# Why $M$ Heads are Better than One:
# Training a Diverse Ensemble of Deep Networks

Stefan Lee[1]    Senthil Purushwalkam[2]    Michael Cogswell[3]    David Crandall[1]    Dhruv Batra[3]

Indiana University[1]        Carnegie Mellon University[2]        Virginia Tech[3]

{steflee, djcran}@indiana.edu    spurushw@andrew.cmu.edu    {cogswell, dbatra}@vt.edu

## Abstract

*Convolutional Neural Networks have achieved state-of-the-art performance on a wide range of tasks. Most benchmarks are led by ensembles of these powerful learners, but ensembling is typically treated as a post-hoc procedure implemented by averaging independently trained models with model variation induced by bagging or random initialization. In this paper, we rigorously treat ensembling as a first-class problem to explicitly address the question: what are the best strategies to create an ensemble? We first compare a large number of ensembling strategies, and then propose and evaluate novel strategies, such as parameter sharing (through a new family of models we call TreeNets) as well as training under ensemble-aware and diversity-encouraging losses. We demonstrate that TreeNets can improve ensemble performance and that diverse ensembles can be trained end-to-end under a unified loss, achieving significantly higher "oracle" accuracies than classical ensembles.*

## 1. Introduction

Convolutional Neural Networks (CNNs) have shown impressive performance on a wide range of computer vision tasks. An important (and perhaps under-acknowledged) fact is that the state-of-the-art models are generally *ensembles* of CNNs, including nearly all of the top performers on the ImageNet Large Scale Visual Recognition Challenge [31]. For example, GoogLeNet [36], one of the best-performing models submitted to the ILSVRC challenge, is an ensemble achieving a five percentage point increase in accuracy over a single base model of the same architecture.

In these ensembles, multiple classifiers are trained to perform the same task and their predictions are averaged to generate a new, typically more accurate, prediction. A number of related justifications have been given for the success of ensembles, including:

i) Bayesian Model Averaging, that ensembles are a finite sample approximation to integration over the model class [9, 26, 27];

ii) Model Combination, that ensembles enrich the space of hypotheses considered by the base model class and are representationally richer [8]; and

iii) Reducing Estimation and Optimization Errors, that ensemble averaging reduces the variance of base models, averaging out variations due to objective function non-convexity, initialization, and stochastic learning [7, 28].

At the heart of these arguments is the idea of *diversity*: if we train multiple learners with decorrelated errors, their predictions can be averaged to improve performance [5]. In this work, we rigorously treat ensembling as a problem in its own right, examining multiple ensembling strategies ranging from standard bagging to parameter sharing and ensemble-aware losses. We compare these methods across multiple datasets and architectures, demonstrating that some standard techniques may not be suitable for deep ensembles and novel approaches improve performance.

***Ensemble-Aware Losses.*** Typically, ensemble members are trained independently with no unifying loss, despite the fact that outputs are combined at test time. It is common in classical literature to view ensemble members as "experts" [18] or "specialists" [17], but in typical practice no effort is made to encourage diversity or specialization. It seems natural then to question whether a ensemble-aware loss might result in better performance. Here we study two ensemble-aware losses: (1) directly training an ensemble to minimize the loss of the ensemble mean, and (2) generalization of Multiple Choice Learning [13] to explicitly encourage diversity.

***Parameter Sharing.*** As a number of papers have demonstrated, initial layers of CNNs tend to learn simple, generic features which vary little between models, while deeper layers learn features specific to a particular task and input distribution [12, 24, 37]. We propose a family of tree-structured deep networks (which we call TreeNets) that exploit the generality in lower layers by sharing them across ensemble members to reduce parameters. We investigate the depth at which sharing should happen, along a spectrum from single models (full sharing) to independent ensembles (no sharing). This coupling of lower layers naturally forces any di-

versity between ensemble members to be concentrated in the deeper, unshared layers. Perhaps somewhat surprisingly, we find that the optimal setting is *not* a classical ensemble, but instead a TreeNet that shares a few (typically 1-2) initial layers. Thus tree-structured networks are a simple way to improve performance while reducing parameters.

***Model-Distributed Training of Coupled Ensembles.*** Unfortunately, both of the above approaches to coupling ensemble members, either at the "top" of the architecture with ensemble-aware losses that operate on outputs from all ensemble members, or at the "bottom" with parameter sharing in TreeNets, create significant computational difficulties. Since networks are not independent, it is no longer possible to train them separately in parallel, and sequential training may require months of GPU time even for relatively small ensembles. Moreover, even if training time is not a concern, larger models often take up most of the available RAM on a GPU, so it is not possible to fit an ensemble on one GPU. To overcome these hurdles, we present and will release a novel MPI-based model-parallel distributed modification to the popular Caffe deep learning framework [19] that implements cross-process communication as layers of the CNN.

We thoroughly evaluate each methodology across multiple datasets and network architectures. These experiments cast new light on ensembles in deep networks, demonstrating the effects of randomization in parameter and data space, parameter sharing, and unified losses on modern scale vision problems. More concretely, we:

  i) rigorously treat CNN ensembling as its own problem,
  ii) introduce a family of models called TreeNets that permit a spectrum of degrees of layer-sharing,
  iii) present ensemble-aware and diversity-encouraging loses, and
  iv) present a distributed model-parallel framework to train deep ensembles.

## 2. Related Work

Neural networks, ensembles, and techniques to improve robustness and diversity of grouped learners have decades of work in machine learning research, but only recently have ensembles of CNNs been studied. Related work can be broadly divided into two categories: ensemble learning for general networks, and its more recent application to CNNs.

***Ensemble Learning Theory.*** Neural networks have been applied in a wide variety settings with many diverse modifications. Much of the theoretical foundation for ensemble learning with neural networks was laid in the 1990s. Krogh *et al.* [22] and Hansen and Salamon [16] provided theoretical and empirical evidence that diversity in error distributions across member models can boost ensemble performance. This led to ensemble methods that averaged predic-

tions from models trained with different initializations [16] and from models trained on different bootstrapped training sets [22, 38]. These methods take an indirect approach to introducing diversity in ensembles. Other work has explicitly trained decorrelated ensembles of neural networks by penalizing positive correlation between error distributions [2, 23, 29]. While effective on shallow networks, these methods have not been applied to deeper architectures.

Although initially proposed for Structured SVMs, the work of Guzman-Rivera et al. [13–15] on Multiple Choice Learning (MCL) provides an attractive alternative that does not require computing correlation between error. Related ideas were studied by Dey et al. [6] in the context of submodular list prediction. We generalize MCL and apply it to CNNs – incorporate it with stochastic gradient descent-based training.

***CNN Ensembles.*** While ensembles of CNNs have been used extensively, little work has focused on improving the ensembling process. Most CNN ensembles use multiple random initializations or training data subsets to inject diversity. For example, popular ensembles of VGG [33] and AlexNet [21] simply retrain with different initializations and average the predictions. GoogLeNet [36] induces diversity with straightforward bagging, training each model with a sampled dataset. Other networks, like Sequence to Sequence RNNs [35], use both approaches simultaneously.

Parameter sharing is not a novel development in CNNs, but its effect on ensembles has not been studied. Recent related work by Bachman *et al.* [3] proposed a general framework called pseudo-ensembles for training robust models. They define a pseudo-ensemble as a group of child models which are instances of a parent model perturbed by some noise process. They explicitly encourage correlation in model parameters through the parent by a Pseudo-Ensemble Agreement (PEA) regularizer. Although outwardly related to parameter sharing, pseudo-ensembles are fundamentally different than the techniques presented here, as they use parameter sharing to train a single robust CNN model rather than to produce an ensemble with fewer parameters. Other recent work by Sercu et al. [32] uses parameter sharing in the context of multi-task learning to build a common representation for multilingual translation. Finally, Dropout [17] can be interpreted as a procedure that trains an exponential number of highly related networks and cheaply combines them into one network, similar to PEA

One relevant recent work is [17], which briefly focuses on ensembles. Members of this types of ensemble are specialists which are trained on subsets of all possible labels with each subset manually designed to include easily confused labels. These models are fine-tuned from one shared generalist and then combined to make a final prediction. In contrast, our diversity-encouraging loss require no human hand-designing of class specialization – our loss naturally

allows members to specialize according to subset of classes or pockets of feature space, providing an end-to-end way of learning diverse ensembles.

## 3. Experimental Design

We first describe the datasets, architectures, and evaluation metrics that we use in our experiments to better understand ensembling in deep networks.

### 3.1. Datasets and Architectures

We evaluate on three popular image classification benchmarks: CIFAR10 [20], CIFAR100 [20], and the 2012 ImageNet Large Scale Visual Recognition Challenge (ILSVRC) [30]. Since our goal is not to present new designs and architectures for these tasks but rather to study the effect of different ensembling techniques, unless otherwise noted we use standard models and training routines. All models are trained via stochastic gradient descent with momentum and weight decay.

***CIFAR10.*** For this dataset, we use Caffe "CIFAR10 Quick" [19] network as our base model. The reference model is trained using a batch size of 350 for 5,000 iterations with a momentum of 0.9, weight decay of 0.004, and an initial learning rate of 0.001 which drops to 0.0001 after 4000 iterations. We refer to this network and training procedure as CIFAR10-Quick.

***CIFAR100.*** We use the Network in Network model by Lin et al. [25] as well as their reference training procedure, which runs for 300,000 iterations with a batch size of 128, momentum of 0.9, weight decay of 0.0001, and an initial learning rate of 0.1. The learning rate decays by a factor of 10 whenever the training loss fails to drop by at least 1% over 20,000 iterations; this occurs twice over the course of a typical training run. Our reference model's accuracy is about 4% lower than reported in [25], because we do not perform their dataset normalization procedure. We refer to this network and training procedure as CIFAR100-NiN.

***ILSVRC2012.*** For this dataset we use both the Network in Network model [25] and CaffeNet (similar to AlexNet [21]). Both networks are trained for 450,000 iterations with an initial learning rate of 0.1, momentum of 0.9, and weight decay of 0.0005. For NiN the batch size is 128 and the learning rate is reduced by a factor of 10 every 200,000 iterations. For CaffeNet the batch size is 256 and the learning rate schedule is accelerated, reducing every 100,000 iterations. We refer to these models as ILSVRC-NiN and ILSVRC-Alex, respectively.

