the specific option, since the equations are otherwise identical for each option. Figure 2 gives an illustration of the proposed “gated” pooling operation.

The resulting pooling operation for this “gated” max-average pooling is:

$$f_{\text{gate}}(\mathbf{x}) = \sigma(\omega^T \mathbf{x}) f_{\text{max}}(\mathbf{x}) + (1 - \sigma(\omega^T \mathbf{x})) f_{\text{avg}}(\mathbf{x}). \quad (4)$$

We can compute the gradient with respect to the internal “gating mask” ω using the same procedure considered previously, yielding

$$\begin{aligned} \frac{\partial E}{\partial \omega} &= \frac{\partial E}{\partial f_{\text{gate}}(\mathbf{x})} \frac{\partial f_{\text{gate}}(\mathbf{x})}{\partial \omega} \\ &= \delta \sigma(\omega^T \mathbf{x}) (1 - \sigma(\omega^T \mathbf{x})) \mathbf{x} (\max_i \mathbf{x}_i - \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i), \end{aligned} \quad (5)$$

and

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{x}_i} &= \frac{\partial E}{\partial f_{\text{gate}}(\mathbf{x}_i)} \frac{\partial f_{\text{gate}}(\mathbf{x}_i)}{\partial \mathbf{x}_i} \\ &= \delta \left[\sigma(\omega^T \mathbf{x}) (1 - \sigma(\omega^T \mathbf{x})) \omega_i (\max_i \mathbf{x}_i - \frac{1}{N} \sum_{i=1}^N \mathbf{x}_i) \right. \\ &\quad \left. + \sigma(\omega^T \mathbf{x}) \cdot \mathbf{1}[\mathbf{x}_i = \max_j \mathbf{x}_j] + (1 - \sigma(\omega^T \mathbf{x})) \frac{1}{N} \right]. \end{aligned} \quad (6)$$

In a head-to-head parameter count, every single mixing proportion parameter a in the “mixed” max-average pooling strategy corresponds to a gating mask ω in the “gated” strategy (assuming they use the same parameter count option). To take a specific example, suppose that we consider a network with 2 pooling layers and pooling regions that are 3×3 . If we use the “mixed” strategy and the per-layer option, we would have a total of $2 = 2 \times 1$ extra parameters relative to standard pooling. If we use the “gated” strategy and the per-layer option, we would have a total of $18 = 2 \times 9$ extra parameters, where 9 is the number of parameters in each gating mask. The “mixed” strategy detailed immediately above uses fewer parameters and is “nonresponsive”; the “gated” strategy involves more parameters and is “responsive”. In our experiments, we find that “mixed” (with one mix per pooling layer) is outperformed by “gated” with one gate per pooling layer. Interestingly, an 18 parameter “gated” network with only one gate per pooling layer also outperforms a “mixed” option with far more parameters (40,960 with one mix per layer/channel/region) — except on the relatively large SVHN dataset. We touch on this below; Section 5 contains details.

3.1.3 Quick comparison: mixed and gated pooling

The results in Table 1 indicate the benefit of learning pooling operations over not learning. Within learned pooling operations, we see that when the number of parameters in the mixed strategy is increased, performance improves; however, parameter count is not the entire story. We see that the “responsive” gated max-avg strategy consistently yields better performance (using 18 extra parameters) than is achieved with the $>40k$ extra parameters in the 1 per layer/rg/ch “non-responsive” mixed max-avg strategy. The relatively larger SVHN dataset provides the sole exception (SVHN has $\approx 600k$ training images versus $\approx 50k$ for MNIST, CIFAR10, and CIFAR100) — we found baseline 1.89%, 50/50 mix 1.84%, mixed (1 per lyr) 1.76%, mixed (1 per lyr/ch/rg) 1.64%, and gated (1 per lyr) 1.74%.

3.2 Tree pooling

The strategies described above each involve combinations of fixed pooling operations; another natural generalization of pooling operations is to allow the pooling operations that are being combined to themselves be learned. These pooling layers remain distinct from convolution layers since pooling is performed separately within each channel; this channel isolation also means that even the option that introduces the largest number of parameters still introduces far fewer parameters than a convolution layer would introduce. The most basic version of this approach would not involve combining learned pooling operations, but simply learning pooling operations in the form of the values in pooling filters. One step further brings us to what we refer to as *tree*

| Method | MNIST | CIFAR10 | CIFAR10 ⁺ | CIFAR100 |
|---|------------------------|-----------------------|-----------------------|------------------------|
| Baseline | 0.39 ± 0.031 | 9.01 ± 0.096 | 7.22 ± 0.099 | 34.38 ± 0.096 |
| w/ Stochastic no learning | 0.38 ± 0.04 | 8.50 ± 0.05 | 7.30 ± 0.07 | 33.48 ± 0.27 |
| w/ 50/50 mix no learning | 0.34 ± 0.012 | 8.11 ± 0.10 | 6.78 ± 0.17 | 33.53 ± 0.16 |
| w/ Mixed 1 per pool layer 2 extra params | 0.33 ± 0.018 | 8.09 ± 0.19 | 6.62 ± 0.21 | 33.51 ± 0.11 |
| w/ Mixed 1 per layer/ch/rg >40k extra params | 0.30 ± 0.012 | 8.05 ± 0.16 | 6.58 ± 0.30 | 33.35 ± 0.19 |
| w/ Gated 1 per pool layer 18 extra params | 0.29 ± 0.016 | 7.90 ± 0.07 | 6.36 ± 0.28 | 33.22 ± 0.16 |

TABLE 1: Classification error (in %) comparison between baseline model (trained with conventional max pooling) and corresponding networks in which max pooling is replaced by the pooling operation listed. A superscripted + indicates the standard data augmentation as in [24], [27], [38]. We report means and standard deviations over 3 separate trials without model averaging.

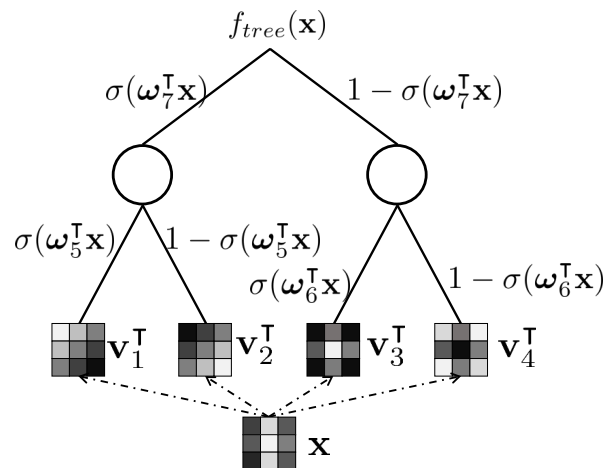


Fig. 3: Illustration of proposed tree pooling operation (3 levels in this figure). We indicate the input being pooled by x , gating masks by ω , and pooling filters by v (subscripted as appropriate).

pooling, in which we learn pooling filters and also learn to responsively combine those learned filters.

