

PITFALLS OF IN-DOMAIN UNCERTAINTY ESTIMATION AND ENSEMBLING IN DEEP LEARNING

Arsenii Ashukha *

Samsung AI Center Moscow, HSE[‡]
aashukha@bayesgroup.ru

Alexander Lyzhov *

Samsung AI Center Moscow, Skoltech[†], HSE[§]
alyzhov@bayesgroup.ru

Dmitry Molchanov *

Samsung AI Center Moscow, HSE[‡]
dmolch@bayesgroup.ru

Dmitry Vetrov

HSE[‡], Samsung AI Center Moscow
dvetrov@bayesgroup.ru

ABSTRACT

Uncertainty estimation and ensembling methods go hand-in-hand. Uncertainty estimation is one of the main benchmarks for assessment of ensembling performance. At the same time, deep learning ensembles have provided state-of-the-art results in uncertainty estimation. In this work, we focus on in-domain uncertainty for image classification. We explore the standards for its quantification and point out pitfalls of existing metrics. Avoiding these pitfalls, we perform a broad study of different ensembling techniques. To provide more insight in this study, we introduce the *deep ensemble equivalent score* (DEE) and show that many sophisticated ensembling techniques are equivalent to an ensemble of only few independently trained networks in terms of test performance.



1 INTRODUCTION

Deep neural networks (DNNs) have become one of the most popular families of machine learning models. The predictive performance of DNNs for classification is often measured in terms of accuracy. However, DNNs have been shown to yield inaccurate and unreliable *probability* estimates, or predictive uncertainty (Guo et al., 2017). This has brought considerable attention to the problem of uncertainty estimation with deep neural networks.

There are many faces to uncertainty estimation. Different desirable uncertainty estimation properties of a model require different settings and metrics to capture them. *Out-of-domain uncertainty of the model is measured on the data that does not follow the same distribution as the training dataset (out-of-domain data)*. Out-of-domain data can include images corrupted with rotations or blurring, adversarial attacks (Szegedy et al., 2013) or data points from a completely different dataset. The model is expected to be resistant to data corruptions and to be more uncertain on out-of-domain data than on in-domain data. On the contrary, *in-domain uncertainty of the model is measured on data taken from the training data distribution, i.e. data from the same domain*. In this setting the model is expected to produce reliable probability estimates, e.g. the model shouldn't be too overconfident in its wrong predictions.

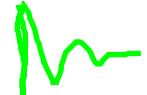
Pitfalls of metrics We show that many common metrics of in-domain uncertainty estimation (e.g. log-likelihood, Brier score, calibration metrics, etc.) are either not comparable across different models or fail to provide a reliable ranking. We address some of the stated pitfalls and point out more reasonable evaluation schemes. For instance, although *temperature scaling* is not a standard for ensembling techniques, it is a must for a fair evaluation. With this in mind, the *calibrated log-likelihood* avoids most of the stated pitfalls and generally is a reasonable metric for in-domain uncertainty estimation task.

*Equal contribution [§]HSE refers to National Research University Higher School of Economics

[†]Skoltech refers to Skolkovo Institute of Science and Technology

[‡]HSE refers to Samsung-HSE Laboratory, National Research University Higher School of Economics

Pitfalls of ensembles Equipped with the proposed evaluation framework, we are revisiting the evaluation of ensembles of DNNs—one of the major tools for uncertainty estimation. We introduce the *deep ensemble equivalent* (DEE) score that measures the number of independently trained models that, when ensembled, achieve the same performance as the ensembling technique of interest. The DEE score allows us to compare ensembling techniques across different datasets and architectures using a unified scale. Our study shows that most of the popular ensembling techniques require averaging predictions across dozens of samples (members of an ensemble), yet are essentially equivalent to an ensemble of only few independently trained models.



Missing part of ensembling In our study, test-time data augmentation (TTA) turned out to be a surprisingly strong baseline for uncertainty estimation and a simple way to improve ensembles. Despite being a popular technique in large-scale classification, TTA seems to be overlooked in the community of uncertainty estimation and ensembling.

2 SCOPE OF THE PAPER

We use standard benchmark problems of image classification which comprise a common setting in research on learning ensembles of neural networks. There are other relevant settings where the correctness of probability estimates can be a priority, and ensembling techniques are used to improve it. These settings include, but are not limited to, regression, language modeling (Gal, 2016), image segmentation (Gustafsson et al., 2019), active learning (Settles, 2012) and reinforcement learning (Buckman et al., 2018; Chua et al., 2018).

We focus on in-domain uncertainty, as opposed to out-of-domain uncertainty. Out-of-domain uncertainty includes detection of inputs that come from a completely different domain or have been corrupted by noise or adversarial attacks. This setting has been thoroughly explored by (Ovadia et al., 2019).

We only consider methods that are trained on clean data with simple data augmentation. Some other methods use out-of-domain data (Malinin & Gales, 2018) or more elaborate data augmentation, e.g. mixup (Zhang et al., 2017) or adversarial training (Lakshminarayanan et al., 2017) to improve accuracy, robustness and uncertainty.

We use conventional training procedures. We use the stochastic gradient descent (SGD) and use batch normalization (Ioffe & Szegedy, 2015), both being the de-facto standards in modern deep learning. We refrain from using more elaborate optimization techniques including works on super-convergence (Smith & Topin, 2019) and stochastic weight averaging (SWA) (Izmailov et al., 2018). These techniques can be used to drastically accelerate training and to improve the predictive performance. Thus, we do not comment on the training time of different ensembling methods since the use of these and other more efficient training techniques would render such a comparison obsolete.

A number of related works study ways of approximating and accelerating prediction in ensembles. The distillation mechanism allows to approximate the prediction of an ensemble by a single neural network (Hinton et al., 2015; Balan et al., 2015; Tran et al., 2020), whereas fast dropout (Wang & Manning, 2013) and deterministic variational inference (Wu et al., 2018) allow to approximate the predictive distribution of specific stochastic computation graphs. We measure the raw power of ensembling techniques without these approximations.

All of the aforementioned alternative settings are orthogonal to the scope of this paper and are promising points of interest for further research.

3 PITFALLS OF IN-DOMAIN UNCERTAINTY ESTIMATION

No single metric measures all the desirable properties of uncertainty estimates obtained by a model of interest. Because of this, the community is using many different metrics in an attempt to capture the quality of uncertainty estimation, such as the Brier score (Brier, 1950), log-likelihood (Quinonero-Candela et al., 2005), metrics of calibration (Guo et al., 2017; Nixon et al., 2019), performance of misclassification detection (Malinin & Gales, 2018), and threshold-accuracy curves (Lakshminarayanan et al., 2017). In the section we highlight the pitfalls of the aforementioned met-

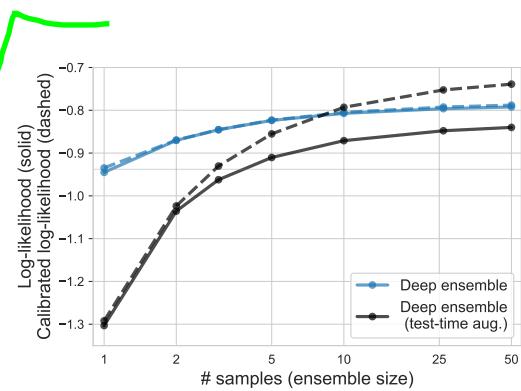


Figure 1: The average log-likelihood of two different ensembling techniques for ResNet50 on ImageNet dataset before (solid) and after (dashed) temperature scaling. Without the temperature scaling, test-time data augmentation decreases the log-likelihood of plain deep ensembles. However, when the temperature scaling is enabled, deep ensembles with test-time data augmentation outperform plain deep ensembles.

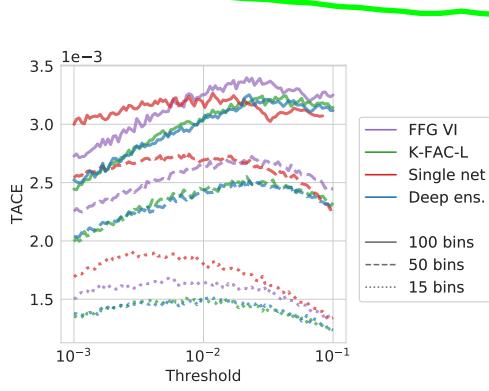


Figure 2: Thresholded adaptive calibration error (TACE) is highly sensitive to the threshold and the number of bins. It does not provide a consistent ranking of different ensembling techniques. Here TACE is reported for VGG16BN model on CIFAR-100 dataset and is evaluated at the optimal temperature.

rics, and demonstrate that these pitfalls can significantly affect evaluation, changing the ranking of the methods.

Notation We consider a classification problem with a dataset that consists of N training and n testing pairs $(x_i, y_i^*) \sim p(x, y)$, where x_i is an object and $y_i^* \in \{1, \dots, C\}$ is a discrete class label. A probabilistic classifier maps an object x_i into a predictive distribution $\hat{p}(y | x_i)$. The predictive distribution $\hat{p}(y | x_i)$ of a deep neural network is typically defined by the softmax function $\hat{p}(y | x) = \text{Softmax}(z(x)/T)$, where $z(x)$ is a vector of logits and T is a scalar parameter standing for the temperature of the predictive distribution. This scalar parameter is usually set to $T = 1$ or is tuned on a validation set (Guo et al., 2017). The maximum probability $\max_c \hat{p}(y = c | x_i)$ is called the confidence of a classifier \hat{p} on an object x_i . $\mathbb{I}[\cdot]$ denotes the indicator function throughout the text.

3.1 LOG-LIKELIHOOD AND BRIER SCORE

The average test log-likelihood $\text{LL} = \frac{1}{n} \sum_{i=1}^n \log \hat{p}(y = y_i^* | x_i)$ is a popular metric for measuring the quality of in-domain uncertainty of deep learning models. It directly penalizes high probability scores assigned to incorrect labels and low probability scores assigned to the correct labels y_i^* .

LL is sensitive to the softmax temperature T . The temperature that has been implicitly learned during training can be far from optimal for the test data. However, a nearly optimal temperature can be found post-hoc by maximizing the log-likelihood on validation data. This approach is called temperature scaling or calibration (Guo et al., 2017). Despite its simplicity, the temperature scaling results in a notable improvement in the LL.

While ensembling techniques tend to have better temperature than single models, the default choice of $T = 1$ is still suboptimal. Comparing the LL with suboptimal temperatures—that is often the case in practice—can potentially produce an arbitrary ranking of different methods.

Comparison of the log-likelihood should only be performed at the optimal temperature.

Empirically, we demonstrate that the overall ordering of methods and also the best ensembling method according to the LL can vary depending on temperature T . While this applies to most ensembling techniques (see Figure 10), this effect is most noticeable on experiments with data augmentation on ImageNet (Figure 1). We will call the log-likelihood at the optimal temperature *the cal-*

ibrated log-likelihood. We show how to obtain an unbiased estimate of the calibrated log-likelihood without a held-out validation set in Section 3.5.

LL also demonstrates a high correlation with accuracy ($\rho > 0.86$), that in case of calibrated LL becomes even stronger ($\rho > 0.95$). That suggests that while (calibrated) LL measures the uncertainty of the model, it still has a significant dependence on the accuracy and vice versa. A model with higher accuracy would likely have a higher log-likelihood. See Figure 9 in Appendix C for more details.

Brier score $BS = \frac{1}{n} \frac{1}{C} \sum_{i=1}^n \sum_{c=1}^C (\mathbb{I}[y_i^* = c] - \hat{p}(y=c|x_i))^2$ has also been known for a long time as a metric for verification of predicted probabilities (Brier, 1950). Similarly to the log-likelihood, the Brier score penalizes low probabilities assigned to correct predictions and high probabilities assigned to wrong ones. It is also sensitive to the temperature of the softmax distribution and behaves similarly to the log-likelihood. While these metrics are not strictly equivalent, they show a high empirical correlation for a wide range of models on CIFAR-10, CIFAR-100 and ImageNet datasets (see Figure 8 in Appendix C).

3.2 MISCLASSIFICATION DETECTION

Detection of wrong predictions of the model, or misclassifications, is a popular downstream problem relevant to the problem of in-domain uncertainty estimation. Since misclassification detection is essentially a binary classification problem, some papers measure its quality using conventional metrics for binary classification such as AUC-ROC and AUC-PR (Malinin & Gales, 2018; Cui et al., 2019; Možejko et al., 2018). These papers use an uncertainty criterion like confidence or predictive entropy $\mathcal{H}[\hat{p}(y|x_i)]$ as a prediction score.