### 3.2. Evaluation Metrics

We evaluate our ensemble performance with respect to two different metrics. **Ensemble-Mean Accuracy** is the accuracy of the "standard" test-time procedure for ensembles – averaging the beliefs of all members and predicting the most confident class. Strong performance on this metric indicates that the ensemble members generally agree on the correct response, with errors reduced by smoothing across members. In contrast, **Oracle Accuracy** is the accuracy of the ensemble if an "oracle" selects the prediction of the *most accurate* ensemble member for each example. Oracle Accuracy demonstrates what the ensemble knows as a collection of specialists, and has been used in prior work to measure ensemble performance [4, 6, 13–15].

## 4. Random Initialization and Bagging

We now present our analysis of different approaches to training CNN ensembles. This section focuses on standard approaches, while Sections 5 and 6 present novel ideas on parameter sharing and ensemble-aware losses.

Randomly initializing network weights and randomly resampling dataset subsets (bagging) are perhaps the most commonly-used methods to create model variation in members of CNN ensembles. Table 1 presents results using three different ensembling techniques: (1) Random Initialization, in which all member models see the same training data but are initialized using different random seeds, (2)Bagging, in which each member uses the same initial weights but trains on a subset of data sampled (with replacement) from the original, and (3) Combined, which uses both techniques. Numbers in the table are accuracies and standard deviations across three trials. The CIFAR ensembles were built with four members while the ILSVRC ensembles had five.

As expected, all ensembles improve performance over their single base model. Somewhat surprisingly, we find that bagging reduces Ensemble-Mean Accuracy compared to random initialization alone, while Oracle Accuracy remains nearly constant. This result suggests that the bagged networks are poorly calibrated, such that confident incorrect responses are negatively impacting results. The individual member networks (not shown in table) also perform worse than those trained on the original dataset. We attribute these

| | Single Model | Random Init. | Bagging | Combined |
|---|---|---|---|---|
| | Accuracy | Ensemble-Mean Accuracy | | |
| **CIFAR10-Quick** ×4 | 77.06 ±0.27 | 80.72 ±0.10 | 78.40 ±0.28 | 78.95 ±0.17 |
| **CIFAR100-NiN** ×4 | 60.19 ±0.49 | 66.51 ±0.27 | 62.11 ±0.24 | 61.73 ±0.16 |
| **ILSVRC-Alex** ×5 | 56.79 ±0.04 | 59.94 ±0.36 | 57.46 ±0.12 | 57.39 ±0.14 |
| **ILSVRC-NiN** ×5 | 58.90 ±0.13 | 64.08 ±0.11 | 55.02 ±0.15 | 60.51 ±0.12 |
| | Accuracy | Oracle Accuracy | | |
| **CIFAR10-Quick** ×4 | 77.06 ±0.27 | 89.89 ±0.17 | 89.94 ±0.27 | 89.28 ±0.25 |
| **CIFAR100-NiN** ×4 | 60.19 ±0.49 | 78.63 ±0.31 | 75.47 ±0.12 | 75.21 ±0.28 |
| **ILSVRC-Alex** ×5 | 56.79 ±0.04 | 70.45 ±0.63 | 69.58 ±0.17 | 69.61 ±0.17 |
| **ILSVRC-NiN** ×5 | 58.90 ±0.13 | 73.60 ±0.07 | 67.79 ±0.02 | 72.92 ±0.00 |

Table 1: Comparison of standard ensembling techniques. All ensembles outperform their base models, but bagging shows smaller gains resulting from reduced training data.

results to the reduction in unique training exemplars that bagging introduces. Given an initial dataset of $M$ examples from which we draw $M$ points with replacement to make a bagged set $B$, the probability of an example $X_i$ being sampled at least once is $P(X_i \in B) = 1 - \left(\frac{(M-1)}{M}\right)^M$. The expected fraction of examples drawn at least once is thus $1 - \left(1 - \frac{1}{M}\right)^M$, which is approximately $1 - 1/e \approx 0.63$ for large $M$; i.e. bagging costs over a third of our unique data points! Not only are we losing 37% of our data, we are also introducing that many duplicated data points. To examine whether these duplicates affect performance, we reran the CIFAR10 experiments with a dataset of 31,500 unique examples (approximately 63% of the original dataset) and found similar reductions in accuracy, indicating that the loss of unique data is the primary negative effect of bagging.

Note that for convex or shallow models, the loss of unique exemplars in bagging is typically acceptable as random parameter initialization is simply insufficient to produce diversity. To the best of our knowledge, this is the first finding to establish **that random initialization may not only be sufficient but** *preferred* **over bagging for deep networks given their large parameter space and the necessity of large training data**.

## 5. Parameter Sharing with TreeNets

Ensembles and single models can be seen as two endpoints on a spectrum of approaches: single models require a careful allocation of parameters to perform well, while ensembles extract as much performance as possible from multiple instances of a base model. Ensemble approaches likely introduce wasteful duplication of parameters in generic lower layers, increasing training time and model size. The hierarchical nature of CNNs makes them well-suited to alternative ensembling approaches where member models benefit from shared information at the lower layers while retaining the advantages of classical ensembling methods.

Motivated by this observation, in this section we present and evaluate a family of tree-structured CNN ensembles called TreeNets, as shown in Figure 1. A TreeNet is an ensemble consisting of zero or more shared initial layers, followed by a branching point and zero or more independent layers. During training, the shared layers above a branch receive gradient information from each child network, which are accumulated according to back-propagation. At test time, each path from root to leaf can be considered an independent network, except that redundant computations at the shared layers need not be performed.

We evaluated our novel TreeNet models on the two larger architectures trained on ImageNet, ILSVRC-Alex and ILSVRC-NiN, and Table 2 presents the results. The table shows the Ensemble-Mean Accuracy (again in terms of means and standard deviations across three trials) achieved
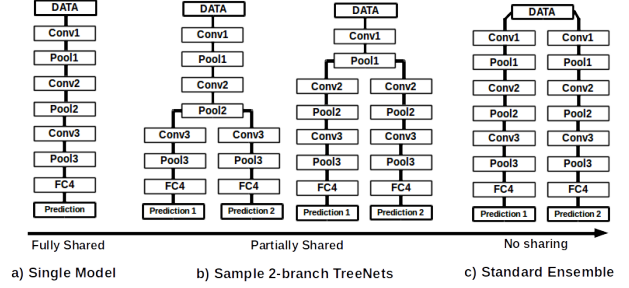


Figure 1: TreeNets exist on a spectrum between single models and fully independent ensembles.

| ILSVRC-Alex ×5 Split Point | Ensemble-Mean Accuracy | ILSVRC-NiN ×4 Split Point | Ensemble-Mean Accuracy |
|---|---|---|---|
| ensemble | 59.47 ±0.45 | ensemble | 64.08 ±0.00 |
| **conv1** | **59.62 ±0.09** | conv1 | 65.50 ±0.24 |
| conv2 | 59.32 ±0.17 | **cccp1** | **65.69 ±0.08** |
| conv3 | 58.39 ±0.10 | cccp2 | 65.64 ±0.11 |
| conv4 | 57.73 ±0.05 | conv2 | 65.64 ±0.07 |
| conv5 | 55.25 ±0.03 | cccp3 | 65.47 ±0.07 |
| single model | 56.79 ±0.04 | cccp4 | 65.62 ±0.01 |
|  |  | single model | 58.90 ±0.13 |

Table 2: Results for TreeNet training at various depths of ILSVRC-Alex and ILSVRC-NiN. Ensemble performance is retained even with substantial parameter sharing.

by TreeNets with splits at different depths. For example, splitting at conv2 means that all layers up to and including conv2 are shared, and all branches are independent afterwards. Since layers that do not contain any parameters (e.g. pooling, nonlinearity) are unaffected by parameter sharing, we only show results for splitting on parameterized layers.

We see that **shared parameter networks not only retain the performance of full ensembles, but can outperform them**. For our best ILSVRC-NiN TreeNet, we improve accuracy over standard ensembles while reducing the parameter count by 7%. It may be that lower layer representations, though simple and generic, still had room for improvement. By sharing low level weights, each weight is updated by multiple sources of supervision, one per branch. This indicates TreeNets could provide regularization which favors slightly better low level representations.

We find further evidence for this claim by looking at individual branches of the TreeNet compared to the independently trained networks of the ensemble. Regardless of split point, each TreeNet branch in our shared ensemble achieved around 2 to 3 percentage points higher accuracy than independent ensemble members. Unlike in classical ensembles where each member model performs about as well as the base architecture, TreeNets seem to boost performance of not only the ensemble but the individual networks as well. We also experimented with multiple splits leading to more complicated "balanced binary" tree structures on ILSVRC-NiN and found similar improvements.

We also tested ILSVRC-Alex TreeNet models trained for

| | Independent Losses | Score-Averaged | Prob-Averaged |
|---|---|---|---|
| | Ensemble-Mean | Accuracy | |
| **CIFAR10-Quick** $\times 4$ | 80.72 ±0.10 | 79.32 ±0.02 | 77.10 ±0.16 |
| **CIFAR100-NiN** $\times 4$ | 66.51 ±0.27 | 65.77 ±0.21 | 62.77 ±0.28 |
| **ILSVRC-Alex** $\times 5$ | 59.94 ±0.13 | 56.56 ±0.10 | 49.81 ±0.18 |
| **ILSVRC-NiN** $\times 5$ | 83.43 ±0.10 | 79.24 ±0.36 | 42.39 ±0.24 |

Table 3: Results of training ensembles to reduce loss over member predictions averaged either over scores or probabilities.

object detection on PASCAL VOC 2007 [10] dataset. We used the Fast R-CNN [12] architecture fine-tuned from our TreeNet models. For the test-time bounding-box regression, we average the results from each member model for an ensemble. We found a statistically significant increase in mean average precision of about 0.7% across multiple runs compared to starting from a standard ensemble. We take these initial experiments to imply TreeNet models are at least as generalizable to other tasks as standard ensembles. More details are provided in the supplementary materials.

To summarize the key results in this section, we found that TreeNets with a few (typically 1-2) initial layers outperform classical ensembles, while also having fewer parameters which may reduce test-time computation time and memory requirements.