Both aspects of this learning are performed within a binary tree (with number of levels that is pre-specified rather than “grown” as in traditional decision trees) in which each leaf is associated with a pooling filter learned during training. As we consider internal nodes of the tree, each parent node is associated with an output value that is the mixture of the child node output values, until we finally reach the root node. The root node corresponds to the overall output produced by the tree and each of the mixtures (by which child outputs are “fused” into a parent output) is responsively learned. Tree pooling is intended (1) to learn pooling filters directly from the data; (2) to learn how to “mix” leaf node pooling filters in a differentiable fashion; (3) to bring together these other characteristics within a hierarchical

tree structure.

When the mixing of the leaf node pooling filters is allowed to be “responsive”, the resulting tree pooling operation becomes an integrated method for learning pooling filters and fusions of those filters that can display a range of different behaviors depending on the characteristics of the region being pooled. To more fully explore hidden structure and to potentially accommodate heterogeneous subspaces in complex data, we are motivated to further incorporate higher order operations into pooling operations. In particular, we propose to use a decision tree structure to reflect hierarchical characteristics of natural images. A decision tree consists of (internal) decision nodes and terminal nodes. Here we do not start from any particular pooling function (such as max or average pooling) but instead use (learnable) pooling filters. Figure 3 gives an illustration of the proposed tree pooling operation.

Each leaf node in our tree is associated with a “pooling filter” that will be learned; for a node with index m , we denote the pooling filter by $\mathbf{v}_m \in \mathbb{R}^N$. If we had a “degenerate tree” consisting of only a single (leaf) node, pooling a region $\mathbf{x} \in \mathbb{R}^N$ would result in the scalar value $\mathbf{v}_m^\top \mathbf{x}$. For (internal) nodes (at which two child values are combined into a single parent value), we proceed in a fashion analogous to the case of gated max-average pooling, with learned “gating masks” denoted (for an internal node m) by $\omega_m \in \mathbb{R}^N$. The “pooling result” at any arbitrary node m is thus

$$f_m(\mathbf{x}) = \begin{cases} \mathbf{v}_m^\top \mathbf{x} & \text{if leaf node} \\ \sigma(\omega_m^\top \mathbf{x}) f_{m,\text{left}}(\mathbf{x}) + (1 - \sigma(\omega_m^\top \mathbf{x})) f_{m,\text{right}}(\mathbf{x}) & \text{if internal node} \end{cases} \quad (7)$$

where $f_{m,\text{left}}(\mathbf{x})$ and $f_{m,\text{right}}(\mathbf{x})$ denote left and right child nodes of $f_m(\mathbf{x})$. The overall pooling operation would thus be the result of evaluating $f_{\text{root_node}}(\mathbf{x})$. The appeal of this tree pooling approach would be limited if one could not train the proposed layer in a fashion that was integrated within the network as a whole. This would be the case if we attempted to directly use a traditional decision tree, since its output presents points of discontinuity with respect to its inputs. The reason for the discontinuity (with respect to input) of traditional decision tree output is that a decision tree makes “hard” decisions; in the terminology we have used above, a “hard” decision node corresponds to a mixing proportion that can only take on the value 0 or 1. The consequence is that this type of “hard” function is not differentiable (nor even continuous with respect to its inputs), and this in turn interferes with any ability to use it in iterative parameter updates during backpropagation. This motivates us to instead use the internal node sigmoid “gate” function $\sigma(\omega_m^\top \mathbf{x}) \in [0, 1]$ so that the tree pooling function as a whole will be differentiable with respect to its parameters and its inputs.

For the specific case of a “2 level” tree (with leaf nodes “1” and “2” and internal node “3”) pooling function $f_{\text{tree}}(\mathbf{x}) = \sigma(\omega_3^\top \mathbf{x}) \mathbf{v}_1^\top \mathbf{x} + (1 - \sigma(\omega_3^\top \mathbf{x})) \mathbf{v}_2^\top \mathbf{x}$, we can use

the chain rule to compute the gradients with respect to the leaf node pooling filters $\mathbf{v}_1, \mathbf{v}_2$ and the internal node gating mask ω_3 :

$$\frac{\partial E}{\partial \mathbf{v}_1} = \frac{\partial E}{\partial f_{\text{tree}}(\mathbf{x})} \frac{\partial f_{\text{tree}}(\mathbf{x})}{\partial \mathbf{v}_1} = \delta \sigma(\omega_3^\top \mathbf{x}) \mathbf{x}, \quad (8)$$

$$\frac{\partial E}{\partial \mathbf{v}_2} = \frac{\partial E}{\partial f_{\text{tree}}(\mathbf{x})} \frac{\partial f_{\text{tree}}(\mathbf{x})}{\partial \mathbf{v}_2} = \delta (1 - \sigma(\omega_3^\top \mathbf{x})) \mathbf{x}, \quad (9)$$

$$\frac{\partial E}{\partial \omega_3} = \frac{\partial E}{\partial f_{\text{tree}}(\mathbf{x})} \frac{\partial f_{\text{tree}}(\mathbf{x})}{\partial \omega_3} = \delta \sigma(\omega_3^\top \mathbf{x}) (1 - \sigma(\omega_3^\top \mathbf{x})) \mathbf{x} (\mathbf{v}_1^\top - \mathbf{v}_2^\top) \mathbf{x}. \quad (10)$$

The error signal to be propagated back to the previous layer is

$$\begin{aligned} \frac{\partial E}{\partial \mathbf{x}} &= \frac{\partial E}{\partial f_{\text{tree}}(\mathbf{x})} \frac{\partial f_{\text{tree}}(\mathbf{x})}{\partial \mathbf{x}} \\ &= \delta [\sigma(\omega_3^\top \mathbf{x}) (1 - \sigma(\omega_3^\top \mathbf{x})) \omega_3 (\mathbf{v}_1^\top - \mathbf{v}_2^\top) \mathbf{x} \\ &\quad + \sigma(\omega_3^\top \mathbf{x}) \mathbf{v}_1 + (1 - \sigma(\omega_3^\top \mathbf{x})) \mathbf{v}_2]. \end{aligned} \quad (11)$$

3.2.1 Quick comparison: tree pooling

Table 2 collects results related to tree pooling. We observe that on all datasets but the comparatively simple MNIST, adding a level to the tree pooling operation improves performance. However, even further benefit is obtained from the use of tree pooling in the first pooling layer and gated max-avg in the second. In Table 4 we compare the results of this configuration against recent comparable methods.

| Method | MNIST | CIFAR10 | CIFAR10 ⁺ | CIFAR100 | SVHN |
|---------------------------|-----------------|-----------------|----------------------|------------------|-----------------|
| Our baseline | 0.39 ± 0.031 | 9.01 ± 0.096 | 7.22 ± 0.099 | 34.38 ± 0.096 | 1.89 ± 0.069 |
| Tree | 0.34 ± 0.028 | 8.52 ± 0.175 | 6.54 ± 0.156 | 33.64 ± 0.285 | 1.81 ± 0.047 |
| 2 level; 1 per pool layer | | | | | |
| Tree | 0.38 ± 0.032 | 8.43 ± 0.091 | 6.38 ± 0.165 | 32.85 ± 0.181 | 1.73 ± 0.096 |
| 3 level; 1 per pool layer | | | | | |
| Tree+Max-Avg | 0.31 ± 0.031 | 7.61 ± 0.121 | 6.02 ± 0.047 | 32.87 ± 0.278 | 1.70 ± 0.069 |
| 1 per pool layer | | | | | |

TABLE 2: Classification error (in %) comparison between our baseline model (trained with conventional max pooling) and proposed methods involving tree pooling. A superscripted ⁺ indicates the standard data augmentation as in [24], [27], [38].