While these metrics can be used to assess the misclassification detection performance of a single model, they cannot be used to directly compare misclassification performance across different models. Correct and incorrect predictions are specific for every model, therefore, every model induces its own binary classification problem. The induced problems can differ significantly, since different models produce different confidences and misclassify different objects. In other words, comparing such metrics implies a comparison of performance of classifiers that solve different classification problems. Such metrics are therefore incomparable.

AUCs for misclassification detection cannot be directly compared between different models.

While comparing AUCs is incorrect in the setting of misclassification detection, it is correct to compare these metrics in many out-of-domain data detection problems. In that case, both objects and targets of the induced binary classification problems remain the same for all models. All out-of-domain objects have a positive label and all in-domain objects have a negative label. Note that this condition does not necessarily hold in the problem of detection of adversarial attacks. Different models generally have different inputs after an adversarial attack, so such AUC-based metrics might still be flawed.

3.3 CLASSIFICATION WITH REJECTION

Accuracy-confidence curves are another way to measure the performance of misclassification detection. These curves measure the accuracy on the set of objects with confidence $\max_c \hat{p}(y=c|x_i)$ above a certain threshold τ (Lakshminarayanan et al., 2017) and ignoring or rejecting the others.

The main problem with accuracy-confidence curves is that they rely too much on calibration and the actual values of confidence. Models with different temperatures have different numbers of objects at each confidence level which does not allow for a meaningful comparison. To overcome this problem, one can switch from thresholding by the confidence level to thresholding by the number of rejected objects. The corresponding curves are then less sensitive to temperature scaling and thus allow to compare the rejection ability in a more meaningful way. Such curves have been known as *accuracy-rejection curves* (Nadeem et al., 2009). In order to obtain a scalar metric for easy comparisons, one can compute the area under this curve, resulting in AU-ARC (Nadeem et al., 2009).

3.4 CALIBRATION METRICS

Informally speaking, a probabilistic classifier is calibrated if any predicted class probability is equal to the true class probability according to the underlying data distribution (see (Vaicenavicius et al., 2019) for formal definitions). Any deviation from the perfect calibration is called miscalibration. For brevity, we will use $\hat{p}_{i,c}$ to denote $\hat{p}(y = c | x_i)$ in the current section.

Expected calibration error (ECE) (Naeini et al., 2015) is a metric that estimates model miscalibration by binning the assigned probability scores and comparing them to average accuracies inside these bins. Assuming B_m denotes the m -th bin and M is overall number of bins, the ECE is defined as follows:

$$\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{n} |\text{acc}(B_m) - \text{conf}(B_m)|, \quad (1)$$

where $\text{acc}(B) = |B|^{-1} \sum_{i \in B} \mathbb{I}[\arg \max_c \hat{p}_{i,c} = y_i^*]$ and $\text{conf}(B) = |B|^{-1} \sum_{i \in B} \hat{p}_{i,y_i^*}$.

A recent line of works on measuring calibration in deep learning (Vaicenavicius et al., 2019; Kumar et al., 2019; Nixon et al., 2019) outlines several problems of the ECE score. Firstly, ECE is a biased estimate of the true calibration. Secondly, ECE-like scores cannot be optimized directly since they are minimized by a model with constant uniform predictions, making the infinite temperature $T = +\infty$ its global optimum. Thirdly, ECE only estimates miscalibration in terms of the maximum assigned probability whereas practical applications may require the full predicted probability vector to be calibrated. Finally, biases of ECE on different models may not be equal, rendering the miscalibration estimates incompatible. Similar concerns are also discussed by Ding et al. (2019).

Thresholded adaptive calibration error (TACE) was proposed as a step towards solving some of these problems (Nixon et al., 2019). TACE disregards all predicted probabilities that are less than a certain threshold (hence *thresholded*), chooses the bin locations adaptively so that each bin has the same number of objects (hence *adaptive*), and estimates miscalibration of probabilities across all classes in the prediction (not just the top-1 predicted class as in ECE). Assuming that B_m^{TA} denotes the m -th thresholded adaptive bin and M is the overall number of bins, TACE is defined as follows:

$$\text{TACE} = \frac{1}{CM} \sum_{c=1}^C \sum_{m=1}^M \frac{|B_m^{\text{TA}}|}{n} |\text{objs}(B_m^{\text{TA}}, c) - \text{conf}(B_m^{\text{TA}}, c)|, \quad (2)$$

where $\text{objs}(B^{\text{TA}}, c) = |B^{\text{TA}}|^{-1} \sum_{i \in B^{\text{TA}}} \mathbb{I}[y_i^* = c]$ and $\text{conf}(B^{\text{TA}}, c) = |B^{\text{TA}}|^{-1} \sum_{i \in B^{\text{TA}}} \hat{p}_{i,c}$.

Although TACE does solve several problems of ECE and is useful for measuring calibration of a specific model, it still cannot be used as a reliable criterion for comparing different models. Theory suggests that it is still a biased estimate of true calibration with different bias for each model (Vaicenavicius et al., 2019). In practice, we find that TACE is sensitive to its two parameters, the number of bins and the threshold, and does not provide a consistent ranking of different models, as shown in Figure 2.

3.5 CALIBRATED LOG-LIKELIHOOD AND TEST-TIME CROSS-VALIDATION

There are two common ways to perform temperature scaling using a validation set when training on datasets that only feature public training and test sets (e.g. CIFARs). The public training set might be divided into a smaller training set and validation set, or the public test set can be split into test and validation parts (Guo et al., 2017; Nixon et al., 2019). The problem with the first method is that the resulting models cannot be directly compared with all the other models that have been trained on the full training set. The second approach, however, provides an unbiased estimate of metrics such as log-likelihood and Brier score but introduces more variance.

In order to reduce the variance of the second approach, we perform a “test-time cross-validation”. We randomly divide the test set into two equal parts, then compute metrics for each half of the test set using the temperature optimized on another half. We repeat this procedure five times and average the results across different random partitions to reduce the variance of the computed metrics.

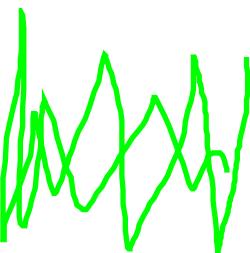
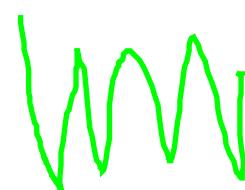
4 A STUDY OF ENSEMBLING & DEEP ENSEMBLE EQUIVALENT

Ensembles of deep neural networks have become a de-facto standard for uncertainty estimation and improving the quality of deep learning models (Hansen & Salamon, 1990; Krizhevsky et al.,

2009; Lakshminarayanan et al., 2017). There are two main directions of training ensembles of DNNs: training stochastic computation graphs and obtaining separate snapshots of neural network parameters.

Methods based on the paradigm of **stochastic computation graphs** introduce some kind of random noise over the weights or activations of deep learning models. When the model is trained, each sample of the noise corresponds to a member of the ensemble. During test time, the predictions are averaged across the noise samples. These methods include (test-time) data augmentation, dropout (Srivastava et al., 2014; Gal & Ghahramani, 2016), variational inference (Blundell et al., 2015; Kingma et al., 2015; Louizos & Welling, 2017), batch normalization (Ioffe & Szegedy, 2015; Teyé et al., 2018; Atanov et al., 2019), Laplace approximation (Ritter et al., 2018) and many more.

Snapshot-based methods aim to obtain sets of weights for deep learning models and then to average the predictions across these weights. The weights can be trained independently (e.g. deep ensembles (Lakshminarayanan et al., 2017)), collected on different stages of a training trajectory (e.g. snapshot ensembles (Huang et al., 2017) and fast geometric ensembles (Garipov et al., 2018)), or obtained from a sampling process (e.g. MCMC-based methods) (Welling & Teh, 2011; Zhang et al., 2019)). These two paradigms can be combined. Some works suggest construction of ensembles of stochastic computation graphs (Tomczak et al., 2018), while others make use of the collected snapshots to construct a stochastic computation graph (Wang et al., 2018; Maddox et al., 2019).

In this paper we consider the following ensembling techniques: deep ensembles (Lakshminarayanan et al., 2017), snapshot ensembles (SSE by Huang et al. (2017)), fast geometric ensembling (FGE by Garipov et al. (2018)), SWA-Gaussian (SWAG by Maddox et al. (2019)), cyclical SGLD (cSGLD by Zhang et al. (2019)), variational inference (VI by Blundell et al. (2015)), K-FAC Laplace approximation (Ritter et al., 2018), dropout (Srivastava et al., 2014) and test-time data augmentation (Krizhevsky et al., 2009). These techniques were chosen to cover a diverse set of approaches keeping their predictive performance in mind.

All these techniques can be summarized as distributions $q_m(\omega)$ over parameters ω of computation graphs, where m stands for the technique. During testing, one can average the predictions across parameters $\omega \sim q_m(\omega)$ to approximate the predictive distribution

$$\hat{p}(y_i | x_i) \approx \int p(y_i | x_i, \omega) q_m(\omega) d\omega \simeq \frac{1}{K} \sum_{k=1}^K p(y_i | x_i, \omega_k), \quad \omega_k \sim q_m(\omega) \quad (3)$$

For example, a deep ensemble of S networks can be represented in this form as a mixture of S Dirac's deltas $q_{DE}(\omega) = \frac{1}{S} \sum_{s=1}^S \delta(\omega - \omega_s)$, centered at independently trained snapshots ω_s . Similarly, a Bayesian neural network with a fully-factorized Gaussian approximate posterior distribution over the weight matrices and convolutional kernels ω is represented as $q_{VI}(\omega) = \mathcal{N}(\omega | \mu, \text{diag}(\sigma^2))$, μ and σ^2 being the optimal variational means and variances respectively.

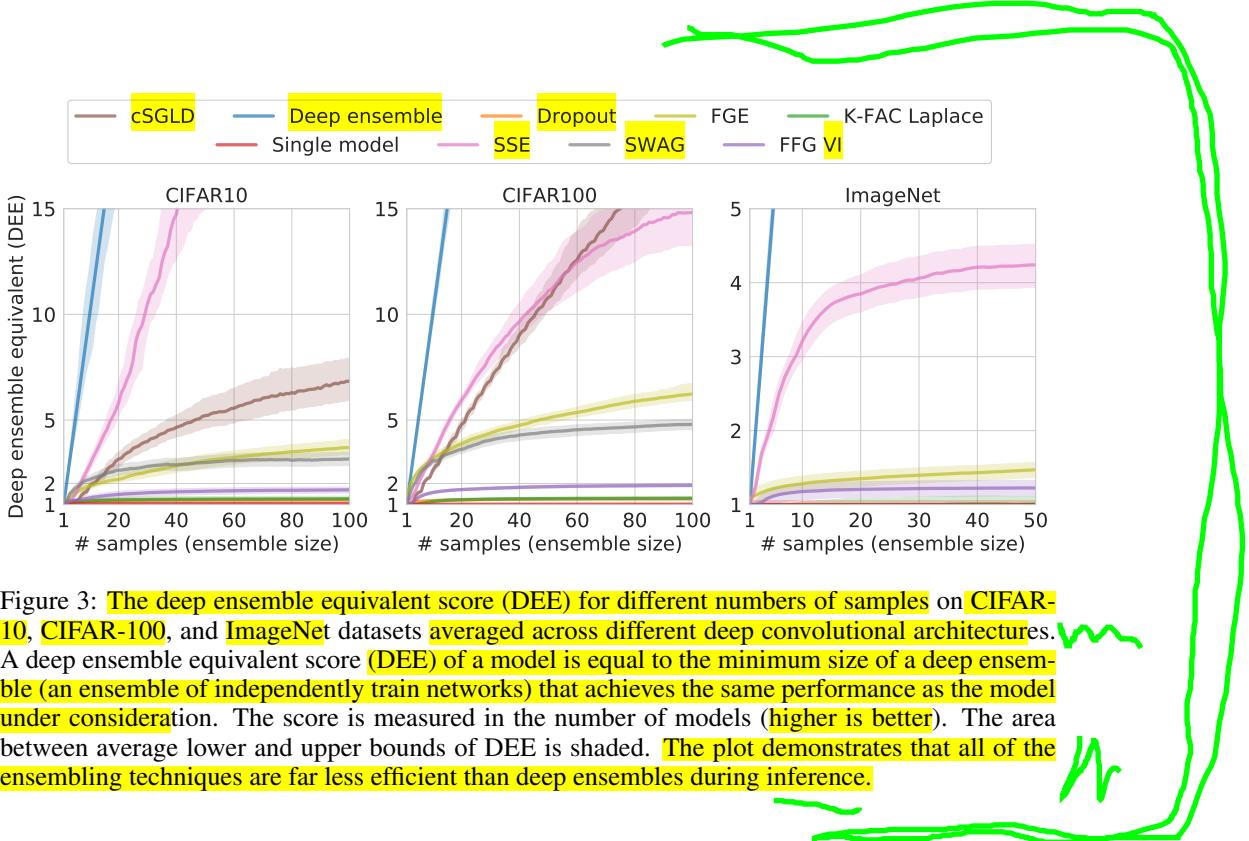
If one considers data augmentation as a part of the computational graph, it can be parameterized by the coordinates of the random crop and the flag for whether to flip the image horizontally or not. Sampling from the corresponding $q_{aug}(\omega)$ would generate different ways to augment the data during inference. However, as data augmentation is present by default during the training of all other mentioned ensembling techniques, it is suitable to study it in combination with these methods and not as a separate ensembling technique. We perform such an evaluation in Section 4.3.