# 6. Training Under Ensemble-Aware Losses

In the two previous sections, each ensemble member was trained with the same objective – independent cross-entropy of each ensemble member. What happens if the objective is aware of the ensemble? We begin by showing a surprising result: the first "natural" idea of simply optimizing the performance of the average-beliefs of the ensemble does *not* work, and we provide intuitions why this is the case (lack of diversity). This negative result shows that a more careful design for ensemble-aware loss functions is crucial. We then propose a diversity-encouraging loss function that shows significantly improved oracle performance.

## 6.1. Directly Optimizing for Model Averaging

For a standard ensemble, test-time classification is typically performed by averaging the output of the member networks, so it is natural to explicitly optimize the performance of the corresponding Ensemble-Mean loss during training. We ran all four ensemble architectures under two settings: (1) Score-Averaged, in which we average the last layer outputs (i.e. the scores that are *inputs* to the softmax function), and (2) Probability-Averaged, in which we average the softmax probabilities of ensemble members. Intuitively, the difference between the two settings is that the former assumes the ensemble members are "calibrated" to produce scores of similar relative magnitudes while the latter does not.

Table 3 shows the results of these experiments, again averaged over three trials. In all cases, network averaging

*reduced* performance, with Probability-Averaged causing greater degradation. This is counter-intuitive: explicitly optimizing for the performance of Ensemble-Mean does *worse* than averaging independently trained models. We attribute this to two problems, which we now discuss: lack of diversity and numerical instability.

***Averaging Outputs Reduces Diversity.*** Unfortunately, averaging scores or probabilities during training has the unintended consequence of eliminating diversity in gradients back-propagated through the ensemble. Consider a generic averaging layer,

$$\mu(\boldsymbol{x_1}, ..., \boldsymbol{x_N}) = \frac{1}{N} \sum_{i=1}^{N} \boldsymbol{x_i},$$

that ultimately contributes to some loss $\ell$, and consider the derivative of $\ell$ with respect to some $\boldsymbol{x_i}$,

$$\frac{\partial \ell}{\partial \boldsymbol{x_i}} = \frac{\partial \ell}{\partial \mu} \frac{\partial \mu}{\partial \boldsymbol{x_i}} = \frac{\partial \ell}{\partial \mu} \frac{1}{N}.$$

This expression does not depend on $i$ — gradients back-propagated into all ensemble members are identical! Due to the averaging layer, responsibility for mistakes is shared, which eliminates gradient diversity. This is different from the behavior of an ensemble of independently trained networks, where each member receives a different gradient depending on individual performance. (The averaging also scales the gradients, so in our experiments we compensate by increasing the learning rate by a factor of $N$; otherwise, we found learning tended to arrive at even worse solutions.)

***Averaging Probabilities Is Unstable.*** We attribute the further loss of accuracy when averaging probabilities (versus scores) to increased numerical instability. The softmax function's derivative with respect to its input is unstable for outputs near 0 or 1. However, when paired with a cross-entropy loss, the derivative of the loss with respect to softmax input reduces to a simple subtraction. Unfortunately, there is no similar simplification for cross-entropy over an average of softmax outputs (see supplemental materials for details). Optimization under these conditions is difficult, causing loss at convergence for Probability-Averaged networks to be nearly twice that of Score-Averaged networks, and about the same as a single network.

Motivated by the finding that decreased diversity from optimizing Ensemble-Mean leads to reduced performance, we next present an explicit diversity-encouraging loss.

## 6.2. Adding Diversity via Multiple Choice Learning

We have so far discussed the role of ensemble diversity in the context of model averaging; however, in many settings, generating multiple plausible hypotheses may be preferred to producing a single answer. Ensembles fit naturally into this space as they produce multiple answers by design.

However, independently trained models typically converge to similar solutions, prompting the need to optimize for diversity directly. In this section, we develop and experiment with diversity encouraging losses and demonstrate their effectiveness at specializing ensembles.

We build on *Multiple Choice Learning* (MCL) [13], which we briefly recap here. Consider a set of predictors $\{\theta_1, ..., \theta_M\}$ such that $\theta_m : x \rightarrow P$ where $P$ is a probability distribution over some set of labels, and a dataset $D=\{(x_1, y_1), ..., (x_N, y_N)\}$, where each feature vector $x_i$ has a ground truth label $y_i$. From the point of view of an oracle that only listens to the most correct $\theta_m$, the loss for an example $(x, y)$ is

$$\mathcal{L}_{set}(x, y) = \min_{m \in [1, M]} \ell\left(\theta_m(x), y\right),$$

which we will call the oracle set-loss. Intuitively, given that the oracle will select the most correct predictor, the loss on any example is the minimum loss over predictors. Alternatively, the oracle loss can be interpreted as allowing a system to guess $M$ times, scoring an example as correct if any guess is correct. Thus an ensemble of $M$ predictors is directly comparable to the commonly used top-$M$ metric used in many benchmarks (e.g. top-5 in ILSVRC [30]).

We adapt this framework to the cross-entropy loss used for training deep classification networks. Given a single predictor $\theta_m$, the cross-entropy loss for example $(x, y)$ is

$$\ell(x, y) = -\log\left(p_y^{\theta_m}\right),$$

where $p_y^{\theta_m}$ is the predicted probability of class $y$. Let $\alpha_{mi}$ be a binary variable indicating whether predictor $\theta_m$ has the lowest loss on example $(x_i, y_i)$. We can then define a cross-entropy oracle set-loss over a dataset $D$,

$$\mathcal{L}_{set}(D) = \frac{1}{|D|} \sum_{(x_i, y_i) \in D} \sum_{m=1}^{M} -\alpha_{mi} \log\left(p_{y_i}^{\theta_m}\right).$$

Notice that just like cross-entropy is an upper-bound on training error, this expression is an upper-bound on the oracle training error [13]. Guzman-Rivera et al. [13] presented a coordinate descent algorithm for optimizing such an objective. Their approach alternates between two stages: first, each data point is assigned to its most accurate predictors, and then models are trained until convergence using only the assigned examples.

Even if done in parallel, training multiple CNNs to convergence for each iteration is intractable. We thus interleave the assignment step with batch updates in stochastic gradient descent. For each batch, we pass the examples through the network, producing probability distributions over the label space from each ensemble member. During the backward pass, the gradient of the loss for each example is computed with respect only to the predictor with the lowest error

| CIFAR10-Quick×4 | Member Networks | | | | Ensemble Accuracy | |
|---|---|---|---|---|---|---|
| k | | Accuracy | | | Ensemble-Mean | Oracle |
| 1 | 24.35 | 27.18 | 27.15 | 19.36 | 28.38 | 93.10 |
| 2 | 50.95 | 27.41 | 55.88 | 34.81 | 75.16 | 92.78 |
| 3 | 65.46 | 40.71 | 64.79 | 70.37 | 79.76 | 92.55 |
| 4 | 77.12 | 76.76 | 77.29 | 76.80 | 80.72 | 89.78 |

Table 4: Increasing the number of predictors each data point is assigned to results in reduced oracle accuracy as the diversifying effect is reduced. Note that $k = 4$ is a standard ensemble.

on that example (with ties broken randomly). Pseudo-code is available in the supplement.

So far we have assumed that the oracle can select only one answer, i.e. $\sum_m \alpha_{mi} = 1$, however this can easily be generalized to select the $k$ predictors with lowest loss such that $\sum_m \alpha_{mi} = k$. Varying $k$ from one to the number of predictors trades off between diversity and the number of training examples each predictor sees, which affects both generalization and convergence.

*Experimental results.* We begin our experiments with MCL on the CIFAR10-Quick network. Table 4 shows the individual network accuracies and the oracle accuracy for MCL trained ensembles of different values of $k$. As $k$ is increased, each member network is exposed to more of the data and we see decreased oracle accuracy in exchange for increased individual member performance. At $k$=4, the oracle-set loss reduces to independent cross-entropy for each member, producing a standard ensemble. The $k$=1 case showcases the degree of model specialization. **Each individual network performs *very* poorly (accuracy of 19-27%); however, taken as an ensemble the oracle accuracy is over 93%!** This clearly shows that the networks have specialized and diversified with each taking responsibility for a subset of examples. To the best of our knowledge, this is the first work to demonstrate such behavior.

To further characterize what the MCL member networks are learning, we tracked which test examples are assigned to each ensemble member by the oracle accuracy metric (i.e. which ensemble member has the lowest error on each example). Figure 2(a)-(d) show the distribution of classes assigned to each ensemble member, and the results are striking: at $k$=1 we see almost complete division of the label space! As $k$ increases we see increased uniformity in these distributions. Note that these divisions emerge from the loss and are not hand-designed or pre-initialized in any way.

In Figure 2(e) we visualize how the ensemble members respond to input images using guided backprop [34], which is similar to the *deconv* visualizations of Zeiler and Fergus [37]. These images can be interpreted as the gradient of the indicated class score with respect to the input image. Features that are clear in these images have the largest influence on the network's output for a given input image. Each row shows these visualizations for a single input image for a standard network and for members of an MCL ensemble.
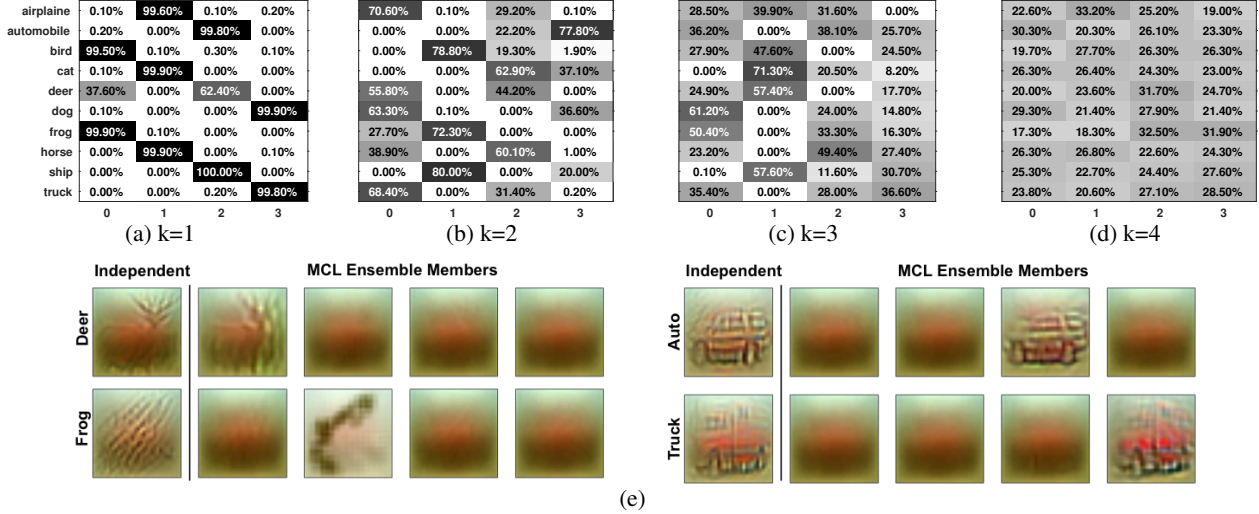
| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| airplane | 0.10% | 99.60% | 0.10% | 0.20% |
| automobile | 0.20% | 0.00% | 99.80% | 0.00% |
| bird | 99.50% | 0.10% | 0.30% | 0.10% |
| cat | 0.10% | 99.90% | 0.00% | 0.00% |
| deer | 37.60% | 0.00% | 62.40% | 0.00% |
| dog | 0.10% | 0.00% | 0.00% | 99.90% |
| frog | 99.90% | 0.10% | 0.00% | 0.00% |
| horse | 0.00% | 99.90% | 0.00% | 0.10% |
| ship | 0.00% | 0.00% | 100.00% | 0.00% |
| truck | 0.00% | 0.00% | 0.20% | 99.80% |