Comparison with making the network deeper using conv layers. To further investigate whether simply adding depth to our baseline network gives a performance boost comparable to that observed for our proposed pooling operations, we report in Table 3 below some additional experiments on CIFAR10 (error rate in percent; no data augmentation). If we count depth by counting any layer with learned parameters as an extra layer of depth (even if there is only 1 parameter), the number of parameter layers in a baseline network with 2 additional standard convolution layers matches the number of parameter layers in our best performing net (although the convolution layers contain many more parameters).

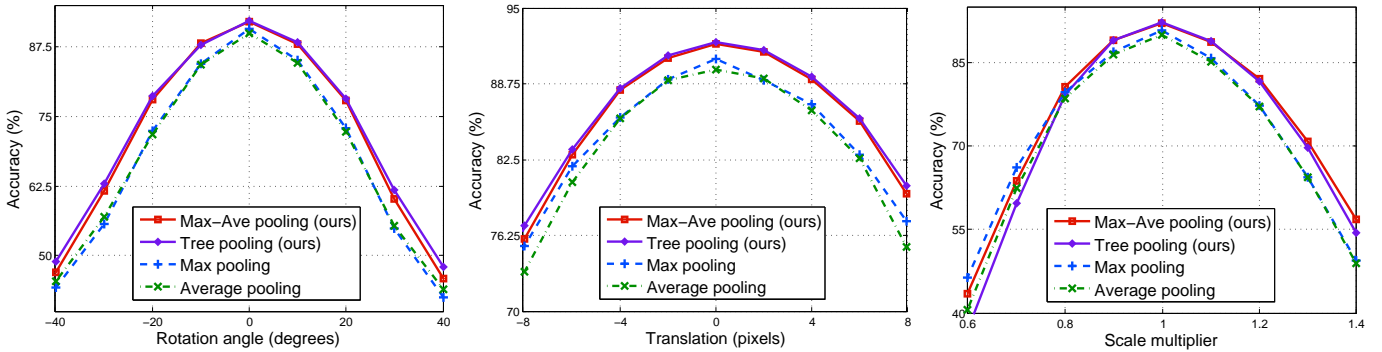


Fig. 4: Controlled experiment on CIFAR10 investigating the relative benefit of selected pooling operations in terms of robustness to three types of data variation. The three kinds of variations we choose to investigate are rotation, translation, and scale. With each kind of variation, we modify the CIFAR10 test images according to the listed amount. We observe that, across all types and amounts of variation (except extreme down-scaling) the proposed pooling operations investigated here (gated max-avg and 2 level tree pooling) provide improved robustness to these transformations, relative to the standard choices of maxpool or avgpool.

Our method requires only 72 extra parameters and obtains state-of-the-art 7.61% error. On the other hand, making networks deeper with conv layers adds many more parameters but yields test error that does not drop below 8.98% in the configuration explored. Since we follow each additional conv layer with a ReLU, these networks correspond to increasing nonlinearity as well as adding depth and adding (many) parameters. These experiments indicate that the performance of our proposed pooling is not accounted for as a simple effect of the addition of depth/parameters/nonlinearity.

We also perform an experiment comparing grouped convolutional layers using ReLU activation function in the form of ¹:

$$\text{gconv}(3 \times 3 \times 5, n) - \text{relu} - \text{gconv}(1 \times 1 \times 3, n) - \text{relu} - \text{gconv}(1 \times 1 \times 1, n), \quad (12)$$

where $\text{gconv}(h \times w \times c, m)$ denotes m groups of convolutions with kernel size $h \times w$, n input channels, and c output channels. This operation performs blocks of convolution in a channel-wise fashion and therefore has a closer behavior to the proposed tree pooling operation. With this experimental setting, we observe better performance than simply adding more convolutional layers, but worse performance than the proposed tree pooling and Tree + (gated) Max-Avg (see Table 3). The proposed tree structure function achieves the lowest error rate in this comparison.

Comparison with alternative pooling layers. To see whether we might find similar performance boosts by replacing the max pooling in the baseline network configuration with alternative pooling operations such as stochastic pooling, “pooling” using a stride 2 convolution layer as pooling (cf All-CNN), or a simple fixed 50/50 proportion in max-avg pooling, we performed another set of experiments on unaugmented CIFAR10.

1. This was recommended by one of the anonymous reviewers.

| Method | % Error | Extra parameters |
|--------------------------------|------------------|------------------|
| Baseline | 9.01 ± 0.096 | 0 |
| w/ 1 extra conv layer (+ReLU) | 8.98 ± 0.058 | 0.6M |
| w/ 2 extra conv layers (+ReLU) | 9.25 ± 0.077 | 1.2M |
| w/ grouped conv layers (+ReLU) | 8.85 ± 0.149 | 63 |
| w/ Tree + (gated) Max-Avg | 7.61 ± 0.121 | 72 |

TABLE 3: Classification error (%) on CIFAR10 (without data augmentation) comparison between networks made deeper with convolution layers and proposed Tree+(gated) Max-Avg pooling.

From the baseline error rate of 9.01%, replacing each of the 2 max pooling layers with stacked stride 2 conv:ReLU (as in [38]) lowers the error to 8.77%, but adds 0.5M extra parameters. Using stochastic pooling [46] adds computational overhead but no parameters and results in 8.50% error. A simple 50/50 mix of max and average is computationally light and yields 8.11% error with no additional parameters. Finally, our tree+gated max-avg configuration adds 72 parameters and achieves a state-of-the-art 7.61% error.

4 QUICK PERFORMANCE OVERVIEW

For ease of discussion, we collect here observations from subsequent experiments with a view to highlighting aspects that shed light on the performance characteristics of our proposed pooling functions.

First, as seen in the experiment shown in Figure 4 replacing standard pooling operations with either gated max-avg or (2 level) tree pooling (each using the “one per layer” option) yielded a boost (relative to max or avg pooling) in CIFAR10 test accuracy as the test images underwent three different kinds of transformations. This boost was observed across the entire range of transformation amounts for each of the transformations (with the exception of extreme downscaling). We already observe improved robustness in this initial experiment

and intend to investigate more instances of our proposed pooling operations as time permits.