Typically, the approximation (equation 3) requires K independent forward passes through a neural network, making the test-time budget directly comparable across all methods.

4.1 DEEP ENSEMBLE EQUIVALENT

Most ensembling techniques under consideration are either bounded to a single mode, or provide positively correlated samples. Deep ensemble, on the other hand, is a simple technique that provides independent samples from different modes of the loss landscape, which, intuitively, should result in a better ensemble. Therefore deep ensembles can be considered a strong baseline for the performance of other ensembling techniques given a fixed test-time computation budget.

Comparing the performance of ensembling techniques is, however, a challenging problem. Different models on different datasets achieve different values of metrics; their dependence on the number of



samples is non-trivial, and varies depending on a specific model and dataset. Values of the metrics are thus lacking in interpretability as the gain in performance has to be compared against a model- and dataset-specific baseline.

Aiming to introduce perspective and interpretability in our study, we introduce the *deep ensemble equivalent* score that employs deep ensembles to measure the performance of other ensembling techniques. Specifically, the deep ensemble equivalent score answers the following question:

What size of deep ensemble yields the same performance as a particular ensembling method?

Following the insights from the previous sections, we base the deep ensemble equivalent on the calibrated log-likelihood (CLL). Formally speaking, we define the deep ensemble equivalent (DEE) for an ensembling method m and its upper and lower bounds as follows:

$$\text{DEE}_m(k) = \min \left\{ l \in \mathbb{R}, l \geq 1 \mid \text{CLL}_{DE}^{\text{mean}}(l) \geq \text{CLL}_m^{\text{mean}}(k) \right\}, \quad (4)$$

$$\text{DEE}_m^{\text{upper}}(k) = \min \left\{ l \in \mathbb{R}, l \geq 1 \mid \text{CLL}_{DE}^{\text{mean}}(l) \mp \text{CLL}_{DE}^{\text{std}}(l) \geq \text{CLL}_m^{\text{mean}}(k) \right\}, \quad (5)$$

where $\text{CLL}_m^{\text{mean/std}}(l)$ are the mean and the standard deviation of the calibrated log-likelihood achieved by an ensembling method m with l samples. We compute $\text{CLL}_{DE}^{\text{mean}}(l)$ and $\text{CLL}_{DE}^{\text{std}}(l)$ for natural numbers $l \in \mathbb{N}_{>0}$ and use linear interpolation to define them for real values $l \geq 1$. In the following plots we report $\text{DEE}_m(k)$ for different methods m with different numbers of samples k , and shade the area between the respective lower and upper bounds $\text{DEE}_m^{\text{lower}}(k)$ and $\text{DEE}_m^{\text{upper}}(k)$.

4.2 EXPERIMENTS

We compute the deep ensemble equivalent (DEE) of various ensembling techniques for four popular deep architectures: VGG16 (Simonyan & Zisserman, 2014), PreResNet110/164 (He et al., 2016), and WideResNet28x10 (Zagoruyko & Komodakis, 2016) on CIFAR-10/100 datasets (Krizhevsky et al., 2009), and ResNet50 (He et al., 2016) on ImageNet dataset (Russakovsky et al., 2015). We use PyTorch (Paszke et al., 2017) for implementation of these models, building upon available public implementations. Our implementation closely matches the quality of methods that has been reported

(cSGLD and SSE both employ a cyclical learn

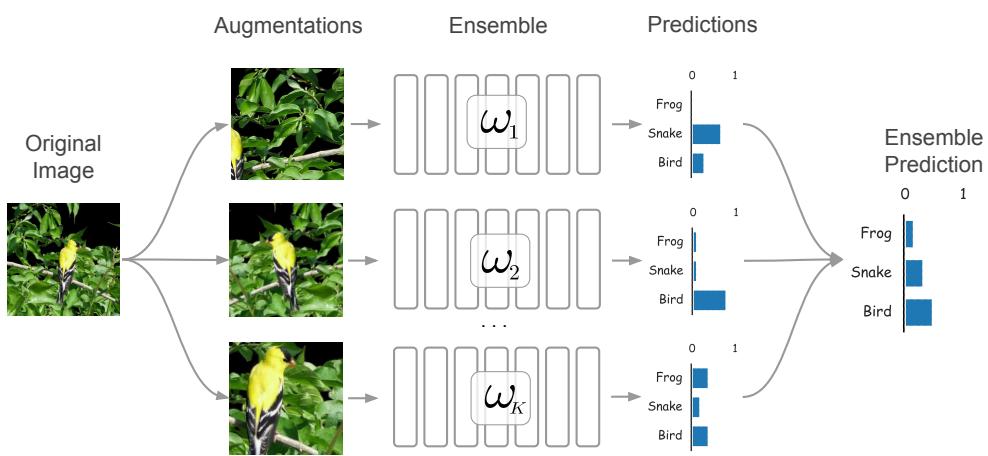


Figure 4: An illustration of test-time augmentation (TTA) for an ensemble. We apply every member of an ensemble to a separate random augmentation of an image. The predictions of all members are averaged to produce a final prediction of an ensemble. In our experiments, TTA leads to a significant boost of the performance for most of the ensembling techniques on ImageNet with a sufficient computational budget (see Figure 5).

in original works. Technical details on training, hyperparameters and implementations can be found in Appendix B. The source code and all computed metrics are available on GitHub¹.

As one can see on Figure 3, ensembling methods clearly fall into three categories. SSE and cSGLD outperform all other techniques except deep ensembles and enjoy a near-linear scaling of DEE with the number of samples on CIFAR datasets. The investigation of weight-space trajectories of cSGLD and SSE (Huang et al., 2017; Zhang et al., 2019) suggests that these methods can efficiently explore different modes of the loss landscape. In terms of the deep ensemble equivalent, these methods do not saturate unlike other methods that are bound to a single mode. We found SSE to still saturate on ImageNet. This is likely due to suboptimal hyperparameters of the cyclic learning rate schedule. More verbose results are presented in Figures 11–13 and in Table 5 and Table 8 in Appendix C.

In our experiments SSE typically outperforms cSGLD. This is mostly due to the fact that SSE has a much larger training budget. The cycle lengths and learning rates of SSE and cSGLD are comparable, however, SSE collects one snapshot per cycle while cSGLD collects three snapshots. This makes samples from SSE less correlated with each other while increasing the training budget threefold. Both SSE and cSGLD can be adjusted to obtain a different trade-off between the training budget and the DEE-to-samples ratio. We reused the schedules provided in the original papers (Huang et al., 2017; Zhang et al., 2019).

Being more “local” methods, FGE and SWAG perform worse than SSE and cSGLD, but still significantly outperform “single-snapshot” methods like dropout, K-FAC Laplace approximation and variational inference. We hypothesize that by covering a single mode with a set of snapshots, FGE and SWAG provide a better fit for the local geometry than models trained as stochastic computation graphs. This implies that the performance of FGE and SWAG should be achievable by single-snapshot methods. However, one might need more elaborate posterior approximations and better inference techniques in order to match the performance of FGE and SWAG by training a stochastic computation graph end-to-end (as opposed to SWAG that constructs a stochastic computation graph post-hoc).

The deep ensemble equivalent curves allow us to notice the common behaviour of different methods, e.g. the relation between deep ensembles, snapshot methods, advanced local methods and single-snapshot local methods. They also allow us to notice inconsistencies that may indicate a suboptimal choice of hyperparameters. For example, we find that SSE on ImageNet quickly saturates, unlike SSE on CIFAR datasets (Figure 3). This may indicate that the hyperparameters used on ImageNet

¹Source code: <https://github.com/bayesgroup/pytorch-ensembles>

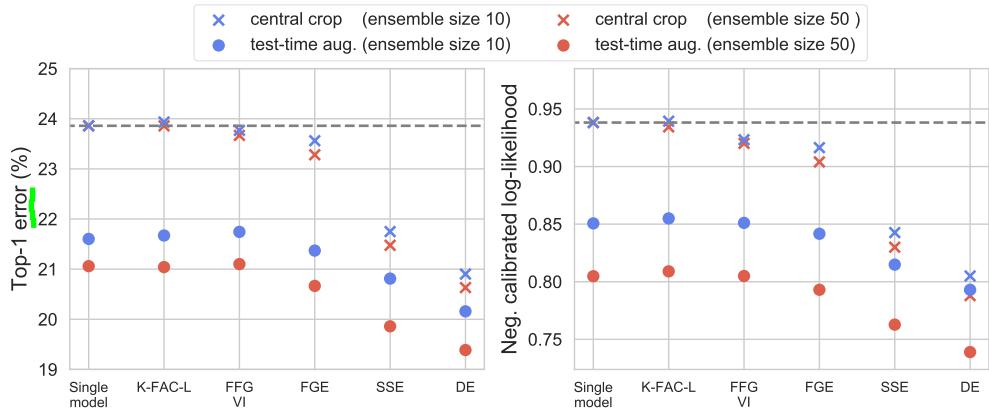


Figure 5: The negative calibrated log-likelihood (lower is better) for different ensembling techniques on ImageNet. We report performance for two regimes. *Central-crop evaluation* ($\times \times$) means every member of an ensemble is applied to a central crop of an image, and *test-time data augmentation* ($\bullet \bullet$) means each member of the ensemble is applied to a separate random augmentation of the image. Given a sufficiently large number of samples, test-time data augmentation significantly improves all ensembling techniques.

are not good enough for efficient coverage of different modes of the loss landscape. We also find that SSE on WideResNet on CIFAR-10 achieves a DEE score of 100 on approx. 70 samples (Figure 12). This may indicate that the members of the deep ensemble for this dataset-architecture pair are underfitted and may benefit from longer training or a different learning rate schedule. Such inconsistencies might be more difficult to spot using plain calibrated log-likelihood plots.

4.3 TEST-TIME DATA AUGMENTATION IMPROVES ENSEMBLES FOR FREE

Data augmentation is a time-honored technique that is widely used in deep learning, and is a crucial component for training modern DNNs. Test-time data augmentation has been used for a long time to improve the performance of convolutional networks. For example, multi-crop evaluation has been a standard procedure for the ImageNet challenge (Simonyan & Zisserman, 2014; Szegedy et al., 2015; He et al., 2016). It is, however, often overlooked in the literature on ensembling techniques in deep learning. In this section, we study the effect of test-time data augmentation on the aforementioned ensembling techniques. To keep the test-time computation budget the same, we sample one random augmentation for each member of an ensemble. Figure 5 reports the calibrated log-likelihood on combination of ensembles and test-time data augmentation for ImageNet. Other metrics and results on CIFAR-10/100 datasets are reported in Appendix C. We have used the standard data augmentation: random horizontal flips and random padded crops for CIFAR-10/100 datasets, and random horizontal flips and random resized crops for ImageNet (see more details in Appendix B).

Test-time data augmentation (Figure 4) consistently improves most ensembling methods, especially on ImageNet, where we see a clear improvement across all methods (Figure 5 and Table 7). The performance gain for powerful ensembles (deep ensembles, SSE and cSGLD) on CIFAR datasets is not as dramatic (Figures 14–15 and Table 4). This is likely due to the fact that CIFAR images are small, making data augmentation limited, whereas images from ImageNet allow for a large number of diverse samples of augmented images. On the other hand, while the performance of “single-snapshot” methods (e.g. variational inference, K-FAC Laplace and dropout) is improved significantly, they perform approximately as good as an augmented version of a *single* model across all datasets.