(a) k=1

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| airplane | 70.60% | 0.10% | 29.20% | 0.10% |
| automobile | 0.00% | 0.00% | 22.20% | 77.80% |
| bird | 0.00% | 78.80% | 19.30% | 1.90% |
| cat | 0.00% | 0.00% | 62.90% | 37.10% |
| deer | 55.80% | 0.00% | 44.20% | 0.00% |
| dog | 63.30% | 0.10% | 0.00% | 36.60% |
| frog | 27.70% | 72.30% | 0.00% | 0.00% |
| horse | 38.90% | 0.00% | 60.10% | 1.00% |
| ship | 0.00% | 80.00% | 0.00% | 20.00% |
| truck | 68.40% | 0.00% | 31.40% | 0.20% |

(b) k=2

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| airplane | 28.50% | 39.90% | 31.60% | 0.00% |
| automobile | 36.20% | 0.00% | 38.10% | 25.70% |
| bird | 27.90% | 47.60% | 0.00% | 24.50% |
| cat | 0.00% | 71.30% | 20.50% | 8.20% |
| deer | 24.90% | 57.40% | 0.00% | 17.70% |
| dog | 61.20% | 0.00% | 24.00% | 14.80% |
| frog | 50.40% | 0.00% | 33.30% | 16.30% |
| horse | 23.20% | 0.00% | 49.40% | 27.40% |
| ship | 0.10% | 57.60% | 11.60% | 30.70% |
| truck | 35.40% | 0.00% | 28.00% | 36.60% |

(c) k=3

| | 0 | 1 | 2 | 3 |
|---|---|---|---|---|
| airplane | 22.60% | 33.20% | 25.20% | 19.00% |
| automobile | 30.30% | 20.30% | 26.10% | 23.30% |
| bird | 19.70% | 27.70% | 26.30% | 26.30% |
| cat | 26.30% | 26.40% | 24.30% | 23.00% |
| deer | 20.00% | 23.60% | 31.70% | 24.70% |
| dog | 29.30% | 21.40% | 27.90% | 21.40% |
| frog | 17.30% | 18.30% | 32.50% | 31.90% |
| horse | 26.30% | 26.80% | 22.60% | 24.30% |
| ship | 25.30% | 22.70% | 24.40% | 27.60% |
| truck | 23.80% | 20.60% | 27.10% | 28.50% |

(d) k=4

(e)

Figure 2: *(a)-(d):* Percentage of test examples of each class assigned to each ensemble member by the oracle (i.e. those with lowest loss). The degree of specialization is very sharp at $k = 1$ and softens to almost uniform at $k = 4$. *(e):* Guided-backprop images for standard and MCL $k = 1$ trained ensemble members. Networks that have not specialized in a given class are agnostic to the image content.

Networks that have not specialized in the given class are agnostic to the image content. See supplementary material for more examples.

***MCL As Label-Space Clustering.*** We have shown that $k = 1$ MCL trained ensembles tend to converge to a label-space clustering where each member focuses on a subset of the labels. The set of possible label-space clusterings is vast, so to put the MCL results into perspective we train hand-designed specialist ensembles with randomized label assignments. For CIFAR10 we randomly split the labels evenly to the four ensemble members and train each with respect to those labels. Over the course of 100 trials, we found oracle-accuracy ranged from 87.62 to 94.65 with a mean of 91.83. This shows that generally the MCL optimization selects high quality label space clusterings with respect to oracle accuracy.

An alternative strategy presented by [17] is to diversify members by dividing labels into clusters of hard to distinguish classes; very briefly described, assignments are generated by clustering the covariance matrix of label scores computed across an input set for a generalist CNN. We trained an ensemble using this clustering method and it led to significantly decreased oracle performance versus MCL on CIFAR10-Quick and ILSVRC-Alex. This is not surprising since they do not optimize for oracle accuracy.

***Overcoming Data Fragmentation.*** Despite not training member networks to convergence in each iteration of co-ordinate descent, our method results in improved oracle accuracy over standard ensembles. However, interleaving the assignment step with stochastic gradient descent results in data fragmentation, with each network seeing only a fraction of each batch (as illustrated by the class-specialization). We find this reduced effective batch size results in noisy gra-

| CIFAR10-Quick×4 | Iterations of Cross-Entropy Pretraining | | | | |
|---|---|---|---|---|---|
| k | 100 | 500 | 1000 | 2000 | 4000 |
| 1 | 94.21 | 94.21 | 94.65 | 95.75 | 96.00 |
| 2 | 92.79 | 93.06 | 93.16 | 93.07 | 93.00 |
| 3 | 92.25 | 92.93 | 91.77 | 90.94 | 90.94 |

Table 5: Increasing the amount of pretraining before fine-tuning with the MCL loss results in increase oracle accuracy.

dients that inhibit learning, especially on larger networks.

Deep networks are especially sensitive to the effects of data fragmentation early in training when errors (and therefore gradients) are typically larger. In Guzman-Rivera et al. [13], initial assignments for the first iteration of training were decided by clustering the data into $k$ clusters. In contrast, assignments in our approach are based on network performance which is initially the result of random initialization. To investigate the effect of this initial phase of learning, we applied our MCL loss to fine-tune a previously trained CIFAR10-Quick ensemble. As shown in Table 5, the benefits of pretraining are most pronounced for lower values of $k$ where data fragmentation is most severe.

While pretraining did stabilize learning, data fragmentation on CIFAR10 is a relatively minor problem whereas training with MCL from scratch on larger networks using standard batch sizes consistently failed to outperform standard ensembles. We attribute this to a combination of data fragmentation and the difficulty of initial learning. To test this hypothesis, we experimented with fine-tuning and gradient accumulation across batches on the ILSVRC-Alex architecture. We accumulated gradients from 5 batches before updating parameters and fine-tuned from a fully-trained ensemble. Table 6 shows the result of 3000 iterations of this fine-tuning experiment for different values of $k$. This setup overcame the data fragmentation problem and we see the

| ILSVRC-Alex×5 | Single Member Accuracy | Ensemble Accuracy | |
|---|---|---|---|
| | | Ensemble-Mean | Oracle |
| k=1 | 46.50 | 55.22 | 74.67 |
| k=2 | 52.48 | 59.21 | 73.40 |
| k=3 | 55.38 | 59.73 | 71.75 |
| k=4 | 56.33 | 60.09 | 70.84 |
| base ensemble | 57.17 | 60.31 | 70.50 |

Table 6: Fine-tuning and gradient accumulation across batches allows larger networks to specialized under the MCL loss.

same trends as in CIFAR10.

These experiments demonstrate MCL's ability to quickly diversify an ensemble. To push this further, we reran the fine-tuning experiment for $k$=1, this time initializing all ensemble members with the same network. Despite starting from an ensemble of identical networks with an oracle accuracy of 56.90%, the ensemble reached an oracle accuracy of 72.67% after only 3000 iterations!

We have demonstrated that the MCL loss is effective at inducing diversity, however the member networks specialize so much that Ensemble-Mean Accuracy suffers. We tried linearly combining the MCL loss with the standard cross-entropy to balance diversity with general performance. We find training under this loss improves CIFAR10 Ensemble-Mean accuracy by 1% over a standard ensemble.

In this section we have developed a novel MCL framework and shown it produces ensembles with substantially improved oracle accuracies when training from scratch and even when fine-tuning from a single network.

## 7. Distributed Ensemble Training

Training an ensemble on a single GPU is prohibitively expensive, so standard practice for large ensembles is to train the multiple networks either sequentially or in parallel. However, any form of model coupling requires communication between learners. To make enable our experiments at scale, we have developed and will release a modification to Caffe, which we call MPI-Caffe, that uses the Message Passing Interface (MPI) [1] standard to enable cross-GPU/machine communication. These communication operations are provided as Caffe model layers, allowing network designers to quickly experiment with distributed networks, where different parts of the model reside on different GPUs and machines. Figure 3 shows how an ensemble of CIFAR10-Quick networks with parameter sharing and model averaging is defined as a single specification and distributed across multiple process. In MPI-Caffe, each process is assigned a identifier (called a rank); by setting the ranks each network layer belongs to, we can easily design distributed ensembles.

The `MPIBroadcast` and `MPIGather` layers provide the core communication functionality. `MPIBroadcast` forwards its input to the other processes during a forward



Figure 3: MPI-Caffe models can be defined by a single network specification and distributed by across multiple GPUs. Dashed lines indicate cross-process communication and not input/output.

pass and accumulates gradients from each during backpropagation. The forward pass for `MPIGather` collects all of its inputs from multiple processes and outputs them to a single network, and the backward pass simply routes the gradients back to the corresponding input.

We tested our MPI-Caffe framework on a large-scale cluster with one Telsa K20 GPU per node and a maximum MPI node interconnect bandwidth of 5.8 GB/sec. To characterize the communication overhead for an ensemble, we measure the time spent sharing various layers of the ILSVRC-Alex×5 architecture. The largest layer we shared was *pool2* which amounts to broadcasting nearly 36 million floats per batch. Despite the layer's size we find only 0.49% of the forward-backward pass time is used by communication. More details are available in the supplement.