Second, the performance that we attain in the experiments reported in Figure 4, Table 1, Table 2, Table 4, and Table 5 is achieved with very modest additional numbers of parameters — e.g. on CIFAR10, our best performance (obtained with the tree+gated max-avg configuration) only uses an additional 72 parameters (above the 1.8M of our baseline network) and yet reduces test error from 9.01% to 7.61%; see the **CIFAR10** Section for details. In our AlexNet experiment, replacing the maxpool layers with our proposed pooling operations gave a 6% relative reduction in test error (top-5, single-view) with only 45 additional parameters (above the >50M of standard AlexNet); see the **ImageNet 2012** Section for details. We also investigate the additional time incurred when using our proposed pooling operations; in the experiments reported in the **Timing** section, this overhead ranges from 5% to 15%.

Testing invariance properties. Before going to the overall classification results, we investigate the invariance properties of networks utilizing either standard pooling operations (max and average) or two instances of our proposed pooling operations (gated max-avg and 2 level tree, each using the “1 per pool layer” option) that we find to yield best performance (see Sec. 5 for architecture details used across each network). We begin by training four different networks on the CIFAR10 training set, one for each of the four pooling operations selected for consideration; training details are found in Sec. 5. We seek to determine the respective invariance properties of these networks by evaluating their accuracy on various transformed versions of the CIFAR10 test set. Figure 4 illustrates the test accuracy attained in the presence of image rotation, (vertical) translation, and scaling of the CIFAR10 test set.

It is interesting to note that (with a few exceptions) the accuracy curves for max pooling and average pooling are quite close, as are the accuracy curves for gated max-average and for 2 level tree pooling. One might expect that the performance of (gated) max-average pooling would perhaps resemble the point-wise best of max or average, the reasoning being that it is a combination of these two. We see that this appears to not capture the story as shown in our experiment — other than in the lower regime of scale multiplier values, our gated max-avg pooling performance is not simply comparable to, but notably better than the performance of either max or average pooling. One might also speculate that because max-average pooling has access to the highly nonlinear max operation, it might be able to leverage this to potentially outperform a 2 level tree pooling operation. These invariance experiments, in contrast, indicate that tree pooling performs as well as (and perhaps slightly better than) gated max-average pooling; the exception is found at the extreme regimes of scale multiplier values. One might explain the tree pooling performance drop

in these extreme scale multiplier value regimes as an indication that the basic pooling masks utilized at the leaf nodes are unable to effectively respond to patterns across significant scaling ranges — at least for this 2 level case.

Timing. In order to evaluate how much additional time is incurred by the use of our proposed learned pooling operations, we measured the average forward+backward time per CIFAR10 image. In each case, the one per layer option is used. We find that the additional computation time incurred ranges from 5% to 15%. More specifically, the baseline network took 3.90 ms; baseline with mixed max-avg took 4.10 ms; baseline with gated max-avg took 4.16 ms; baseline with 2 level tree pooling took 4.25 ms; finally, baseline with tree+gated max-avg took 4.46 ms.

5 EXPERIMENTS

We evaluate the proposed max-average pooling and tree pooling approaches on five standard benchmark datasets: MNIST [22], CIFAR10 [18], CIFAR100 [18], SVHN [30] and ImageNet [34]. To control for the effect of differences in data or data preparation, we match our data and data preparation to that used in [24]. Please refer to [24] for the detailed description. We also seek to control for the effect of hyperparameter settings by using the same hyperparameter settings. We will use the term hyperparameters to collectively refer to number of layers, number of channels, dimensions of pooling regions, dropout rate, learning rate, learning rate schedule, momentum, weight decay, and parameter initialization.

We now describe the basic network architecture and then will specify the various hyperparameter choices. The basic experiment architecture contains six 3×3 standard convolutional layers (named conv1 to conv6) and three mlpconv layers (named mlpconv1 to mlpconv3) [27], placed after conv2, conv4, and conv6, respectively. We chose the number of channels at each layer to be analogous to the choices in [24], [27]; the specific numbers are provided in the sections for each dataset. We follow every one of these conv-type layers with ReLU activation functions. One final mlpconv layer (mlpconv4) is used to reduce the dimension of the last layer to match the total number of classes for each different dataset, as in [27]. The overall model has parameter count analogous to [24], [27]. The proposed max-average pooling and tree pooling layers with 3×3 pooling regions are used after mlpconv1 and mlpconv2 layers². We provide a detailed listing of the network configurations in Table 6.

Moving on to the hyperparameter settings, dropout with rate 0.5 is used after each pooling layer. We also use hidden layer supervision to ease the training process as in [24]. The learning rate is decreased whenever the validation error stops decreasing; we use the schedule

2. There is one exception: on the very small images of the MNIST dataset, the second pooling layer uses 2×2 pooling regions.

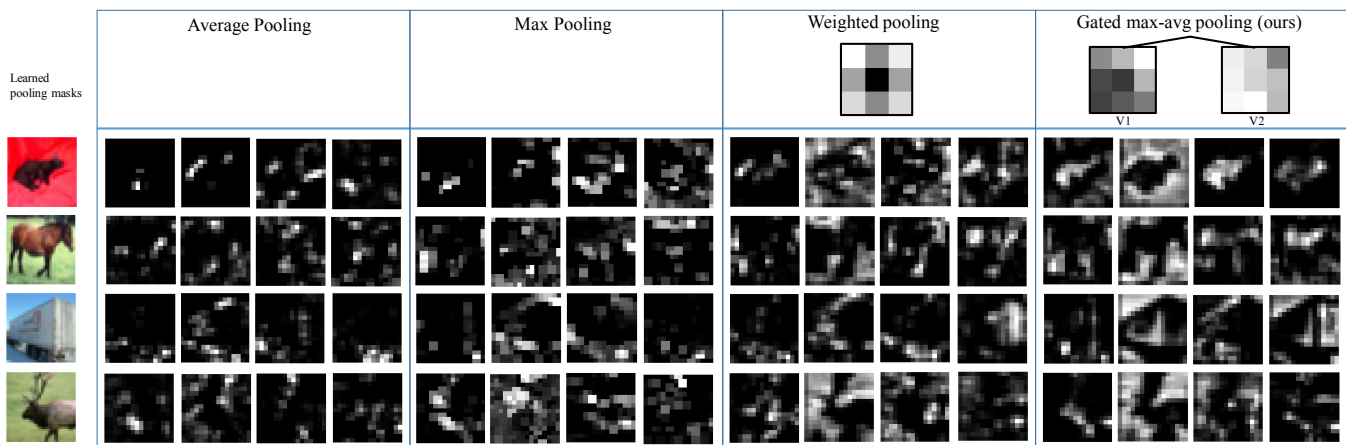


Fig. 5: Visualization of the learned pooling masks of weighted pooling and the proposed gated max-average pooling function on CIFAR-10 dataset. Here we denote weighted pooling to a single learned pooling filter without the tree structure (i.e., a singleton leaf node containing 9 parameters; one such singleton leaf node per pooling layer). We also visualize the output feature maps from different pooling methods, including max pooling, average pooling, weighted pooling, and gated max-average pooling. We can see that the feature responses of learnable pooling functions (weight and gated pooling) encode much of the structure in the image, as some of it is lost when pooling without learning (average or max pooling) is used.