Interestingly, test-time data augmentation on ImageNet improves accuracy but decreases the uncalibrated log-likelihood of deep ensembles (Table 7 in Appendix C). Test-time data augmentation breaks the nearly optimal temperature of deep ensembles and requires temperature scaling to reveal the actual performance of the method, as discussed in Section 3.1. The experiment demonstrates

that ensembles may be highly miscalibrated by default while still providing superior predictive performance after calibration.

5 DISCUSSION & CONCLUSION

We have explored the field of in-domain uncertainty estimation and performed an extensive evaluation of modern ensembling techniques. Our main findings can be summarized as follows:

- Temperature scaling is a must even for ensembles. While ensembles generally have better calibration out-of-the-box, they are not calibrated perfectly and can benefit from the procedure. A comparison of log-likelihoods of different ensembling methods without temperature scaling might not provide a fair ranking, especially if some models happen to be miscalibrated.
- Many common metrics for measuring in-domain uncertainty are either unreliable (ECE and analogues) or cannot be used to compare different methods (AUC-ROC, AUC-PR for misclassification detection; accuracy-confidence curves). In order to perform a fair comparison of different methods, one needs to be cautious of these pitfalls.
- Many popular ensembling techniques require dozens of samples for test-time averaging, yet are essentially equivalent to a handful of independently trained models. Deep ensembles dominate other methods given a fixed test-time budget. The results indicate, in particular, that exploration of different modes in the loss landscape is crucial for good predictive performance.
- Methods that are stuck in a single mode are unable to compete with methods that are designed to explore different modes of the loss landscape. Would more elaborate posterior approximations and better inference techniques shorten this gap?
- Test-time data augmentation is a surprisingly strong baseline for in-domain uncertainty estimation. It can significantly improve other methods without increasing training time or model size since data augmentation is usually already present during training.

Our takeaways are aligned with the take-home messages of Ovadia et al. (2019) that relate to in-domain uncertainty estimation. We also observe a stable ordering of different methods in our experiments, and observe that deep ensembles with few members outperform methods based on stochastic computation graphs.

A large number of unreliable metrics inhibits a fair comparison of different methods. Because of this, we urge the community to aim for more reliable benchmarks in the numerous setups of uncertainty estimation.

ACKNOWLEDGMENTS

Dmitry Vetrov and Dmitry Molchanov were supported by the Russian Science Foundation grant no. 19-71-30020. This research was supported in part through computational resources of HPC facilities at NRU HSE.

REFERENCES

- Andrei Atanov, Arsenii Ashukha, Dmitry Molchanov, Kirill Neklyudov, and Dmitry Vetrov. Uncertainty estimation via stochastic batch normalization. In *International Symposium on Neural Networks*, pp. 261–269. Springer, 2019.
- Anoop Korattikara Balan, Vivek Rathod, Kevin P Murphy, and Max Welling. Bayesian dark knowledge. In *Advances in Neural Information Processing Systems*, pp. 3438–3446, 2015.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. Weight uncertainty in neural networks. *arXiv preprint arXiv:1505.05424*, 2015.
- Glenn W Brier. Verification of forecasts expressed in terms of probability. *Monthly weather review*, 78(1):1–3, 1950.

- Jacob Buckman, Danijar Hafner, George Tucker, Eugene Brevdo, and Honglak Lee. Sample-efficient reinforcement learning with stochastic ensemble value expansion. In *Advances in Neural Information Processing Systems*, pp. 8224–8234, 2018.
- Kurtland Chua, Roberto Calandra, Rowan McAllister, and Sergey Levine. Deep reinforcement learning in a handful of trials using probabilistic dynamics models. In *Advances in Neural Information Processing Systems*, pp. 4754–4765, 2018.
- Yufei Cui, Wuguannan Yao, Qiao Li, Antoni B Chan, and Chun Jason Xue. Accelerating monte carlo bayesian inference via approximating predictive uncertainty over simplex. *arXiv preprint arXiv:1905.12194*, 2019.
- Yukun Ding, Jinglan Liu, Jinjun Xiong, and Yiyu Shi. Evaluation of neural network uncertainty estimation with application to resource-constrained platforms. *arXiv preprint arXiv:1903.02050*, 2019.
- Yarin Gal. *Uncertainty in deep learning*. PhD thesis, PhD thesis, University of Cambridge, 2016.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pp. 1050–1059, 2016.
- Timur Garipov, Pavel Izmailov, Dmitrii Podoprikhin, Dmitry P Vetrov, and Andrew G Wilson. Loss surfaces, mode connectivity, and fast ensembling of dnns. In *Advances in Neural Information Processing Systems*, pp. 8789–8798, 2018.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 1321–1330. JMLR. org, 2017.
- Fredrik K Gustafsson, Martin Danelljan, and Thomas B Schön. Evaluating scalable bayesian deep learning methods for robust computer vision. *arXiv preprint arXiv:1906.01620*, 2019.
- Lars Kai Hansen and Peter Salamon. Neural network ensembles. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, (10):993–1001, 1990.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.
- Geoffrey Hinton, Oriol Vinyals, and Jeff Dean. Distilling the knowledge in a neural network. *arXiv preprint arXiv:1503.02531*, 2015.
- Gao Huang, Yixuan Li, Geoff Pleiss, Zhuang Liu, John E Hopcroft, and Kilian Q Weinberger. Snapshot ensembles: Train 1, get m for free. *arXiv preprint arXiv:1704.00109*, 2017.
- Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. *arXiv preprint arXiv:1502.03167*, 2015.
- Pavel Izmailov, Dmitrii Podoprikhin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018.
- Durk P Kingma, Tim Salimans, and Max Welling. Variational dropout and the local reparameterization trick. In *Advances in Neural Information Processing Systems*, pp. 2575–2583, 2015.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. Technical report, Citeseer, 2009.
- Ananya Kumar, Percy S Liang, and Tengyu Ma. Verified uncertainty calibration. In *Advances in Neural Information Processing Systems*, pp. 3787–3798, 2019.
- Balaji Lakshminarayanan, Alexander Pritzel, and Charles Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In *Advances in Neural Information Processing Systems*, pp. 6402–6413, 2017.

- Christos Louizos and Max Welling. Multiplicative normalizing flows for variational bayesian neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 2218–2227. JMLR. org, 2017.
- Wesley Maddox, Timur Garipov, Pavel Izmailov, Dmitry Vetrov, and Andrew Gordon Wilson. A simple baseline for bayesian uncertainty in deep learning. *arXiv preprint arXiv:1902.02476*, 2019.
- Andrey Malinin and Mark Gales. Predictive uncertainty estimation via prior networks. In *Advances in Neural Information Processing Systems*, pp. 7047–7058, 2018.
- James Martens and Roger Grosse. Optimizing neural networks with kronecker-factored approximate curvature. In *International conference on machine learning*, pp. 2408–2417, 2015.
- Dmitry Molchanov, Arsenii Ashukha, and Dmitry Vetrov. Variational dropout sparsifies deep neural networks. In *Proceedings of the 34th International Conference on Machine Learning-Volume 70*, pp. 2498–2507. JMLR. org, 2017.
- Marcin Możejko, Mateusz Susik, and Rafał Karczewski. Inhibited softmax for uncertainty estimation in neural networks. *arXiv preprint arXiv:1810.01861*, 2018.
- Malik Sajjad Ahmed Nadeem, Jean-Daniel Zucker, and Blaise Hanczar. Accuracy-rejection curves (arcs) for comparing classification methods with a reject option. In *Machine Learning in Systems Biology*, pp. 65–81, 2009.
- Mahdi Pakdaman Naeini, Gregory Cooper, and Milos Hauskrecht. Obtaining well calibrated probabilities using bayesian binning. In *Twenty-Ninth AAAI Conference on Artificial Intelligence*, 2015.
- Jeremy Nixon, Mike Dusenberry, Linchuan Zhang, Ghassen Jerfel, and Dustin Tran. Measuring calibration in deep learning. 2019.
- Yaniv Ovadia, Emily Fertig, Jie Ren, Zachary Nado, D Sculley, Sebastian Nowozin, Joshua V Dillon, Balaji Lakshminarayanan, and Jasper Snoek. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. *arXiv preprint arXiv:1906.02530*, 2019.
- Adam Paszke, Sam Gross, Soumith Chintala, Gregory Chanan, Edward Yang, Zachary DeVito, Zeming Lin, Alban Desmaison, Luca Antiga, and Adam Lerer. Automatic differentiation in PyTorch. In *NIPS Autodiff Workshop*, 2017.
- Joaquin Quinonero-Candela, Carl Edward Rasmussen, Fabian Sinz, Olivier Bousquet, and Bernhard Schölkopf. Evaluating predictive uncertainty challenge. In *Machine Learning Challenges Workshop*, pp. 1–27. Springer, 2005.
- Hippolyt Ritter, Aleksandar Botev, and David Barber. A scalable laplace approximation for neural networks. 2018.
- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- Burr Settles. Active learning. *Synthesis Lectures on Artificial Intelligence and Machine Learning*, 6(1):1–114, 2012.
- Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. *arXiv preprint arXiv:1409.1556*, 2014.
- Leslie N Smith and Nicholay Topin. Super-convergence: Very fast training of neural networks using large learning rates. In *Artificial Intelligence and Machine Learning for Multi-Domain Operations Applications*, volume 11006, pp. 1100612. International Society for Optics and Photonics, 2019.
- Casper Kaae Sønderby, Tapani Raiko, Lars Maaløe, Søren Kaae Sønderby, and Ole Winther. How to train deep variational autoencoders and probabilistic ladder networks. In *33rd International Conference on Machine Learning (ICML 2016)*, 2016.

- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- Christian Szegedy, Wojciech Zaremba, Ilya Sutskever, Joan Bruna, Dumitru Erhan, Ian Goodfellow, and Rob Fergus. Intriguing properties of neural networks. *arXiv preprint arXiv:1312.6199*, 2013.
- Christian Szegedy, Wei Liu, Yangqing Jia, Pierre Sermanet, Scott Reed, Dragomir Anguelov, Dumitru Erhan, Vincent Vanhoucke, and Andrew Rabinovich. Going deeper with convolutions. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 1–9, 2015.
- Mattias Teye, Hossein Azizpour, and Kevin Smith. Bayesian uncertainty estimation for batch normalized deep networks. In *International Conference on Machine Learning (ICML)*, 2018.
- Marcin B Tomczak, Siddharth Swaroop, and Richard E Turner. Neural network ensembles and variational inference revisited. In *1st Symposium on Advances in Approximate Bayesian Inference*, pp. 1–11, 2018.
- Linh Tran, Bastiaan S Veeling, Kevin Roth, Jakub Swiatkowski, Joshua V Dillon, Jasper Snoek, Stephan Mandt, Tim Salimans, Sebastian Nowozin, and Rodolphe Jenatton. Hydra: Preserving ensemble diversity for model distillation. *arXiv preprint arXiv:2001.04694*, 2020.
- Karen Ullrich, Edward Meeds, and Max Welling. Soft weight-sharing for neural network compression. *arXiv preprint arXiv:1702.04008*, 2017.
- Juozas Vaicenavicius, David Widmann, Carl Andersson, Fredrik Lindsten, Jacob Roll, and Thomas B Schon. Evaluating model calibration in classification. *arXiv preprint arXiv:1902.06977*, 2019.
- Kuan-Chieh Wang, Paul Vicol, James Lucas, Li Gu, Roger Grosse, and Richard Zemel. Adversarial distillation of bayesian neural network posteriors. In *International Conference on Machine Learning*, pp. 5177–5186, 2018.
- Sida Wang and Christopher Manning. Fast dropout training. In *international conference on machine learning*, pp. 118–126, 2013.
- Max Welling and Yee W Teh. Bayesian learning via stochastic gradient langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 681–688, 2011.
- Anqi Wu, Sebastian Nowozin, Edward Meeds, Richard E Turner, José Miguel Hernández-Lobato, and Alexander L Gaunt. Deterministic variational inference for robust bayesian neural networks. 2018.
- Sergey Zagoruyko. 92.45 on cifar-10 in torch, 2015. URL <http://torch.ch/blog/2015/07/30/cifar.html>.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint arXiv:1605.07146*, 2016.
- Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- Ruqi Zhang, Chunyuan Li, Jianyi Zhang, Changyou Chen, and Andrew Gordon Wilson. Cyclical stochastic gradient mcmc for bayesian deep learning. *arXiv preprint arXiv:1902.03932*, 2019.