## 8. Discussion and Conclusion

There is a running theme behind all of the ideas presented in this paper: diversity. Our experiments on bagging demonstrate that the diversity induced in ensemble members by random parameter initializations is more useful than that introduced by bags with duplicated examples. Our experiment on explicitly training for Ensemble-Mean performance show that averaging beliefs of ensemble members before computing losses has the unintended effect of removing diversity in gradients. Our novel diversity-inducing MCL loss shows that encouraging diversity in ensemble members can significantly improve performance. Finally, our novel TreeNet architecture shows that diversity is important in high-level representations while low-level filters might be better off without it. Training these large-scale architectures is made practical by our MPI-Caffe framework.

In future work, we would like to adapt the MCL framework to structured predictions problems. In a structured context where the space of "good" solutions is quite large, we feel diverse models can have an even greater benefit.

# References

[1] MPI: A Message-Passing Interface Standard Version 3.1. Retrieved from http://www.mpi-forum.org/docs/mpi-3.1/mpi31-report.pdf, 2015. 8

[2] M. Alhamdoosh and D. Wang. Fast decorrelated neural network ensembles with random weights. *Information Sciences*, 264:104–117, 2014. 2

[3] P. Bachman, O. Alsharif, and D. Precup. Learning with pseudo-ensembles. In *Advances in Neural Information Processing Systems*, pages 3365–3373, 2014. 2

[4] D. Batra, P. Yadollahpour, A. Guzman-Rivera, and G. Shakhnarovich. Diverse $m$-best solutions in Markov Random Fields. In *Proceedings of the European Conference on Computer Vision*, pages 1–16. Springer, 2012. 3

[5] C. Bishop. *Pattern Recognition and Machine Learning*. Springer, 2007. 1

[6] D. Dey, T. Liu, M. Hebert, and J. A. Bagnell. Contextual sequence prediction with application to control library optimization. In *Robotics Science and Systems (RSS)*, 2012. 2, 3

[7] T. G. Dietterich. Ensemble methods in machine learning. In *Proceedings of the First International Workshop on Multiple Classifier Systems*, MCS '00, 2000. 1

[8] P. Domingos. Why does bagging work? A Bayesian account and its implications. In *Proceedings of the International Conference on Knowledge Discovery and Data Mining (KDD)*, 1997. 1

[9] P. Domingos. Bayesian averaging of classifiers and the overfitting problem. In *Proceedings of the International Conference on Machine Learning (ICML)*, 2000. 1

[10] M. Everingham, L. Van Gool, C. K. I. Williams, J. Winn, and A. Zisserman. The PASCAL Visual Object Classes Challenge 2007 (VOC2007) Results. http://www.pascal-network.org/challenges/VOC/voc2007/workshop/index.html. 5

[11] R. Girshick. Fast r-cnn. *arXiv preprint arXiv:1504.08083*, 2015. 10

[12] R. Girshick, J. Donahue, T. Darrell, and J. Malik. Rich feature hierarchies for accurate object detection and semantic segmentation. In *IEEE Conference on Computer Vision and Pattern Recognition*, 2014. 1, 5

[13] A. Guzman-Rivera, D. Batra, and P. Kohli. Multiple Choice Learning: Learning to Produce Multiple Structured Outputs. In *Advances in Neural Information Processing Systems*, 2012. 1, 2, 3, 6, 7

[14] A. Guzman-Rivera, P. Kohli, D. Batra, and R. Rutenbar. Efficiently enforcing diversity in multi-output structured prediction. In *AISTATS*, 2014. 2, 3

[15] A. Guzman-Rivera, P. Kohli, B. Glocker, J. Shotton, T. Sharp, A. Fitzgibbon, and S. Izadi. Multi-output learning for camera relocalization. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2014. 2, 3

[16] L. K. Hansen and P. Salamon. Neural network ensembles. *IEEE Transactions on Pattern Analysis and Machine Intelligence*, 12(10):993–1001, 1990. 2

[17] G. E. Hinton, O. Vinyals, and J. Dean. Distilling the knowledge in a neural network. In *NIPS 2014 Deep Learning Workshop*, 2014. 1, 2, 7

[18] R. A. Jacobs, M. I. Jordan, S. J. Nowlan, and G. E. Hinton. Adaptive mixtures of local experts. *Neural Computation*, 3(1):79–87, Mar. 1991. 1

[19] Y. Jia. Caffe: An open source convolutional architecture for fast feature embedding. http://caffe.berkeleyvision.org/, 2013. 2, 3

[20] A. Krizhevsky. Learning multiple layers of features from tiny images, Univeristy of Toronto Tech Report. 3

[21] A. Krizhevsky, I. Sutskever, and G. Hinton. Imagenet classification with deep convolutional neural networks. In *Advances in Neural Information Processing Systems*, pages 1097–1105, 2012. 2, 3

[22] A. Krogh, J. Vedelsby, et al. Neural network ensembles, cross validation, and active learning. *Advances in Neural Information Processing Systems*, pages 231–238, 1995. 2

[23] H. Lee, E. Kim, and W. Pedrycz. A new selective neural network ensemble with negative correlation. *Applied intelligence*, 37(4):488–498, 2012. 2

[24] K. Lenc and A. Vedaldi. Understanding image representations by measuring their equivariance and equivalence. *arXiv preprint arXiv:1411.5908*, 2014. 1

[25] M. Lin, Q. Chen, and S. Yan. Network in network. *CoRR*, abs/1312.4400, 2013. 3

[26] T. P. Minka. Bayesian model averaging is not model combination. http://research-srv.microsoft.com/en-us/um/people/minka/papers/minka-bma-isnt-mc.pdf, 2002. 1

[27] K. Monteith, J. Carroll, K. Seppi, and T. Martinez. Turning bayesian model averaging into bayesian model combination. In *International Joint Conference on Neural Networks (IJCNN)*, 2011. 1

[28] R. Polikar. Ensemble learning. In *Ensemble Machine Learning*, pages 1–34. Springer, 2012. 1

[29] B. E. Rosen. Ensemble learning using decorrelated neural networks. *Connection Science*, 8(3-4):373–384, 1996. 2

[30] O. Russakovsky, J. Deng, J. Krause, A. Berg, and L. Fei-Fei. The ImageNet Large Scale Visual Recognition Challenge 2012 (ILSVRC2012). http://www.image-net.org/challenges/LSVRC/2012/. 3, 6

[31] 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. *arXiv preprint arXiv:1409.0575*, 2014. 1

[32] T. Sercu, C. Puhrsch, B. Kingsbury, and Y. LeCun. Very deep multilingual convolutional neural networks for LVCSR. *CoRR*, abs/1509.08967, 2015. 2

[33] K. Simonyan and A. Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014. 2

[34] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. A. Riedmiller. Striving for simplicity: The all convolutional net. *CoRR*, abs/1412.6806, 2014. 6, 13

[35] I. Sutskever, O. Vinyals, and Q. V. Le. Sequence to sequence learning with neural networks. In *Advances in Neural Information Processing Systems*, pages 3104–3112, 2014. 2

[36] C. Szegedy, W. Liu, Y. Jia, P. Sermanet, S. Reed, D. Anguelov, D. Erhan, V. Vanhoucke, and A. Rabinovich. Going deeper with convolutions. *arXiv preprint arXiv:1409.4842*, 2014. 1, 2

[37] M. D. Zeiler and R. Fergus. Visualizing and understanding convolutional networks. In *European Conference on Computer Vision*, pages 818–833. Springer, 2014. 1, 6

[38] Z.-H. Zhou, J. Wu, and W. Tang. Ensembling neural networks: many could be better than all. *Artificial intelligence*, 137(1):239–263, 2002. 2

# Appendix A. TreeNet Object Detection Results on PASCAL VOC 2007

As briefly described in Section 5 of the main paper, the ILSVRC-Alex TreeNet architecture was also evaluated for object detection using the PASCAL VOC 2007 dataset, which includes labeled ground-truth bounding-box annotations for 20 object classes. For this task, we used Fast R-CNNs [11]. During training, Fast RCNNs finetune a model pretrained on ImageNet for classification under two losses, one over the predicted class of an object proposal, and one with bounding box regression. For our ensembles, we average both the class prediction as well as the bounding box coordinates from ensemble member models.

To evaluate TreeNets and standard ensembles, we fine-tune four different instances for each under the Fast R-CNN framework and compute the mean and standard deviation of the classwise average precisions (APs) as well as the mean APs over all classes. Table 7 presents these results for various models with the averaged bounding boxes – a standard ensemble, a TreeNets split after conv1, conv2, and conv3, as well as a single model. We remind the reader that non-parameterized layers are irrelevant with respect to splits so we do not report results for those layers. We also evaluate without the bounding-box regression, instead using the initial selective search proposals directly. Table 8 shows these results.