{0.025, 0.0125, 0.0001} for all experiments. The momentum of 0.9 and weight decay of 0.0005 are fixed for all datasets as another regularizer besides dropout. All the initial pooling filters and pooling masks have values sampled from a Gaussian distribution with zero mean and standard deviation 0.5. We use these hyperparameter settings for all experiments reported in Tables 1, 2, and 3. No model averaging is done at test time.

5.1 Classification results

Tables 1 and 2 show our overall experimental results. Our baseline is a network trained with conventional max pooling. *Mixed* refers to the same network but with a max-avg pooling strategy in both the first and second pooling layers (both using the mixed strategy); *Gated* has a corresponding meaning. *Tree* (with specific number of levels noted below) refers to the same again, but with our tree pooling in the first pooling layer only; we do not see further improvement when tree pooling is used for both pooling layers. This observation motivated us to consider following a tree pooling layer with a gated max-avg pooling layer: *Tree+Max-Average* refers to a network configuration with (2 level) tree pooling for the first pooling layer and gated max-average pooling for the second pooling layer. All results are produced from the same network structure and hyperparameter settings — the only difference is in the choice of pooling function. See Table 6 for details.

MNIST. The MNIST dataset consists of 28×28 gray scale images from 10 different classes (the digits 0-9) with 60,000 training and 10,000 testing samples. Our MNIST model has {128, 128, 192, 192, 256, 256} channels for conv1 to conv6 and {128, 192, 256} channels for mlpcnv1 to mlpcnv3, respectively. Our only preprocessing

is mean subtraction. Tables 4, 1, and 2 show previous best results and those for our proposed pooling methods.

CIFAR10. The CIFAR10 dataset consists of 32×32 color images with 50,000 training examples and 10,000 testing examples. The dataset is preprocessed by using global contrast normalization and ZCA whitening as in [6], [24], [27]. Our CIFAR10 model has {128, 128, 192, 192, 256, 256} channels for conv1 to conv6 and {128, 192, 256} channels for mlpcnv1 to mlpcnv3, respectively. We also performed an experiment in which we learned a single pooling filter without the tree structure (i.e., a singleton leaf node containing 9 parameters; one such singleton leaf node per pooling layer) and obtained 0.3% improvement over the baseline model. We refer this single learned pooling filter to weighted pooling. Our results indicate that performance improves when the pooling filter is learned, and further improves when we also learn how to combine learned pooling filters.

In Figure 5 we visualize learned pooling masks of weighted pooling and the proposed gated max-average pooling function on CIFAR-10 dataset. The weighted pooling mask mimics a low-pass filter while gated pooling masks capture oriented edges. The gating mask (function) will then learn and determine how to combine the pooling responses based on the input feature maps (Equation. 4). We also visualize the output feature maps from different pooling methods, including average pooling, max pooling, weighted pooling, and gated max-average pooling in Figure 5. We can see that the feature responses of learnable pooling functions (weight and gated pooling) encode much of the structure in the image, as some of it is lost when pooling without learning (average or max pooling) is used.

The All-CNN method in [38] uses convolutional lay-

| Method | MNIST | CIFAR10 | CIFAR10 ⁺ | CIFAR100 | SVHN |
|---------------------|-------------|-------------|----------------------|--------------|-------------|
| CNN [15] | 0.53 | - | - | - | - |
| Stoch. Pooling [46] | 0.47 | 15.13 | - | 42.51 | 2.80 |
| Maxout Networks [6] | 0.45 | 11.68 | 9.38 | 38.57 | 2.47 |
| Prob. Maxout [39] | - | 11.35 | 9.39 | 38.14 | 2.39 |
| Tree Priors [40] | - | - | - | 36.85 | - |
| DropConnect [25] | 0.57 | 9.41 | 9.32 | - | 1.94 |
| FitNet [33] | 0.51 | - | 8.39 | 35.04 | 2.42 |
| NiN [27] | 0.47 | 10.41 | 8.81 | 35.68 | 2.35 |
| DSN [24] | 0.39 | 9.69 | 7.97 | 34.57 | 1.92 |
| NiN + LA units [1] | - | 9.59 | 7.51 | 34.40 | - |
| dasNet [42] | - | 9.22 | - | 33.78 | - |
| All-CNN [38] | - | 9.08 | 7.25 | 33.71 | - |
| R-CNN [26] | 0.31 | 8.69 | 7.09 | 31.75 | 1.77 |
| Our baseline | 0.39 | 9.01 | 7.22 | 34.38 | 1.89 |
| Our Tree+Max-Avg | 0.31 | 7.61 | 6.02 | 32.87 | 1.70 |

TABLE 4: Classification error (in %) reported by recent comparable publications on four benchmark datasets with a single model and no data augmentation, unless otherwise indicated. A superscripted ⁺ indicates the standard data augmentation as in [24], [27], [38]. A “-” indicates that the cited work did not report results for that dataset. A fixed network configuration using the proposed tree+max-avg pooling (1 per pool layer option) yields state-of-the-art performance on all datasets (with the exception of CIFAR100).

ers in place of pooling layers in a CNN-type network architecture. However, a standard convolutional layer requires many more parameters than a gated max-average pooling layer (only 9 parameters for a 3×3 pooling region kernel size in the 1 per pooling layer option) or a tree-pooling layer (27 parameters for a 2 level tree and 3×3 pooling region kernel size, again in the 1 per pooling layer option). The pooling operations in our tree+max-avg network configuration use $7 \times 9 = 63$ parameters for the (first, 3 level) tree-pooling layer — 4 leaf nodes and 3 internal nodes — and 9 parameters in the gating mask used for the (second) gated max-average pooling layer, while the best result in [38] contains a total of nearly 500,000 parameters in layers performing “pooling like” operations; the relative CIFAR10 accuracies are 7.61% (ours) and 9.08% (All-CNN).

For the data augmentation experiment, we followed the standard data augmentation procedure [24], [27], [38]. When training with augmented data, we observe the same trends seen in the “no data augmentation” experiments. We note that [7] reports a 4.5% error rate with extensive data augmentation (including translations, rotations, reflections, stretching, and shearing operations) in a much wider and deeper 50 million parameter network — 28 times more than are in our networks.