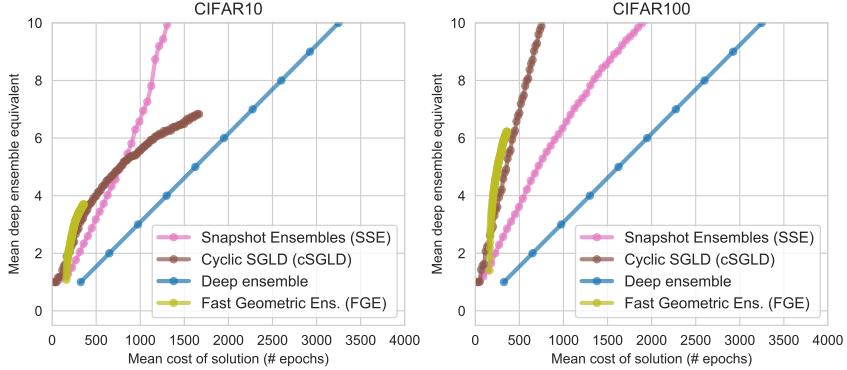


Figure 6: The mean cost of training of an ensemble vs. the mean deep ensemble equivalent score. Each marker on the plot denotes one snapshot of weights.

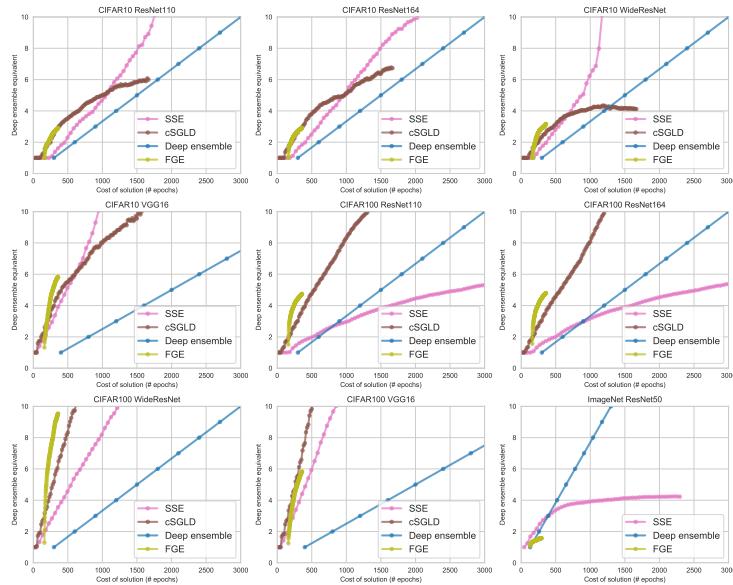


Figure 7: Cost of training of an ensemble vs. its quality (DEE). Each marker on a plot denotes one snapshot of weights.

A IS "EFFICIENT" TRAINING OF ENSEMBLES EFFICIENT AT ALL?

Yes! If you care most about training time efficiency. All snapshot based methods (SSE, cSGLD and FGE) on average (Figure 6) tend to achieve the same performance as deep ensembles **using 2-5 \times less training epochs on CIFAR datasets.**

The gain comes at a cost of inference efficiency and memory consumption. "Efficient" snapshot-based methods require to store much more samples of weights compared to deep ensembles **making inference significantly more expensive (up to $\times 25$) given the same predictive performance.**

You need to get lucky with hyperparameter choice. While on average "efficient" snapshot-based methods require less training resources, they might be completely inefficient if your hyperparameter choice is sub-optimal (see Figure 7). **Such hyperparameters as the maximum learning rate, the length of learning rate decay cycle, the snapshot saving schedule can all impact the performance significantly.**

This analysis assumes that we use only conventional training procedure. Many models can most likely be trained, stored and executed much more efficiently with methods like super-convergence, stochastic weight averaging, compression, distillation and others. These methods are out of the scope of the paper but are interesting topics for future research. The current choice of hyperparameters may also not be optimal. We reuse the hyperparameters used in the original papers.

B EXPERIMENTAL DETAILS

Implementations of deep ensembles, SWAG, FGE and K-FAC Laplace are heavily based on the original PyTorch implementations of stochastic weight averaging (SWA)² and SWAG³. Implementations of cyclical MCMC and snapshot ensembles are based on the original implementation of cyclical MCMC⁴. We hypothesize that the optimal hyperparameters of ensembling methods may vary widely depending on the computational budget and the number of samples in the ensemble. Searching for the optimal values for each configuration is outside the scope of this paper so we stick to the originally proposed hyperparameters whenever possible.

Implied probabilistic model Conventional neural networks for classification are usually trained using the average cross-entropy loss function with weight decay regularization hidden inside an optimizer in a deep learning framework like PyTorch. The underlying optimization problem can be written as follows:

$$L(w) = -\frac{1}{N} \sum_{i=1}^N \log \hat{p}(y_i^* | x_i, w) + \frac{\lambda}{2} \|w\|^2 \rightarrow \min_w, \quad (6)$$

where $\{(x_i, y_i^*)\}_{i=1}^N$ is the training dataset of N objects x_i with corresponding labels y_i^* , λ is the weight decay scale and $\hat{p}(j | x_i, w)$ denotes the probability that a neural network with parameters w assigns to class j when evaluated on object x_i .

The cross-entropy loss defines a likelihood function $p(y^* | x, w)$ and weight decay regularization, or L_2 regularization corresponding to a certain Gaussian prior distribution $p(w)$. The whole optimization objective then corresponds to the *maximum a posteriori* inference in the following probabilistic model:

$$p(y^*, w | x) = p(y^* | x, w)p(w), \quad (7)$$

$$\log p(y^* | x, w) = \log \prod_{i=1}^N p(y_i^* | x_i, w) = \sum_{i=1}^N \log \hat{p}(y_i^* | x_i, w), \quad (8)$$

$$\log p(w) = \frac{-N\lambda}{2} \|w\|^2 + \text{const} \iff p(w) = \mathcal{N}(w | 0, (N\lambda)^{-1}I) \quad (9)$$

In order to make the results comparable across all ensembling techniques, we used the same probabilistic model for all methods, choosing fixed weight decay parameters for each architecture. We used the softmax-based likelihood for all models. We also use the fully-factorized zero-mean Gaussian prior distribution with variances $\sigma^2 = (N\lambda)^{-1}$ where the number of objects N and the weight decay scale λ are dictated by the particular datasets and neural architectures as defined in the following paragraph.

Conventional networks To train a single network on CIFAR-10/100, we used SGD with batch size of 128, momentum 0.9 and model-specific parameters, i.e. the initial learning rate (lr_{init}), the weight decay coefficient (wd), and the number of optimization epochs (epoch). Specific hyperparameters are shown in Table 1. The models were trained with a unified learning rate scheduler that is shown in equation 10. All models have been trained using data augmentation that consists of horizontal flips and a random crop of 32 pixels with a padding of 4 pixels⁵. The standard data normalization has also been applied. Weight decays, initial learning rates, and the learning rate scheduler were taken from (Garipov et al., 2018) paper. Compared with hyperparameters of (Garipov et al., 2018), we increased the number of optimization epochs since we found that all models were underfitted. While the original WideResNet28x10 network includes a number of dropout layers with $p = 0.3$ and is trained for 200 epochs, we find that the WideResNet28x10 underfits in this setting and requires a longer training. Therefore, we used $p = 0$ which reduces training time while bearing

²<https://github.com/timgaripov/swa>

³https://github.com/wjmaddox/swa_gaussian

⁴<https://github.com/ruqizhang/csgmcmc/tree/master/experiments>

⁵`Compose([RandomHorizontalFlip(), RandomCrop(32, padding=4)])`

Model	lr_{init}	epochs	wd
VGG	0.05	400	5e-4
PreResNet110	0.1	300	3e-4
PreResNet164	0.1	300	3e-4
WideResNet28x10	0.1	300	5e-4

Table 1: Hyperparameters of models trained on CIFARs for single-model evaluation

no significant effect on final model performance in our experiments.

$$lr(i) = \begin{cases} lr_{init}, & i \in [0, 0.5 \cdot \text{epochs}] \\ lr_{init} \cdot (1.0 - 0.99 * (i/\text{epochs} - 0.5)/0.4), & i \in [0.5 \cdot \text{epochs}, 0.9 \cdot \text{epochs}] \\ lr_{init} \cdot 0.01, & \text{otherwise} \end{cases} \quad (10)$$

On ImageNet dataset we used ResNet50 with default hyperparameters taken from PyTorch examples⁶. Specifically, we used SGD with momentum 0.9, batch size of 256, initial learning rate 0.1, weight decay $1e - 4$. Training included data augmentation⁷ (scaling, random crops of size 224×224 , horizontal flips), normalization and learning rate scheduler $lr = lr_{init} \cdot 0.1^{\text{epoch}/30}$ where $/$ denotes integer division. We only deviated from standard parameters by increasing the number of training epochs from 90 to 130. Our models achieve top-1 error of 23.81 ± 0.15 that closely matches the accuracy of ResNet50 provided by PyTorch which is 23.85 ⁸. Training of one model on a single NVIDIA Tesla V100 GPU takes approximately 5.5 days.

Deep ensembles Deep ensembles (Lakshminarayanan et al., 2017) average the predictions across networks trained independently starting from different initializations. To obtain a deep ensemble we repeat the described procedure of training standard networks 128 times for all architectures on CIFAR-10 and CIFAR-100 datasets (1024 networks over all) and 50 times for ImageNet dataset. Every member of deep ensembles was trained with exactly the same hyperparameters as conventional models of the same architecture.

Dropout Binary dropout (or MC dropout) (Srivastava et al., 2014; Gal & Ghahramani, 2016) is one of the most widely known ensembling techniques. It involves putting a multiplicative Bernoulli noise with a parameter p over the activations of either a fully connected layer or a convolutional layer, averaging predictions of the network w.r.t. noise at test-time. **Dropout layers were applied to VGG and WideResNet networks on CIFAR-10 and CIFAR-100 datasets. Dropout for VGG was applied to fully connected layers with $p = 0.5$. Two dropout layers were applied: one before the first fully connected layer and one before the second one.** While the original version of VGG for CIFARs (Zagoruyko, 2015) exploits more dropout layers, we observed that any additional dropout layer deteriorates the performance of the model in either deterministic or stochastic mode. Dropout for WideResNet was applied in accordance with the original paper (Zagoruyko & Komodakis, 2016) with $p = 0.3$. Dropout usually increases the time needed to achieve convergence. Because of this, WideResNet networks with dropout were trained for 400 epochs instead of 300 epochs for deterministic case, and VGG networks have always been trained with dropout. All the other hyperparameters were the same as in the case of conventional models.

Variational Inference Variational Inference (VI) approximates the true posterior distribution over weights $p(w | Data)$ with a tractable variational approximation $q_\theta(w)$ by maximizing a so-called variational lower bound \mathcal{L} (eq. 11) w.r.t. the parameters of variational approximation θ . **We used a fully-factorized Gaussian approximation $q(w)$ and Gaussian prior distribution $p(w)$.**

$$\mathcal{L}(\theta) = \mathbb{E}_q \log p(y^* | x, w) - KL(q_\theta(w) || p(w)) \rightarrow \max_{\theta} \quad (11)$$

$$q(w) = \mathcal{N}(w | \mu, \text{diag}(\sigma^2)) \quad p(w) = N(w | 0, \text{diag}(\sigma_p^2)), \quad \text{where } \sigma_p^2 = (N \cdot \text{wd})^{-1} \quad (12)$$

⁶<https://github.com/pytorch/examples/tree/ee964a2/imagenet>

⁷`Compose([RandomResizedCrop(224), RandomHorizontalFlip()])`

⁸<https://pytorch.org/docs/stable/torchvision/models.html>

Architecture	Optimal noise scale			
	CIFAR-10	CIFAR-10-aug	CIFAR-100	CIFAR-100-aug
VGG16BN	0.042	0.042	0.100	0.100
PreResNet110	0.213	0.141	0.478	0.401
PreResNet164	0.120	0.105	0.285	0.225
WideResNet28x10	0.022	0.018	0.022	0.004

Table 2: Optimal noise scale for K-FAC Laplace for different datasets and architectures. For ResNet50 on ImageNet, the optimal scale found was 2.0 with test-time augmentation and 6.8 without test-time augmentation.