In both tasks we see that TreeNets outperform the standard ensembles and single models by significant margins. We note that we see similar gains in accuracy when using the regressed bounding boxes for both single models and ensembles, implying that the bounding box averaging procedure for ensembles is reasonable.

| $\mu \pm \sigma$ | aeroplane | bicycle | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | motorbike | person | plant | sheep | sofa | train | tv | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ensemble | 67.49 ±1.42 | 68.89 ±0.57 | 52.00 ±1.25 | 38.25 ±2.38 | 16.74 ±0.90 | 68.14 ±1.90 | 70.68 ±0.37 | 67.56 ±2.65 | 26.45 ±0.16 | 63.78 ±1.63 | 61.93 ±1.59 | 61.74 ±1.69 | 73.31 ±0.68 | 67.18 ±0.56 | 56.53 ±0.30 | 23.84 ±0.58 | 50.45 ±0.82 | 54.81 ±1.45 | 69.04 ±1.45 | 59.26 ±0.88 | 55.90 ±0.21 |
| conv1 | **67.83** ±1.34 | 68.47 ±0.85 | **52.66** ±1.41 | **37.90** ±0.98 | **18.01** ±0.81 | 69.53 ±1.85 | 70.83 ±0.51 | **67.93** ±0.49 | **27.30** ±0.84 | 61.59 ±2.31 | **62.96** ±0.65 | 62.06 ±1.16 | **74.89** ±1.14 | 67.97 ±0.48 | 57.43 ±0.33 | 23.80 ±0.50 | 50.94 ±2.06 | **56.42** ±0.55 | 70.57 ±1.53 | 59.79 ±0.91 | 56.44 ±0.23 |
| conv2 | 67.30 ±1.05 | **69.29** ±0.79 | 52.62 ±1.41 | 37.58 ±1.41 | 17.90 ±0.82 | 68.81 ±0.73 | **71.04** ±0.28 | 68.54 ±1.50 | 27.13 ±1.05 | **63.66** ±1.06 | 62.37 ±1.83 | **62.20** ±0.62 | 74.60 ±0.56 | **68.45** ±0.23 | **57.52** ±0.46 | **24.20** ±1.24 | **52.53** ±0.59 | 55.00 ±0.73 | **71.31** ±1.64 | **60.66** ±0.19 | **56.64** ±0.23 |
| conv3 | 66.29 ±1.49 | 67.35 ±1.05 | 50.22 ±1.40 | 36.23 ±2.54 | 16.65 ±0.83 | 67.77 ±2.21 | 70.22 ±0.73 | 67.73 ±0.48 | 25.01 ±1.62 | 60.91 ±3.27 | 61.80 ±1.03 | 61.99 ±0.69 | 73.64 ±0.51 | 68.38 ±0.92 | 56.42 ±1.02 | 21.57 ±0.42 | 49.75 ±1.95 | 55.54 ±1.26 | 70.03 ±0.77 | 58.81 ±0.58 | 55.32 ±0.54 |
| Single Model | 62.53 ±0.24 | 65.25 ±0.49 | 41.30 ±2.12 | 32.36 ±2.47 | 11.98 ±1.11 | 62.56 ±1.27 | 66.89 ±0.80 | 61.76 ±2.14 | 20.72 ±0.58 | 56.16 ±3.32 | 56.91 ±1.58 | 55.14 ±1.56 | 69.53 ±1.71 | 64.48 ±1.27 | 51.39 ±1.10 | 21.06 ±0.36 | 45.20 ±0.77 | 47.91 ±1.44 | 65.70 ±0.89 | 54.35 ±1.94 | 50.66 ±0.30 |

Table 7: Average Precision for Object Detection using different TreeNet models with the Fast R-CNN framework, when the coordinates of bounding boxes from each member model are averaged.
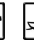
| $\mu \pm \sigma$ | aeroplane | bicycle | bird | boat | bottle | bus | car | cat | chair | cow | table | dog | horse | motorbike | person | plant | sheep | sofa | train | tv | mean |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ensemble | **64.48** ±0.75 | 63.99 ±0.94 | 45.42 ±0.71 | 34.17 ±1.14 | 15.60 ±0.26 | 63.90 ±1.46 | 67.81 ±0.34 | 62.55 ±0.90 | 24.20 ±0.27 | **58.36** ±2.09 | 55.91 ±1.34 | **56.81** ±1.14 | 62.61 ±1.24 | **67.76** ±0.52 | 50.35 ±0.35 | 22.25 ±0.84 | 45.79 ±1.68 | 49.47 ±1.76 | 65.20 ±1.57 | 57.76 ±0.72 | 51.72 ±0.15 |
| conv1 | 64.32 ±1.01 | **64.80** ±0.73 | **46.45** ±1.04 | **36.04** ±1.69 | **17.03** ±1.02 | 65.61 ±1.43 | 68.24 ±0.88 | 63.21 ±1.43 | **25.06** ±0.36 | 56.79 ±2.57 | 56.90 ±0.99 | 56.80 ±1.33 | **62.32** ±1.90 | 66.58 ±1.31 | 50.68 ±0.38 | **22.68** ±0.81 | 46.66 ±1.99 | **49.76** ±0.97 | 66.03 ±0.75 | 58.24 ±1.53 | **52.21** ±0.42 |
| conv2 | 64.22 ±1.04 | 63.68 ±1.51 | 45.85 ±1.24 | 35.31 ±0.63 | 16.43 ±0.70 | **65.89** ±0.25 | **68.78** ±0.25 | 63.21 ±1.60 | 24.57 ±0.86 | 56.37 ±1.40 | 57.10 ±1.50 | 56.57 ±1.28 | 62.28 ±0.77 | 66.69 ±1.49 | **51.04** ±0.56 | 21.72 ±0.96 | **48.19** ±0.80 | 48.86 ±1.31 | 66.24 ±2.19 | **58.71** ±1.06 | 52.08 ±0.33 |
| conv3 | 63.56 ±1.10 | 63.91 ±0.75 | 44.18 ±0.67 | 34.24 ±2.78 | 16.01 ±0.96 | 64.36 ±1.57 | 67.71 ±0.86 | 63.00 ±1.41 | 23.27 ±0.82 | 56.99 ±1.71 | **57.21** ±0.46 | 56.31 ±1.49 | 61.05 ±1.65 | 66.04 ±1.18 | 50.22 ±0.85 | 20.22 ±0.29 | 45.03 ±2.39 | 49.40 ±0.73 | **66.26** ±1.80 | 56.92 ±0.56 | 51.29 ±0.44 |
| Single Model | 59.11 ±1.49 | 60.82 ±1.04 | 35.84 ±1.86 | 28.84 ±1.20 | 10.73 ±0.86 | 58.91 ±2.04 | 63.98 ±0.69 | 55.84 ±1.31 | 19.40 ±0.78 | 50.03 ±2.91 | 51.65 ±2.91 | 47.95 ±2.48 | 58.38 ±2.49 | 61.44 ±1.06 | 44.81 ±0.99 | 18.61 ±0.38 | 41.64 ±2.17 | 42.74 ±1.32 | 61.73 ±1.38 | 51.89 ±0.63 | 46.22 ±0.26 |

Table 8: Average Precision for Object Detection using different TreeNet models with the Fast R-CNN framework, when predicted bounding boxes are not used.

## Appendix B. Instability of Averaged Softmax Outputs

As discussed in Section 6.1 of the main paper, training under a cross-entropy loss over averaged softmax outputs results in reduced performance compared to both standard ensembles and score-averaged ensembles. We find that this is because averaging softmax outputs prior to the cross-entropy loss has less stable gradients compared to standard cross-entropy over softmax outputs. Let us consider the standard case first and formulate the derivative of the cross-entropy loss with respect to softmax inputs. The cross-entropy loss and softmax function are defined as:

$$\ell(x, y) = -\log(p_y) \quad \text{and} \quad p_y = \frac{e^{s_y}}{\sum_i e^{s_i}}$$

The derivative of the softmax probability $p_y$ with respect to some score $s_j$ is

$$\frac{\partial p_y}{\partial s_j} = p_y(I_{jy} - p_j)$$

where $I_{jy}$ is 1 if $y = j$ and 0 otherwise. This derivative requires multiplying probabilities which can be quite small, leading to underflow errors. Taking the derivative of the cross-entropy loss with respect to some $s_j$ results in a more stable solution:

$$\frac{\partial \ell}{\partial s_j} = \frac{\partial \ell}{\partial p_y} \frac{\partial p_y}{\partial s_j} = -\frac{1}{p_y} p_y(I_{jy} - p_j) = p_j - I_{jy}$$

Let us now consider the case where $p_y$ is averaged over $M$ predictors such that

$$p_y = \frac{1}{M} \sum_m p_y^m \quad \text{and} \quad p_y^m = \frac{e^{s_y^m}}{\sum_i e^{s_i^m}}$$

The derivative of this new $p_y$ with respect to the score of one predictor $s_j^m$ is then

$$\frac{\partial p_y}{\partial s_j^m} = \frac{\partial p_y}{\partial p_y^m} \frac{\partial p_y^m}{\partial s_j^m} = \frac{1}{M} \frac{\partial p_y^m}{\partial s_j^m} = \frac{1}{M} p_y^m(I_{jy} - p_j^m)$$

Again computing the derivative of the loss with respect to a score $s_j^m$ we see

$$\frac{\partial \ell}{\partial s_j^m} = \frac{\partial \ell}{\partial p_y} \frac{\partial p_y}{\partial p_y^m} \frac{\partial p_y^m}{\partial s_j^m} = \frac{1}{p_y} \frac{1}{M} p_y^m(p_j^m - I_{jy}) = \frac{p_y^m}{\sum_i p_y^i}(p_j^m - I_{jy})$$

The rightmost term in this result is identical to the standard case presented above; however, the first term acts to weight the gradient for each predictor and can be shown to range from 0 to M. The product of this term and the probability $p_j^m$ can be prone to underflow errors when $p_y^m$ is less than $p_y$. On the other hand, when $p_y^m$ is greater than $p_y$ the gradients are increased in magnitude which can result in overshooting good minima.

This scaling of the gradients has an interesting similarity with MCL. If a predictor puts little mass into the correct class compared to the other predictors, the weighting factor and thus the gradient go to zero – meaning worse performing members are less encouraged to improve than strong performers. This is similar behavior to what a soft-assignment variant of MCL might induce. However, we do not notice improved oracle accuracy relative to base ensembles for models trained with probability-averaged losses, implying the predictors are making relatively similar predictions.

# Appendix C. Pseudo-code for Stochastic Gradient Descent with MCL

We describe the classical MCL algorithm and our approach to integrate MCL coordinate descent with stochastic gradient descent in Section 6.2 of the main paper. Here we provide psuedocode for both algorithms to highlight the differences and provide additional clarity.