CIFAR100. The CIFAR100 dataset consists of 32×32 color images with 50,000 training and 10,000 testing images, but with 100 classes rather than 10. The number of images for each class is thus 500 instead of 5,000 as in CIFAR10. We preprocess the dataset by global contrast normalization and ZCA whitening as in [6]. Our CIFAR100 model has 192 channels for all convolutional layers and $\{96, 192, 192\}$ channels for mlpconv1 to mlpconv3, respectively. The fourth column of Table 4 shows recent comparable results.

Street view house numbers. The Street View House

Numbers (SVHN) dataset consists of 32×32 color images from Google Street View images: 73,257 digits for training, 26,032 digits for testing, and 531,131 extra training samples. We follow the same training procedure in [6] that we select 400 samples per class from the regular training set and 200 samples per class from the extra training set as the validation set. We preprocess the dataset by Local Contrast Normalization (LCN) as in [46]. Our SVHN model has $\{128, 128, 320, 320, 384, 384\}$ channels for conv1 to conv6 and $\{96, 256, 256\}$ channels for mlpconv1 to mlpconv3, respectively. In terms of amount of data, SVHN has a larger training data set (>600k versus the ≈ 50 k of most of the other benchmark datasets). The much larger amount of training data motivated us to explore what performance we might observe if we pursued the one per layer/channel/region option, which even for the simple mixed max-avg strategy results in a huge increase in total the number of parameters to learn in our proposed pooling layers: specifically, from a total of 2 in the mixed max-avg strategy, 1 parameter per pooling layer option, we increase to 40,960.

Using this one per layer/channel/region option for the mixed max-avg strategy, we observe test error (in %) of 0.30 on MNIST, 8.05 on CIFAR10, 6.58 on CIFAR10⁺, 33.35 on CIFAR100, and 1.64 on SVHN. Interestingly, for MNIST, CIFAR10, CIFAR10⁺, and CIFAR100 this mixed max-avg (1 per layer/channel/region) performance is between mixed max-avg (1 per layer) and gated max-avg (1 per layer); The SVHN result using mixed max-avg (1 per layer/channel/region) sets a new state of the art.

ImageNet 2012. The ImageNet 2012 dataset consists of 1.2 million training images, 50,000 validation, and 100,000 testing. In this experiment we do not directly compete with the best performing result in the challenge (since the winning methods [43] involve many additional aspects beyond pooling operations), but rather to provide an illustrative comparison of the relative benefit

of the proposed pooling methods versus conventional max pooling on this dataset. We use the same network structure and parameter setup as in [19] (no hidden layer supervision) but simply replace the first max pooling with the (proposed 2 level) tree pooling (2 leaf nodes and 1 internal node for $27 = 3 \times 9$ parameters) and replace the second and third max pooling with gated max-average pooling (2 gating masks for $18 = 2 \times 9$ parameters). Relative to the original AlexNet, this adds 45 more parameters (over the $>50M$ in the original) and achieves relative error reduction of 6% (for top-5, single-view) and 5% (for top-5, multi-view). Our GoogLeNet configuration uses 4 gated max-avg pooling layers, for a total of 36 extra parameters over the 6.8 million in standard GoogLeNet. Table 5 shows a direct comparison (in each case we use single net predictions rather than ensemble). Table 6 gives a summary of the network configurations used in the experiments.

| Method | top-1 s-view | top-5 s-view | top-1 m-view | top-5 m-view |
|---------------------|-----------------|-----------------|-----------------|-----------------|
| AlexNet [19] | 43.1 | 19.9 | 40.7 | 18.2 |
| AlexNet w/ ours | 41.4 | 18.7 | 39.3 | 17.3 |
| GoogLeNet [43] | - | 10.07 | - | 9.15 |
| GoogLeNet w/BN | 28.68 | 9.53 | 27.81 | 9.09 |
| GoogLeNet w/BN+ours | 28.02 | 9.16 | 27.60 | 8.93 |

TABLE 5: ImageNet 2012 test error (in %). BN denotes Batch Normalization [13].

5.2 Observations from experiments

In each experiment, using any of our proposed pooling operations boosted performance. A fixed network configuration using the proposed tree+max-avg pooling (1 per pool layer option) yields state-of-the-art performance on MNIST, CIFAR10 (with and without data augmentation), and SVHN. We observed boosts in tandem with data augmentation, multi-view predictions, batch normalization, and several different architectures — network in network style, deeply-supervised nets style, the $>50M$ parameter AlexNet, and the 22-layer GoogLeNet.

5.3 Visualization of network internal representations

To gain additional qualitative understanding of the pooling methods we are considering, we use the popular t-SNE [44] algorithm to visualize embeddings of some internal feature responses from pooling operations. Specifically, we again use four networks (one utilizing each of the selected types of pooling) trained on the CIFAR10 training set (see Sec. 5 for architecture details used across each network). We extract feature responses for a randomly chosen 800-image subset of the CIFAR10 test set at the first (i.e., earliest) and second pooling layers of each network. These feature response vectors are then embedded into 2-D using t-SNE; see Figure 6.

The first column shows the embeddings of the internal activations immediately after the first pooling operation; the second column shows embeddings of activations immediately after the second pooling operation. From top to bottom, we plot the t-SNE embeddings of the pooling activations within networks that are trained with average, max, gated max-avg, and (2 level) tree pooling. We can see that certain classes such as “0” (airplane), “2” (bird), and “9” (truck) are more separated with the proposed methods than they are with the conventional average and max pooling functions. We can also see that the embeddings of the second-pooling-layer activations are generally more separable than the embeddings of first-pooling-layer activations.

5.4 Relationship to gated convolutional layer

In the previous sections we demonstrated that aspects of mixed, gated, and tree behaviors can be incorporated into the pooling functions of a CNN framework, recalling that a pooling function operates on spatially-local regions in a channel-by-channel fashion; that is, spatial information is condensed but channels remain separate. Another layer type, specifically a convolutional layer, is intended to exploit local correlation by computing feature response maps across different channels within corresponding receptive fields; that is, spatial information remains separate, but channels can interact.

It can be a challenging task to design CNN models that possess good performance without excessive model size. Factors such as the number of layers and number of output channels require careful consideration during the CNN design process because even modest increases in these factors can lead to much greater increases in the dimensions of the output responses in subsequent layers. In light of this, rather than simply increasing the depth and the number of channels of a CNN model, we instead investigate a method to improve performance by integrating aspects of gated behavior into convolutional layers constructed so as to incorporate both learning of the convolutional features and learning how to fuse those convolutional features. In keeping with our earlier terminology, we call these layers gated convolutional (gated_conv) layers. We consider a gated convolutional layer of the form:

$$f_{\text{gate_conv}}(\mathbf{x}) = \sigma(W\mathbf{x}) \circ V_1\mathbf{x} + (\mathbf{1} - \sigma(W\mathbf{x})) \circ V_2\mathbf{x}, \quad (13)$$

where \mathbf{x} denote the vectorized feed-forward and input patches, $\{W, V_1, V_2\} \in \mathbb{R}^{N_{\text{out}} \times N_{\text{in}}}$ are the kernel matrices interacting with input \mathbf{x} . Here N_{out} is the number of output channels and N_{in} is the number of input channels, \circ denotes the Hadamard product, and $\mathbf{1}$ represents the vector in which all elements are 1. The “gating function” $\sigma(W\mathbf{x}) = 1/(1 + \exp\{-W\mathbf{x}\}) \in [0, 1]^{N_{\text{out}}}$ specifies the weights to use when adding the different responses $V_1\mathbf{x}$ and $V_2\mathbf{x}$ to produce the layer response result. The derivatives of the function $f_{\text{gate_conv}}(\mathbf{x})$ w.r.t its weights are analogous to what we saw for the gated and tree pooling