In the case of such a prior, the probabilistic model remains consistent with the conventional training which corresponds to MAP inference in the same probabilistic model. We used variational inference for both convolutional and fully-connected layers where variances of the weights were parameterized by $\log \sigma$. For fully-connected layers we applied the local reparameterization trick (LRT, (Kingma et al., 2015)).

While variational inference provides a theoretically grounded way to approximate the true posterior, it tends to underfit deep learning models in practice (Kingma et al., 2015). The following tricks are applied to deal with it: pre-training (Molchanov et al., 2017) or equivalently annealing of β (Sønderby et al., 2016), and scaling β down (Kingma et al., 2015; Ullrich et al., 2017).

During pre-training we initialize μ with a snapshot of weights of a pre-trained conventional model, and initialize $\log \sigma$ with a model-specific constant $\log \sigma_{init}$. The KL-divergence – except for the term corresponding to the weight decay – is scaled with a model-specific parameter β . The weight decay term is implemented as a part of the optimizer. We used a fact that the KL-divergence between two Gaussian distributions can be rewritten as two terms, one of which is equivalent to the weight decay regularization.

On CIFAR-10 and CIFAR-100 we used β equal to 1e-4 for VGG, ResNet100 and ResNet164 networks, and β equal to 1e-5 for WideResNet. Log-variance $\log \sigma_{init}$ was initialized with -5 for all models. Parameters μ were optimized with SGD in the same manner as in the case of conventional networks except that the initial learning rate lr_{init} was set to 1e-3. We used a separate Adam optimizer with a constant learning rate of 1e-3 to optimize log-variances of the weights $\log \sigma$. Pre-training was done for 300 epochs, and after that the remaining part of training was done for 100 epochs. On ImageNet we used $\beta = 1e-3$, $lr_{init} = 0.01$, $\log \sigma_{init} = -6$, and trained the model for 45 epochs after pre-training.

K-FAC Laplace The Laplace approximation uses the curvature information of the appropriately scaled loss function to construct a Gaussian approximation to the posterior distribution. Ideally, one would use the Hessian of the loss function as a covariance matrix and use the maximum a posteriori estimate w^{MAP} as a mean of the Gaussian approximation:

$$\log p(w | x, y^*) = \log p(y^* | x, w) + \log p(w) + \text{const} \quad (13)$$

$$w^{MAP} = \arg \max_w \log p(w | x, y^*); \quad \Sigma = -\nabla \nabla \log p(w | x, y^*) \quad (14)$$

$$p(w | x, y^*) \approx \mathcal{N}(w | w^{MAP}, \Sigma) \quad (15)$$

In order to keep the method scalable, we use the Fisher Information Matrix as an approximation to the true Hessian (Martens & Grosse, 2015). For K-FAC Laplace, we use the whole dataset to construct an approximation to the empirical Fisher Information Matrix, and use the π correction to reduce the bias (Ritter et al., 2018; Martens & Grosse, 2015). Following (Ritter et al., 2018), we find the optimal noise scale for K-FAC Laplace on a held-out validation set by averaging across five random initializations. We then reuse this scale for networks trained without a hold-out validation set. We report the optimal values of scales in Table 2. Note that the optimal scale is different depending on whether we use test-time data augmentation or not. Since the data augmentation also introduces some amount of additional noise, the optimal noise scale for K-FAC Laplace with data augmentation is lower.

Snapshot ensembles Snapshot ensembles (SSE) (Huang et al., 2017) is a simple example of an array of methods which collect samples from a training trajectory of a network in weight space to construct an ensemble. Samples are collected in a cyclical manner: during each cycle the learning rate goes from a large value to near-zero and snapshot of weights of the network is taken at the end of the cycle. SSE uses SGD with a cosine learning schedule defined as follows:

$$\alpha(t) = \frac{\alpha_0}{2} \left(\cos \left(\frac{\pi \operatorname{mod}(t-1, \lceil T/M \rceil)}{\lceil T/M \rceil} \right) + 1 \right), \quad (16)$$

where α_0 is the initial learning rate, T is the total number of training iterations and M is the number of cycles.

For all datasets and models hyperparameters from the original SSE paper are reused. For CIFAR-10/100 length of the cycle is 40 epochs, maximum learning rate is 0.2, batch size is 64. On ResNet50 and ImageNet length of the cycle is 45 epochs, maximum learning rate is 0.1, batch size is 256.

Cyclical SGLD Cyclical Stochastic Gradient Langevin Dynamics (cSGLD) (Zhang et al., 2019) is a state-of-the-art ensembling method for deep neural networks pertaining to stochastic Markov Chain Monte Carlo family of methods. It bears similarity to SSE, e.g. it employs SGD with a learning rate schedule described with the equation 16 and training is cyclical in the same manner. Its main differences from SSE are the introduction of gradient noise and the capturing of several snapshots per cycle, both of which can aid in sampling from posterior distribution over neural network weights efficiently.

Some parameters from the original paper are reused: length of cycle is 50 epochs, maximum learning rate is 0.5, batch size is 64. Number of epochs with gradient noise per cycle is 3 epochs. This was found to yield much higher predictive performance and better uncertainty estimation compared to the original paper choice of 10 epochs for CIFAR-10 and 3 epochs for CIFAR-100.

Finally, the results of cyclical Stochastic Gradient Hamiltonian Monte Carlo (SGHMC), (Zhang et al., 2019) which reportedly has marginally better performance compared with cyclical SGLD, could not be reproduced with any value of SGD momentum term. Because of this, we only include cyclical SGLD in our benchmark.

FGE Fast Geometric Ensembling (FGE) is an ensembling method that is similar to SSE in that it collects weight samples from a training trajectory to construct an ensemble. Its main differences from SSE are pretraining, short cycle length and a piecewise-linear learning rate schedule:

$$\alpha(i) = \begin{cases} (1 - 2t(i))\alpha_1 + 2t(i)\alpha_2 & 0 < t(i) \leq \frac{1}{2} \\ (2 - 2t(i))\alpha_2 + (2t(i) - 1)\alpha_1 & \frac{1}{2} < t(i) \leq 1 \end{cases}. \quad (17)$$

Hyperparameters of the original implementation of FGE are reused. Model pretraining is done with SGD for 160 epochs according to the standard learning rate schedule described in equation 10 with maximum learning rates from Table 1. After that, a desired number of FGE cycles is done with one snapshot per cycle collected. For VGG the learning rate is changed with parameters $\alpha_1 = 1e-2$, $\alpha_2 = 5e-4$ in a cycle with cycle length of 2 epochs. For other networks the learning rate is changed with parameters $\alpha_1 = 5e-2$, $\alpha_2 = 5e-4$ with cycle length of 4 epochs. Batch size is 128.

SWAG SWA-Gaussian (SWAG) (Maddox et al., 2019) is an ensembling method based on fitting a Gaussian distribution to model weights on the SGD training trajectory and sampling from this distribution to construct an ensemble.

Like FGE, SWAG has a pretraining stage which is done according to the standard learning rate schedule described in equation 10 with maximum learning rates from Table 1. After that, training continues with a constant learning rate of 1e-2 for all models except for PreResNet110 and PreResNet164 on CIFAR-100 where it continues with a constant learning rate of 5e-2 in accordance with the original paper. Rank of the empirical covariance matrix which is used for estimation of Gaussian distribution parameters is set to be 20.

C ADDITIONAL EXPERIMENTAL RESULTS

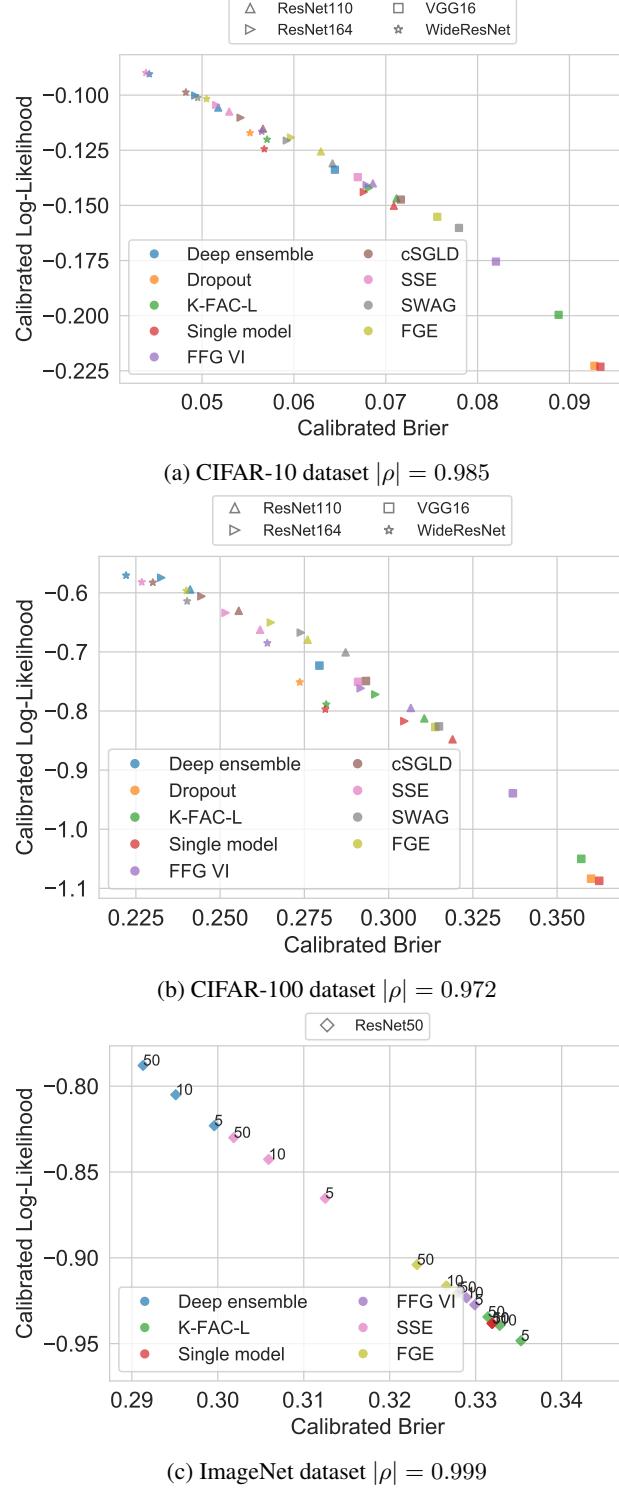


Figure 8: The average log-likelihood vs the Brier score on a test dataset for different ensemble methods on CIFAR-10 (a) and CIFAR-100 (b) and ImageNet (c) datasets. While not being equivalent, these metrics demonstrate a strong linear correlation. The correlation coefficient is denoted as ρ .

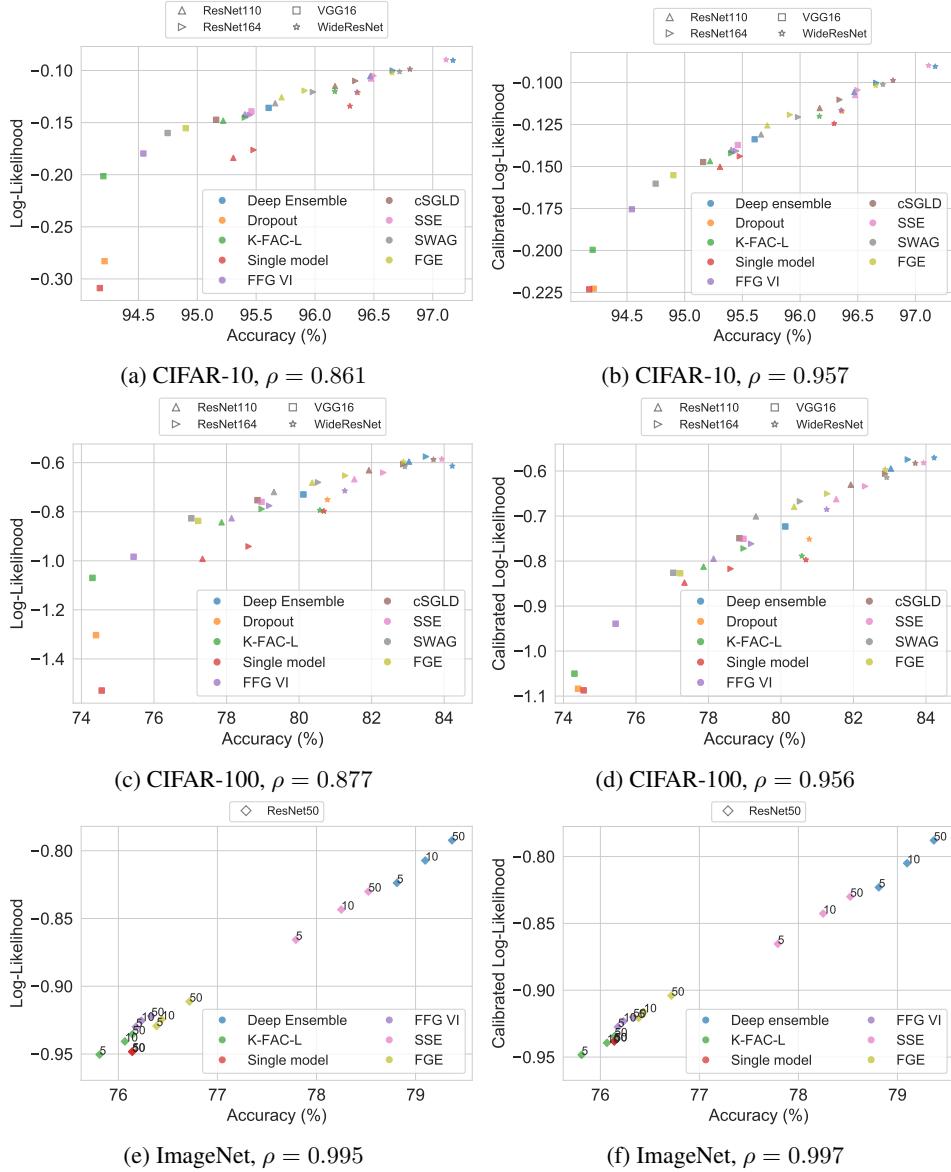


Figure 9: Log-likelihood vs accuracy for different ensembles before (a, c, e) and after (b, d, f) calibration. **Both plain log-likelihood and especially calibrated log-likelihood are highly correlated with accuracy.**