**Data**: Dataset $(x_i, y_i) \in D$ and loss $\mathcal{L}$
**Result**: Predictor parameters $\theta_1 \cdots \theta_M$
**Initialization:**
 $D_0^* = \{D_1 \cdots D_M\} \leftarrow$ k-means$(D, k = M)$
 $t \leftarrow 0$
**while** $D_t^* \neq D_{t-1}^*$ **do**
    **Step 1: Train each predictor to completion using only its corresponding subset of the data**
        $\theta_m \leftarrow$ Train$(D_m)$
    **Step 2: Reassign each example to its least-loss predictor**
        $D_m = \{d \in D | \forall \theta_j, \mathcal{L}(d; \theta_m) \leq \mathcal{L}(d; \theta_j)\}$
    $t \leftarrow t + 1$
**end**

**Algorithm 1:** Classical MCL

**Data**: Dataset $(x_i, y_i) \in D$, SGD parameters $\eta, \lambda$, and loss $\mathcal{L} = \sum_i \ell$
**Result**: Network parameters $\theta_1 \cdots \theta_M$
**Initialization:**
 Randomly initialize $\theta_1 \cdots \theta_M$
 $t \leftarrow 0$
**while** $\mathcal{L}_{set}^t < \mathcal{L}_{set}^{t-1}$ **do**
    $t \leftarrow t + 1$
    Sample batch $B \subset D$
    **Step 1: Forward pass**
      For $b \in B$, compute forward-pass and losses $\ell(b; \theta_1) \cdots \ell(b; \theta_M)$
      Partition $B$ by updating indicator variables $\alpha$ as:
        $\alpha_{mi} = 1[[m = argmin_{m' in 1:M} \ell(b; theta'_m)]]$
      $\mathcal{L}_m = \sum_i \alpha_{mi} \ell(b_i; \theta_m)$
      $\mathcal{L}_{set}^t = \sum \mathcal{L}_m$
    **Step 2: Backward pass**
      For each $\theta_m$ apply gradient descent update using only the subset of examples on which it achieves the lowest loss
      $\theta_m \leftarrow \theta_m - \eta \nabla \mathcal{L}_m - \lambda \Delta \theta_m$
**end**

**Algorithm 2:** Integrating MCL coordinate descent with SGD steps

# Appendix D. Visualizations for MCL Trained Ensembles

In this section we present additional insight into how MCL ensemble training differs from the behavior of standard ensembles. To show how the distribution of class examples changes over training for MCL we have produced a video showing the proportion of each CIFAR10 class assigned to each predictor at test time and how it changes over training iterations. The intensity of each class icon is proportional to the fraction of class examples assigned to a predictor. Figure 4 shows a sample early and later frame from the video.
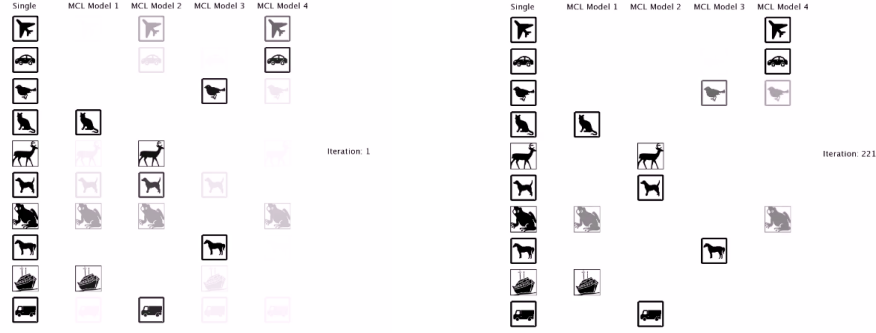


Figure 4: The left image shows the class distribution early in training – notice how many classes are split between multiple predictors. The right frame shows the evolution of this distribution after 220 additional iterations. Many of the classes have stabilized.

We also present additional guided-backprop [34] visualizations described in Section 6.2 of the main paper for different layers in members of traditional and MCL ensembles. These images visualize how the ensemble members respond to input images. These images can be interpreted as the gradient of a neuron output with respect to the input image. Features that are clear in these images have the largest influence on the network's output for a given input image. Figure 5 shows these visualizations taken for an input image with respect to its true class label. Notice that ensemble members are agnostic to classes that they are not specialized in. The input images are those that produce the highest correct response on the ensemble model. Visualisations of the same neurons in Figure 6 are generated independently for each model using the image that gives the highest activation. We note that while there is a greater response for non-specialized ensemble members, they remain largely indifferent to image content. We see similar patterns of indifference in lower convolutional layer visualizations as well shown in Figures 7, 8, and 9.



Figure 5: Reconstructions using features from the output layer using the images that give highest activation for the single model. Column1: Model from Standard Ensemble, Column2-4: Members of ensemble trained using MCL Loss
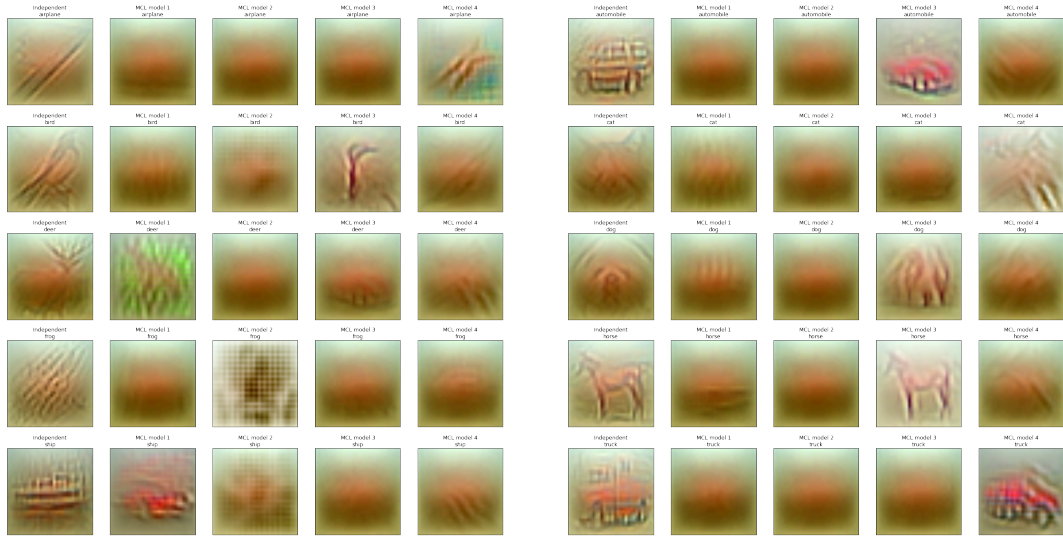
Figure 6: Reconstructions using features from the output layer using the images that give highest activation for each model independently. Column1: Model from Standard Ensemble, Column2-4: Members of ensemble trained using MCL Loss
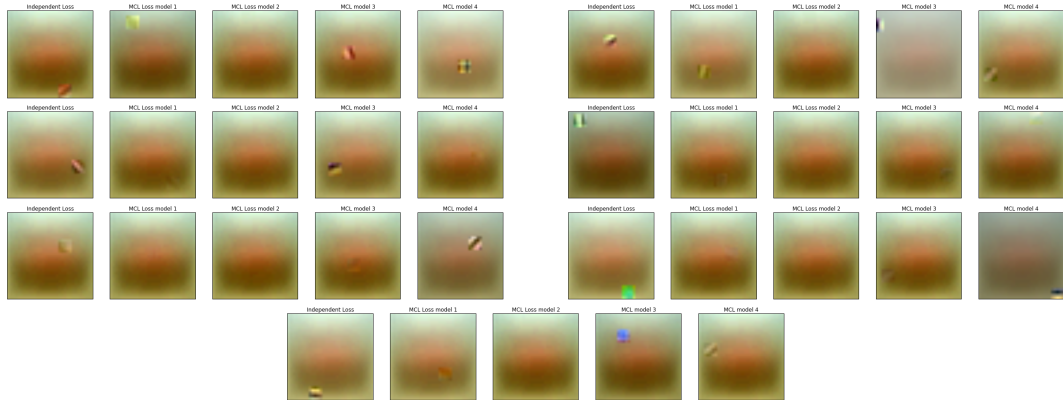


Figure 7: Reconstructions using the conv1 layer. Column1: Model from Standard Ensemble, Column2-4: Members of ensemble trained using MCL Loss
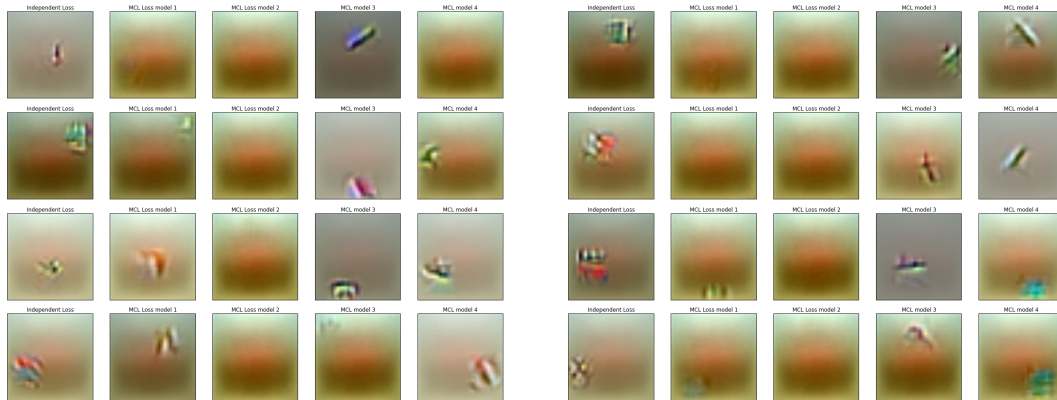
Figure 8: Reconstructions using the conv2 layer. Column1: Model from Standard Ensemble, Column2-4: Members of ensemble trained using MCL Loss
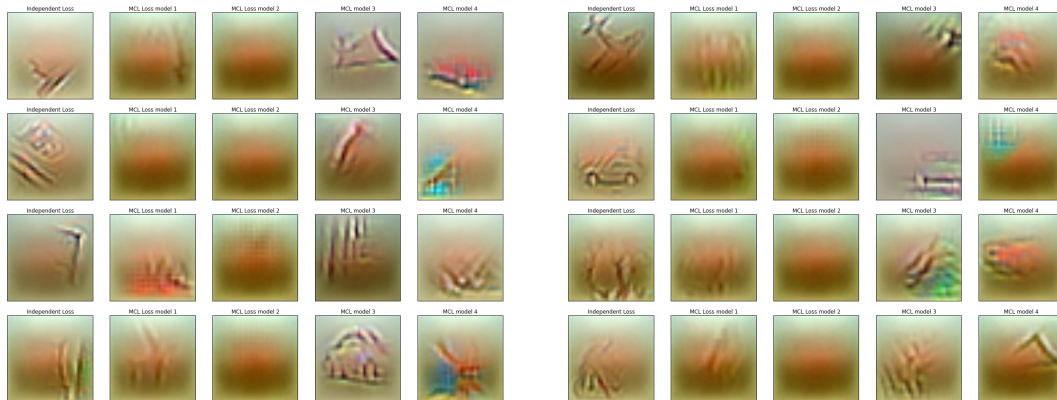


Figure 9: Reconstructions using the conv3 layer. Column1: Model from Standard Ensemble, Column2-4: Members of ensemble trained using MCL Loss

# Appendix E. MPI-Caffe

`MPI-Caffe` is a modification of the popular `Caffe` deep learning framework that enables cross-GPU/cross-machine communication on MPI enabled systems as model layers. Providing these MPI operations as layers allows network designers the flexibility to quickly experiment with distributed networks while abstracting away much of the communication logic. This enables experimentation with extremely large (i.e. larger than can be held in a single GPU) networks as well as ensemble-aware model parallelism schemes. This document explains the function of these layers as well as providing example usage. The core functionality in `MPI-Caffe` is provided by

- the `MPIBroadcast` layer discussed in Section Appendix E.1.1

- and the `MPIGather` layer discussed in Section Appendix E.1.2.