Fig. 6: t-SNE embeddings of the output responses from different pooling operations on the CIFAR10 test set (with classes indicated). From top to bottom: average, max, gated max-avg, and (2 level) tree pooling. The first and the second columns show the first and the second pooling layers, respectively. Best viewed in color.

| Network layer configurations reported in Tables 1, 2, and 4. | | | | | |
|--|--------------------------|--------------------------|------------------------------|------------------------------|--------------------------------|
| DSN (baseline) | mixed max-avg | gated max-avg | 2 level tree pool | 3 level tree pool | tree+gated max-avg pool |
| 3x3 (standard) conv | | | | | |
| 3x3 (standard) conv | | | | | |
| 1x1 mlpconv | | | | | |
| 3x3 maxpool | 3x3 mixed max-avg | 3x3 gated max-avg | 3x3 2 level tree pool | 3x3 3 level tree pool | 3x3 2/3 level tree pool |
| 3x3 (standard) conv | | | | | |
| 3x3 (standard) conv | | | | | |
| 1x1 mlpconv | | | | | |
| 3x3 maxpool | 3x3 mixed max-avg | 3x3 gated max-avg | 3x3 maxpool | 3x3 maxpool | 3x3 gated max-avg |
| 3x3 (standard) conv | | | | | |
| 3x3 (standard) conv | | | | | |
| 1x1 mlpconv | | | | | |
| 1x1 mlpconv | | | | | |
| 8x8 global vote | | | | | |

TABLE 6: Here we provide explicit statement of the experimental conditions (specifically, network layer configurations) explored in Tables 1, 2, and 4. We list all conv-like layers and pool-like layers, but ReLUs are suppressed to lighten the amount of text; these follow each standard conv layer. Also, all network configurations incorporate deep supervision after each standard convolution layer; this is also suppressed for clarity. We bold the changes made to the baseline DSN layer configuration. We now describe the meaning of entries in the table. Each column in the table lists the sequence of layer types used in that network configuration. When a row cell spans multiple columns (i.e. configurations), this indicates that the layer type listed in that cell is kept the same across the corresponding network configurations. Thus, every network in our experiments begins with a stacked pair of 3x3 (standard) conv layers followed by a 1x1 mlpconv layer. For a specific example, let us consider the network configuration in the column headed “mixed max-avg” - the sequence of layers in this configuration is: 3x3 (standard) conv, 3x3 (standard) conv, 1x1 mlpconv, **3x3 mixed max-avg pool**, 3x3 (standard) conv, 3x3 (standard) conv, 1x1 mlpconv, **3x3 mixed max-avg pool**, 3x3 (standard) conv, 3x3 (standard) conv, 1x1 mlpconv, 1x1 mlpconv, 8x8 global vote (cf. [27]) (we again omit mention of ReLUs and deep supervision). CIFAR100 uses (2 level) tree+max-avg; CIFAR10 uses (3 level) tree+max-avg. As a final note: for the MNIST experiments only, the second pooling operation uses 2x2 regions instead of the 3x3 regions used on the other datasets.

layers, and so we omit the details. One can imagine that the function $f_{\text{gate_conv}}(\mathbf{x})$ performs two individual convolutional operations via V_1 and V_2 matrices, and then fuses the two separate activation maps via W to produce the final output.

| Method | MNIST | CIFAR10 | CIFAR10 ⁺ | CIFAR100 |
|----------------------|-------|---------|----------------------|----------|
| Max pooling | 0.39 | 9.10 | 7.32 | 34.21 |
| Gated conv pooling | 0.35 | 8.59 | 6.98 | 33.58 |
| Gated Max-Avg (ours) | 0.29 | 7.90 | 6.36 | 32.37 |

TABLE 7: Classification error (in %) comparison between the baselines (max pooling and gated convolutional layer with stride 2 used in place of pooling) and the proposed gated max-avg pooling. A superscripted ⁺ indicates the standard data augmentation as in [24], [27], [38]. We can see that pooling by combining the responses from two sets of convolutional operations (gated conv pooling) reduces the error rates on all four cases. However, this cross-channel pooling operation does not outperform the channel-by-channel gated max-average pooling.

To see whether we might again find similar performance boosts by replacing the max pooling in the baseline network configuration with the introduced gated convolutional layer with stride 2, we performed another set of experiments on MNIST, CIFAR-10, and CIFAR-100 shown in Table 7. We observe that pooling by combining the responses from two sets of convolutional operations reduces the error rates on all four cases. However, this

cross-channel pooling operation does not outperform the channel-by-channel gated max-average pooling. Furthermore, each gated convolutional pooling requires around extra 0.3M parameters while each gated max-average pooling only adds 9 parameters. Gated max-average pooling is also computationally light compared to gated convolutional pooling.

ACKNOWLEDGMENTS

This work is supported by NSF awards IIS-1216528 (IIS-1360566), IIS-0844566(IIS-1360568), IIS-1618477, and a Northrop Grumman Contextual Robotics grant. We are grateful for the generous donation of the GPUs by NVIDIA. We thank the anonymous reviewers for their constructive comments, and in particular about the suggestion to compare with grouped convolutions.

REFERENCES

- [1] F. Agostinelli, M. Hoffman, P. Sadowski, and P. Baldi, “Learning activation functions to improve deep neural networks,” in *ICLR*, 2015.
- [2] Y. Boureau, N. Le Roux, F. Bach, J. Ponce, and Y. LeCun, “Ask the locals: multi-way local pooling for image recognition,” in *ICCV*, 2011.
- [3] Y. Boureau, J. Ponce, and Y. LeCun, “A Theoretical Analysis of Feature Pooling in Visual Recognition,” in *ICML*, 2010.
- [4] S. R. Buló and P. Kotschieder, “Neural Decision Forests for Semantic Image Labelling,” in *CVPR*, 2014.
- [5] A. Coates and A. Y. Ng, “Selecting receptive fields in deep networks,” in *NIPS*, 2011.