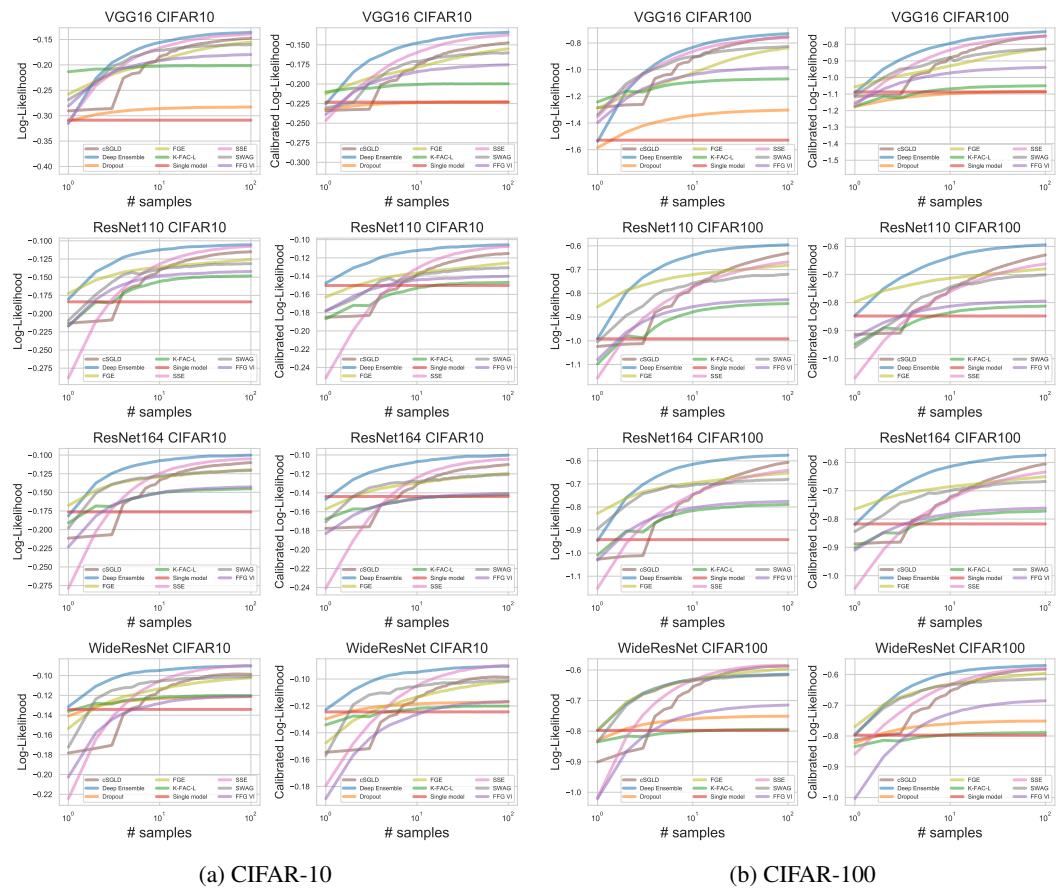


Figure 10: A side-by-side comparison of log-likelihood and calibrated log-likelihood on CIFAR-10 (a) and CIFAR-100 (b) datasets. On CIFAR-10 the performance of one network becomes close to dropout, variational inference (vi), and K-FAC Laplace approximation (kfac) after calibration on all models except VGG. On CIFAR-100 deep ensembles move to the first position in the ranking after calibration on WideResNet and VGG. See Section 3.1 for details on the calibrated log-likelihood.

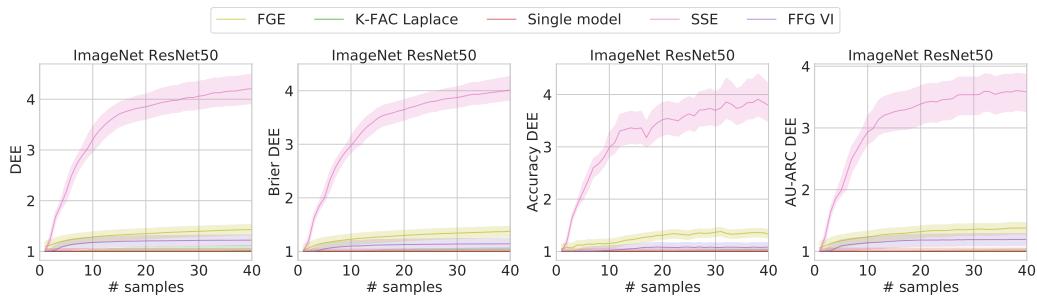


Figure 11: The deep ensemble equivalent of various ensembling techniques on ImageNet. Solid lines: mean DEE for different methods and architectures. Area between $\text{DEE}^{\text{lower}}$ and $\text{DEE}^{\text{upper}}$ is shaded. Columns 2–4 correspond to DEE based on other metrics, defined similarly to the log-likelihood-based DEE. The results are consistent for all metrics

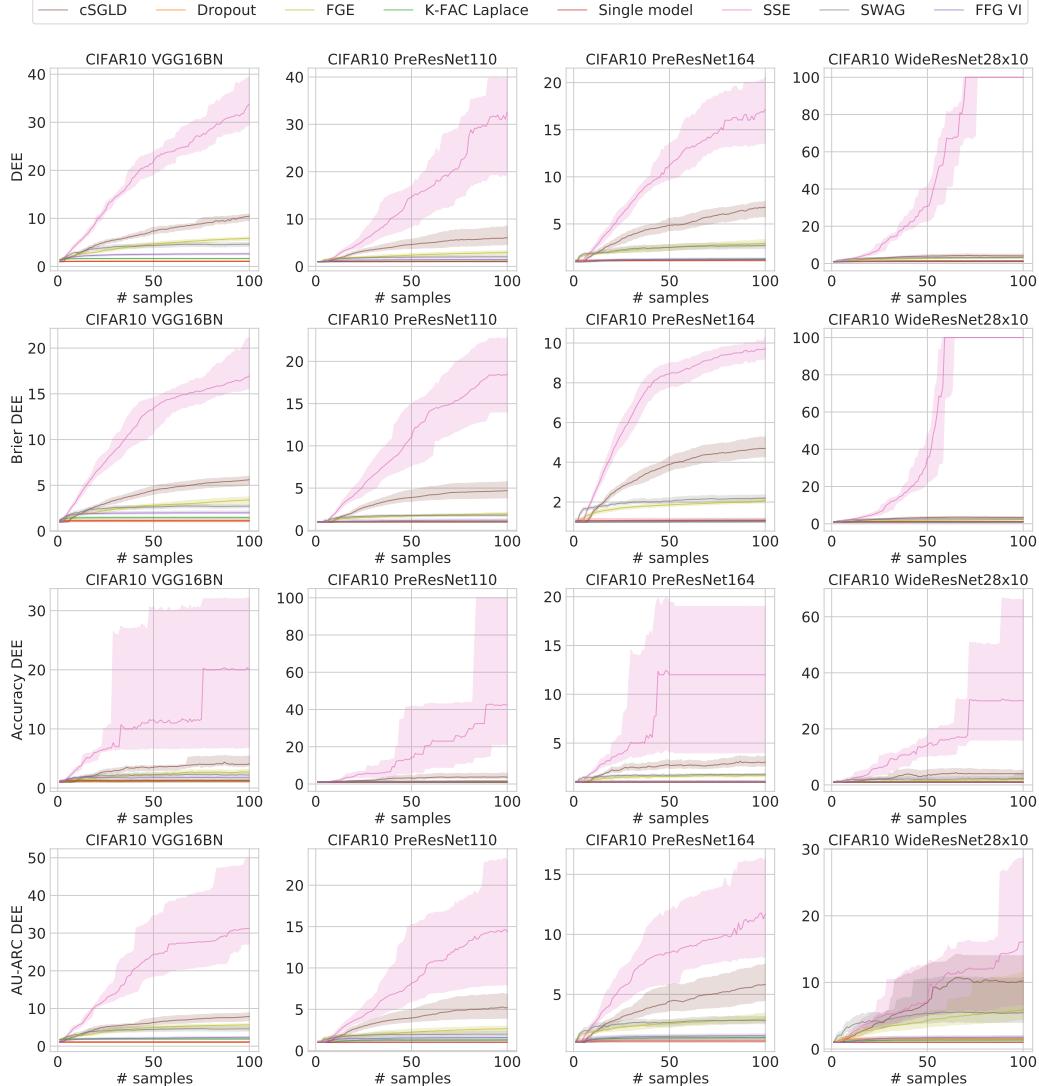


Figure 12: The deep ensemble equivalent of various ensembling techniques on **CIFAR-10**. Solid lines: mean DEE for different methods and architectures. Area between $\text{DEE}^{\text{lower}}$ and $\text{DEE}^{\text{upper}}$ is shaded. Lines 2–4 correspond to DEE based on other metrics, defined similarly to the log-likelihood-based DEE. Note that while the actual scale of DEE varies from metric to metric, the ordering of different methods and the overall behaviour of the lines remain the same.

SSE outperforms deep ensembles on CIFAR-10 on the WideResNet architecture. It possibly indicates that the cosine learning rate schedule and longer training of SSE are more suitable for this architecture than the piecewise-linear learning rate schedule and the number of epochs used in deep ensembles.