The primary file defining the interface of the MPI layers is `MPILayers.hpp`. There are also many supporting modifications in the source that should be noted in case anyone tries to modify or update the base `Caffe` version. The network initialization code in `net.cpp` has been substantially altered to accommodate the distributed framework. Some other changes occur in `layer.hpp`, `solver.cpp`, and `caffe.cpp` among others.

## Appendix E.1. A Toy Example

Let's start with a toy example to build context for the MPI layer descriptions. Suppose we want to train an TreeNet ensemble of CIFAR10-Quick and we want train it under a score-averaged loss. Figure 10 shows how we might modify the LeNet structure using `MPI-Caffe` to implement this model across three processes/GPUs in an MPI enabled cluster. We will go through this example to explain the function and parametrization of the `MPIBroadcast` and `MPIGather` layers.
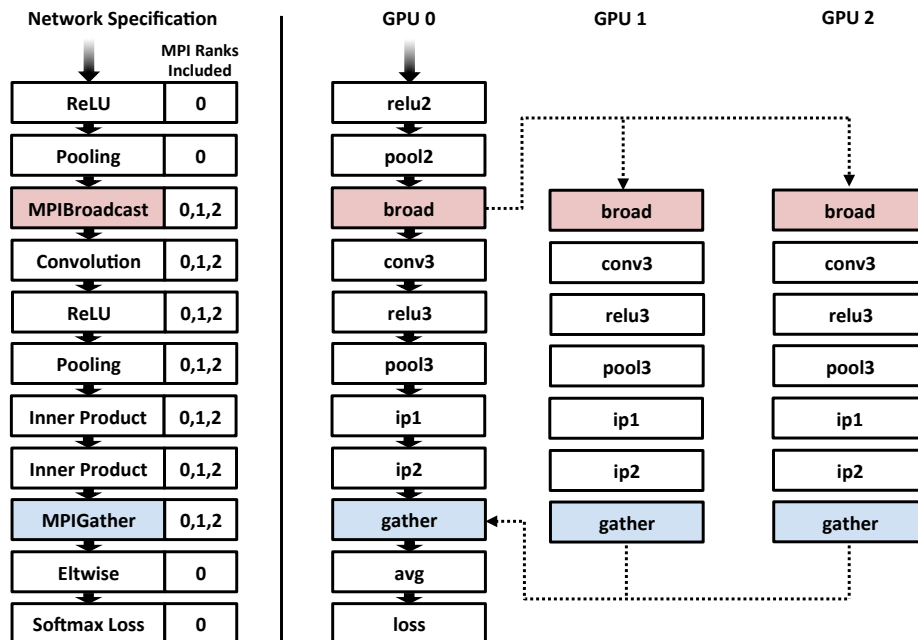


Figure 10: **Left)** Example network specification for `MPI-Caffe` enabled CIFAR10-Quick ensemble with parameter sharing and model averaging. **Right)** Distributed model resulting from the specification, visualized for each participating process. Dashed lines indicate MPI cross process communication and not layer input/outputs.

## Appendix E.1.1   MPIBroadcast

The first layer we discuss is `MPIBroadcast` (highlighted in red in Figure 10). The `MPIBroadcast` layer broadcasts a copy of its input blob to each process in its communication group during its forward pass. During the backward pass, the gradients from each copy are summed and passed back to the input blob. The communication group consists of all processes that carry a copy of a particular broadcast layer. By default a communication group contains all processes; however, adding `mpi_rank:n` rules in either the `include` or `exclude` layer parameters can alter this group.

```
1 layer{
2   name: broad
3   type: MPIBroadcast
4   bottom: pool2
5   top: pool2_b
6   mpi_param{
7     root: 0
8   }
9   include{
10     mpi_rank: 0
11     mpi_rank: 1
12     mpi_rank: 2
13   }
14 }
```
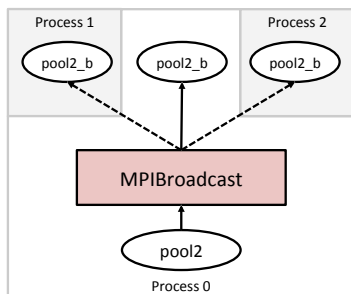


Figure 11: An example `MPIBroadcast` layer definition and a diagram of its forward-pass behavior.

Figure 11 shows the `MPIBroadcast` layer definition from our example and the corresponding forward-pass behavior. Going step by step through the definition: we

- [Line 2-5] declare this layer to be a MPI Broadcast layer named "broad" with input blob `pool2` and output blob `pool2_b`,

- [Line 6-8] set the `mpi_param` root value to 0 indicating that process 0 will be initiating the broadcast,

- and [Line 9-13] establish a communication group consisting of processes 0, 1, and 2.

During a forward pass, the `MPIBroadcast` layer on process 0 will send a copy of `pool2` to processes 1 and 2 as well as retain a copy for itself. For the example, we would also need to modify the `ip1` layer to take `pool2_b` as input rather than `pool2`.

It is important to note the effect the choice of `mpi_param{ root }` has on network structure. As shown in the example in Figure 10, each process parses the entire network structure and retains only the layers that include its MPI rank. For process 0, this includes the entire network, but for processes 1 and 2 the network *starts* with the `MPIBroadcast` layer. In order to allow this behavior, non-root processes have the input blob (`pool2` in our example) stripped out during network parsing. Additionally for this example we need to average these top blobs before sending the result into the softmax loss.

### Appendix E.1.2    MPIGather

If the purpose of a broadcast layer is to take some data and push copies into multiple process spaces, the `MPIGather` layer can be thought of as the opposite. In a forward pass, it takes multiple copies of a blob from multiple process spaces and collects them in the root process. During a backward pass, the gradients for each top blob are routed back to the corresponding input blob and process. Similar to the previous section, Figure 12 shows the layer definition from out example and a diagram of the forward pass behavior.

The `mpi_param{root}` parameter in the gather layer defines which process will be receiving the gathered blobs and producing top blobs. In analogy to the broadcast layer parsing, gather layers in non-root processes are pruned of the top blobs during network parsing (see Figure 10).

There are some restrictions to the gather layer's use. First, the bottom blob (`ip2` in our example) must be defined in all communication group processes. Second, the number of top blobs must equal the number of processes in the communication group. Both of these conditions are checked by the source and will report an error if not satisfied.

```
1 layer{
2   name: ExampleLayer
3   type: MPIGather
4   bottom: ip2
5   top: ip2_0
6   top: ip2_1
7   top: ip2_2
8   mpi_param{
9     root: 0
10   }
11 include{
12     mpi_rank: 0
13     mpi_rank: 1
14     mpi_rank: 2
15   }
16 }
```
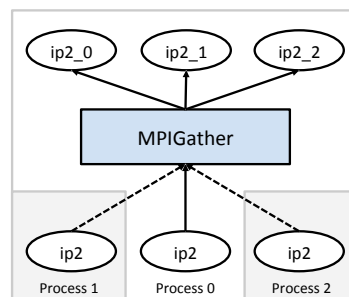


Figure 12: An example `MPIGather` layer definition and a diagram of its forward-pass behavior.

### Appendix E.2. Notes and Other Examples

It is worth noting a few other use points about `MPI-Caffe`:

- the `MPIBroadcast` layer can be used to construct a very large single-path network spanned across multiple GPU's

- the `MPIGather` layer can be used to allow more sophisticated ensemble losses

- there is no limit on the number or order of MPI layers such that complex distributed networks are possible

- in situations where network latency is lower than reading from disk, the `MPIBroadcast` layer can be used to train multiple independent networks more quickly

## Appendix E.3. Communication Cost Analysis

We tested our MPI-Caffe framework on a large-scale cluster with one Tesla K20 GPU per node and a maximum MPI node interconnect bandwidth of 5.8 GB/sec. To characterize the communication overhead for an ensemble, we measure the time spent sharing various layers of the ILSVRC-Alex×5 architecture. Each network was run on a separate node (with one node also holding the shared layers). Figure 13 shows the communication time to share a given layer as a fraction of the forward-backward pass. The x-axis indicates the number of floats broadcast per batch for each layer. We note that for these layers overhead appears approximately linear and even the largest layer incurs very little overhead for communication.
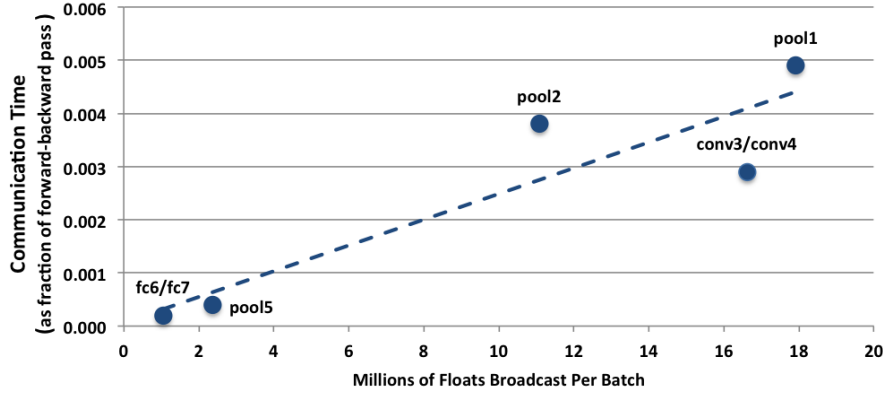


Figure 13: Fraction of forward-backward pass time used for TreeNets sharing various layers against the size of the layers. The overhead from communication is quite small and scales approximately linearly with the size of the layer being shared.