- [6] I. J. Goodfellow, D. Warde-Farley, M. Mirza, A. C. Courville, and Y. Bengio, "Maxout Networks," in *ICML*, 2013.
- [7] B. Graham, "Fractional Max-Pooling," *arXiv preprint arXiv:1412.6071*, 2014.
- [8] C. Gulcehre, K. Cho, R. Pascanu, and Y. Bengio, "Learned-norm pooling for deep feedforward and recurrent neural networks," in *MLKDD*, 2014.
- [9] K. He, X. Zhang, S. Ren, and J. Sun, "Spatial pyramid pooling in deep convolutional networks for visual recognition," in *ECCV*, 2014.
- [10] G. Hinton, S. Osindero, and Y.-W. Teh, "A fast learning algorithm for deep belief nets," *Neural Computation*, vol. 18, no. 7, pp. 1527–1554, 2006.
- [11] D. H. Hubel and T. N. Wiesel, "Receptive fields, binocular interaction and functional architecture in the cat's visual cortex," *Journal of Physiology*, vol. 160, no. 1, pp. 106–154, 1962.
- [12] Y. Ioannou, D. Robertson, D. Zikic, P. Kotschieder, J. Shotton, M. Brown, and A. Criminisi, "Decision forests, convolutional networks and the models in-between," *arXiv preprint arXiv:1603.01250*, 2016.
- [13] S. Ioffe and C. Szegedy, "Batch normalization: Accelerating deep network training by reducing internal covariate shift," in *ICML*, 2015.
- [14] O. Irsoy and E. Alpaydin, "Autoencoder Trees," in *NIPS Deep Learning Workshop*, 2014.
- [15] K. Jarrett, K. Kavukcuoglu, M. Ranzato, and Y. LeCun, "What is the best multi-stage architecture for object recognition?" in *ICCV*, 2009.
- [16] Y. Jia, C. Huang, and T. Darrell, "Beyond spatial pyramids," in *CVPR*, 2012.
- [17] P. Kotschieder, M. Fiterau, A. Criminisi, and S. Rota Bulo, "Deep neural decision forests," in *ICCV*, 2015.
- [18] A. Krizhevsky, "Learning Multiple Layers of Features from Tiny Images," *CS Dept., U Toronto, Tech. Rep.*, 2009.
- [19] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "ImageNet Classification with Deep Convolutional Neural Networks," in *NIPS*, 2012.
- [20] Y. LeCun, B. Boser, J. S. Denker, D. Henderson, R. Howard, W. Hubbard, and L. Jackel, "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, no. 4, pp. 541–551, 1989.
- [21] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [22] Y. LeCun and C. Cortes, "The MNIST database of handwritten digits," 1998.
- [23] C.-Y. Lee, P. W. Gallagher, and Z. Tu, "Generalizing pooling functions in convolutional neural networks: Mixed, gated, and tree," in *AISTATS*, 2016.
- [24] C.-Y. Lee, S. Xie, P. Gallagher, Z. Zhang, and Z. Tu, "Deeply-Supervised Nets," in *AISTATS*, 2015.
- [25] W. Li, M. Zeiler, S. Zhang, Y. LeCun, and R. Fergus, "Regularization of NNs using DropConnect," in *ICML*, 2013.
- [26] M. Liang and X. Hu, "Recurrent CNNs for Object Recognition," in *CVPR*, 2015.
- [27] M. Lin, Q. Chen, and S. Yan, "Network in network," in *ICLR*, 2013.
- [28] J. Minker, *Logic-Based Artificial Intelligence*. Springer Science & Business Media, 2000, vol. 597.
- [29] D. Mishkin, N. Sergievskiy, and J. Matas, "Systematic evaluation of cnn advances on the imagenet," *arXiv preprint arXiv:1606.02228*, 2016.
- [30] Y. Netzer, T. Wang, A. Coates, A. Bissacco, B. Wu, and A. Y. Ng, "Reading Digits in Natural Images with Unsupervised Feature Learning," in *NIPS Workshop on Deep Learning and Unsupervised Feature Learning*, 2011.
- [31] J. R. Quinlan, "C4. 5: Programming for machine learning," *Morgan Kaufmann*, vol. 38, 1993.
- [32] M. Ranzato, Y.-L. Boureau, and Y. LeCun, "Sparse Feature Learning for Deep Belief Networks," in *NIPS*, 2007.
- [33] A. Romero, N. Ballas, S. E. Kahou, A. Chassang, C. Gatta, and Y. Bengio, "FitNets: Hints for Thin Deep Nets," in *ICLR*, 2015.
- [34] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *IJCV*, vol. 115, no. 3, pp. 211–252, 2015.
- [35] D. Scherer, A. Müller, and S. Behnke, "Evaluation of pooling operations in convolutional architectures for object recognition," *ICANN*, pp. 92–101, 2010.
- [36] T. Serre, L. Wolf, S. Bileschi, M. Riesenhuber, and T. Poggio, "Robust object recognition with cortex-like mechanisms," *PAMI*, vol. 29, no. 3, 2007.
- [37] K. Simonyan and A. Zisserman, "Very deep convolutional networks for large-scale image recognition," in *ICLR*, 2015.
- [38] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, "Striving for Simplicity," in *ICLR*, 2015.
- [39] J. T. Springenberg and M. Riedmiller, "Improving deep neural networks with probabilistic maxout units," in *ICLR*, 2014.
- [40] N. Srivastava and R. R. Salakhutdinov, "Discriminative transfer learning with tree-based priors," in *NIPS*, 2013.
- [41] R. K. Srivastava, K. Greff, and J. Schmidhuber, "Highway networks," in *ICML Deep Learning Workshop*, 2015.
- [42] M. Stollenga, J. Masci, F. J. Gomez, and J. Schmidhuber, "Deep Networks with Internal Selective Attention through Feedback Connections," in *NIPS*, 2014.
- [43] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich, "Going deeper with convolutions," in *CVPR*, 2015.
- [44] L. Van der Maaten and G. Hinton, "Visualizing data using t-SNE," *JMLR*, vol. 9, no. Nov, pp. 2579–2605, 2008.
- [45] J. Wang, Z. Wei, T. Zhang, and W. Zeng, "Deeply-fused nets," *arXiv preprint arXiv:1605.07716*, 2016.
- [46] M. D. Zeiler and R. Fergus, "Stochastic Pooling for Regularization of Deep Convolutional Neural Networks," *arXiv preprint arXiv:1301.3557*, 2013.



Chen-Yu Lee Chen-Yu Lee received the Bachelor and Master degrees from National Chiao Tung University. He received the PhD degree from the University of California, San Diego (UCSD). He is now researcher at Magic Leap, Inc. His research interests are in computer vision, machine learning, and deep learning.



Patrick Gallagher Patrick Gallagher received the PhD degree in Cognitive Science from the University of California, San Diego in 2014, where he also spent a year as a post-doctoral researcher. His research interests include machine learning, deep learning, and mathematical optimization.



Zhuowen Tu Zhuowen Tu received the PhD degree in computer science from Ohio State University. He received the BE degree from Beijing Information Technology Institute and the ME degree from Tsinghua University. He is an associate professor of cognitive science at University of California, San Diego (UCSD). His main research interests include computer vision, machine learning, and neural computation.