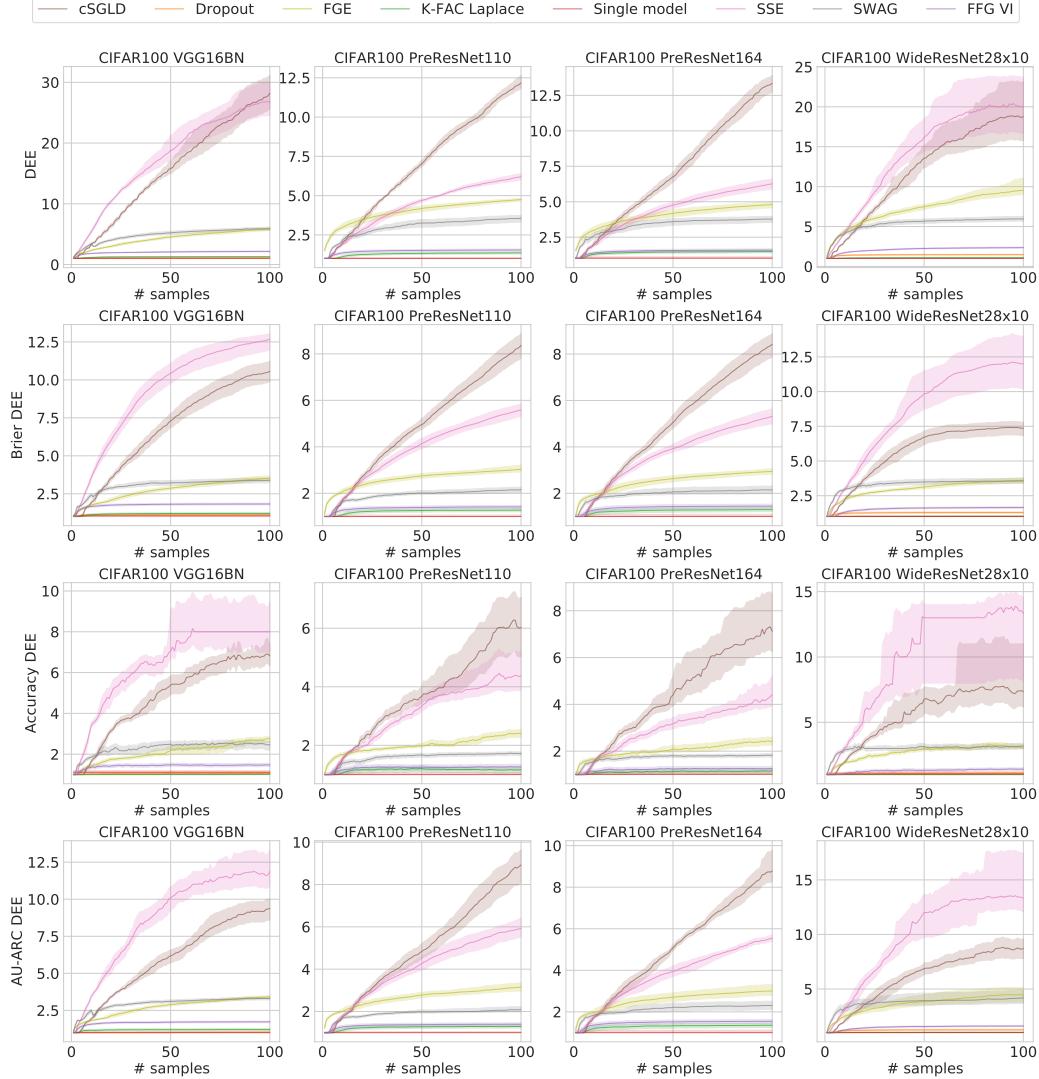


Figure 13: The deep ensemble equivalent of various ensembling techniques on **CIFAR-100**. Solid lines: mean DEE for different methods and architectures. Area between $\text{DEE}^{\text{lower}}$ and $\text{DEE}^{\text{upper}}$ is shaded. Lines 2–4 correspond to DEE based on other metrics, defined similarly to the log-likelihood-based DEE. Note that while the actual scale of DEE varies from metric to metric, the ordering of different methods and the overall behaviour of the lines remain the same.

Model	Method	Deep ensemble equivalent score				
		1 sample	5 samples	10 samples	50 samples	100 samples
VGG16 CIFAR-10	Deep Ensembles	1.0	5.0	10.0	50.0	100.0
	Snapshot Ensembles	1.0	2.6	4.7	21.9	33.8
	Cyclic SGLD	1.0	1.7	2.6	7.4	10.5
	SWA-Gaussian	1.0	2.2	2.9	4.4	4.6
	Fast Geometric Ens.	1.3	2.0	2.4	4.7	5.8
	Dropout	1.0	1.0	1.0	1.1	1.1
	Variational Inf. (FFG)	1.0	1.8	2.0	2.5	2.6
	KFAC-Laplace	1.4	1.6	1.6	1.6	1.6
ResNet110 CIFAR-10	Single model	1.1	1.1	1.1	1.1	1.1
	Deep Ensembles	1.0	5.0	10.0	48.0	94.6
	Snapshot Ensembles	1.0	1.0	2.0	14.6	32.6
	Cyclic SGLD	1.0	1.0	1.6	4.6	6.0
	SWA-Gaussian	1.0	1.3	1.5	1.9	2.0
	Fast Geometric Ens.	1.0	1.4	1.7	2.4	2.9
	Variational Inf. (FFG)	1.0	1.0	1.2	1.4	1.5
	KFAC-Laplace	1.0	1.0	1.0	1.0	1.1
ResNet164 CIFAR-10	Single model	1.0	1.0	1.0	1.0	1.0
	Deep Ensembles	1.0	5.0	10.0	43.7	100.0
	Snapshot Ensembles	1.0	1.1	2.3	11.3	17.2
	Cyclic SGLD	1.0	1.0	1.8	4.8	6.7
	SWA-Gaussian	1.0	1.8	1.9	2.5	2.7
	Fast Geometric Ens.	1.0	1.6	1.8	2.5	2.9
	Variational Inf. (FFG)	1.0	1.0	1.0	1.3	1.3
	KFAC-Laplace	1.0	1.0	1.0	1.2	1.2
WideResNet CIFAR-10	Single model	1.1	1.1	1.1	1.1	1.1
	Deep Ensembles	1.0	5.0	10.0	48.0	86.0
	Snapshot Ensembles	1.0	1.3	2.5	30.6	100.0
	Cyclic SGLD	1.0	1.0	1.6	4.0	4.1
	SWA-Gaussian	1.0	1.8	2.5	3.1	3.3
	Fast Geometric Ens.	1.0	1.2	1.7	2.7	3.2
	Dropout	1.0	1.2	1.3	1.4	1.4
	Variational Inf. (FFG)	1.0	1.0	1.0	1.3	1.4
VGG16 CIFAR-100	KFAC-Laplace	1.0	1.0	1.0	1.2	1.2
	Single model	1.0	1.0	1.0	1.0	1.0
	Deep Ensembles	1.0	5.0	10.0	50.0	100.0
	Snapshot Ensembles	1.0	2.9	5.4	18.8	26.8
	Cyclic SGLD	1.0	1.8	3.3	15.9	28.2
	SWA-Gaussian	1.0	2.3	3.4	5.2	5.9
	Fast Geometric Ens.	1.2	1.9	2.3	4.5	5.8
	Dropout	1.0	1.0	1.0	1.0	1.1
ResNet110 CIFAR-100	Variational Inf. (FFG)	1.0	1.6	1.8	2.1	2.2
	KFAC-Laplace	1.0	1.0	1.2	1.3	1.3
	Single model	1.0	1.0	1.0	1.0	1.0
	Deep Ensembles	1.0	5.0	10.0	50.0	99.0
	Snapshot Ensembles	1.0	1.3	1.9	4.7	6.2
	Cyclic SGLD	1.0	1.3	2.0	7.0	12.2
	SWA-Gaussian	1.0	1.7	2.2	3.3	3.6
	Fast Geometric Ens.	1.5	2.6	3.0	4.2	4.7
ResNet164 CIFAR-100	Variational Inf. (FFG)	1.0	1.2	1.3	1.5	1.5
	KFAC-Laplace	1.0	1.0	1.1	1.3	1.4
	Single model	1.0	1.0	1.0	1.0	1.0
	Deep Ensembles	1.0	5.0	10.0	50.0	100.0
	Snapshot Ensembles	1.0	1.3	1.9	4.8	6.3
	Cyclic SGLD	1.0	1.3	2.0	6.8	13.3
	SWA-Gaussian	1.0	2.1	2.6	3.6	3.8
	Fast Geometric Ens.	1.6	2.6	3.0	4.2	4.8
WideResNet CIFAR-100	Variational Inf. (FFG)	1.0	1.2	1.4	1.6	1.6
	KFAC-Laplace	1.0	1.1	1.3	1.5	1.5
	Single model	1.0	1.0	1.0	1.0	1.0
	Deep Ensembles	1.0	5.0	10.0	49.0	95.9
	Snapshot Ensembles	1.0	2.5	4.2	16.2	20.0
	Cyclic SGLD	1.0	1.9	3.2	13.6	18.8
	SWA-Gaussian	1.0	3.2	4.2	5.7	6.0
	Fast Geometric Ens.	1.3	3.2	4.4	7.5	9.5
CIFAR-100	Dropout	1.0	1.3	1.4	1.5	1.5
	Variational Inf. (FFG)	1.0	1.3	1.7	2.2	2.4
	KFAC-Laplace	1.0	1.0	1.0	1.1	1.1
	Single model	1.0	1.0	1.0	1.0	1.0

Table 5: Deep ensemble equivalent score for CIFAR-10/100.

Model	Method	Error (%)				Negative calibrated log-likelihood			
		1	5	10	50	1	5	10	50
ResNet50	Fast Geometric Ens.	23.71 \pm 0.00	23.61 \pm 0.00	23.56 \pm 0.00	23.28 \pm 0.00	0.929 \pm 0.000	0.921 \pm 0.000	0.916 \pm 0.000	0.904 \pm 0.000
	Deep Ensembles	23.79 \pm 0.14	21.19 \pm 0.14	20.90 \pm 0.08	20.63 \pm NA	0.935 \pm 0.007	0.823 \pm 0.002	0.805 \pm 0.000	0.788 \pm NA
	Single model	23.86 \pm 0.20	23.86 \pm 0.20	23.86 \pm 0.20	23.86 \pm 0.20	0.938 \pm 0.006	0.938 \pm 0.006	0.938 \pm 0.006	0.938 \pm 0.006
	Variational Inf. (FFG)	24.50 \pm 0.06	23.82 \pm 0.03	23.77 \pm 0.04	23.67 \pm 0.00	0.957 \pm 0.001	0.927 \pm 0.000	0.923 \pm 0.001	0.920 \pm 0.000
	KFAC-Laplace	25.01 \pm 0.49	24.19 \pm 0.29	23.93 \pm 0.20	23.86 \pm 0.16	0.988 \pm 0.022	0.948 \pm 0.013	0.939 \pm 0.011	0.934 \pm 0.008
	Snapshot Ensembles	24.92 \pm NA	22.21 \pm NA	21.75 \pm NA	21.48 \pm NA	0.983 \pm NA	0.865 \pm NA	0.843 \pm NA	0.830 \pm NA

Table 6: Classification error and negative calibrated log-likelihood for different numbers of samples on ImageNet.

Model	Method	Error (%)			Negative calibrated log-likelihood		
		5	10	50	5	10	50
ResNet50	Fast Geometric Ens.	23.61 vs 22.21 \downarrow	23.56 vs 21.37 \downarrow	23.28 vs 20.67 \downarrow	0.921 vs 0.894 \downarrow	0.916 vs 0.842 \downarrow	0.904 vs 0.793 \downarrow
	Deep Ensembles	21.19 vs 21.20 \downarrow	20.90 vs 20.16 \downarrow	20.63 vs 19.39 \downarrow	0.823 vs 0.855 \uparrow	0.805 vs 0.793 \downarrow	0.788 vs 0.739 \downarrow
	Single model	23.86 vs 22.39 \downarrow	23.86 vs 21.60 \downarrow	23.86 vs 21.06 \downarrow	0.938 vs 0.900 \downarrow	0.938 vs 0.851 \downarrow	0.938 vs 0.805 \downarrow
	Variational Inf. (FFG)	23.82 vs 22.58 \downarrow	23.77 vs 21.74 \downarrow	23.67 vs 21.10 \downarrow	0.927 vs 0.905 \downarrow	0.923 vs 0.851 \downarrow	0.920 vs 0.805 \downarrow
	KFAC-Laplace	24.19 vs 22.50 \downarrow	23.93 vs 21.67 \downarrow	23.86 vs 21.04 \downarrow	0.948 vs 0.906 \downarrow	0.939 vs 0.855 \downarrow	0.934 vs 0.809 \downarrow
	Snapshot Ensembles	22.21 vs 21.99 \downarrow	21.75 vs 20.81 \downarrow	21.48 vs 19.86 \downarrow	0.865 vs 0.879 \uparrow	0.843 vs 0.815 \downarrow	0.830 vs 0.763 \downarrow

Table 7: Classification error and negative calibrated log-likelihood before vs. after test-time augmentation on ImageNet.

Model	Method	Deep ensemble equivalent score				
		1 sample	5 samples	10 samples	25 samples	50 samples
ResNet50	Deep Ensembles	1.0	5.0	10.0	25.0	50.0
	Snapshot Ensembles	1.0	2.2	3.2	4.0	4.2
	Fast Geometric Ens.	1.1	1.2	1.3	1.4	1.5
	Variational Inf. (FFG)	1.0	1.1	1.2	1.2	1.2
	KFAC-Laplace	1.0	1.0	1.0	1.0	1.0
	Single model	1.0	1.0	1.0	1.0	1.0

Table 8: Deep ensemble equivalent score for ImageNet.

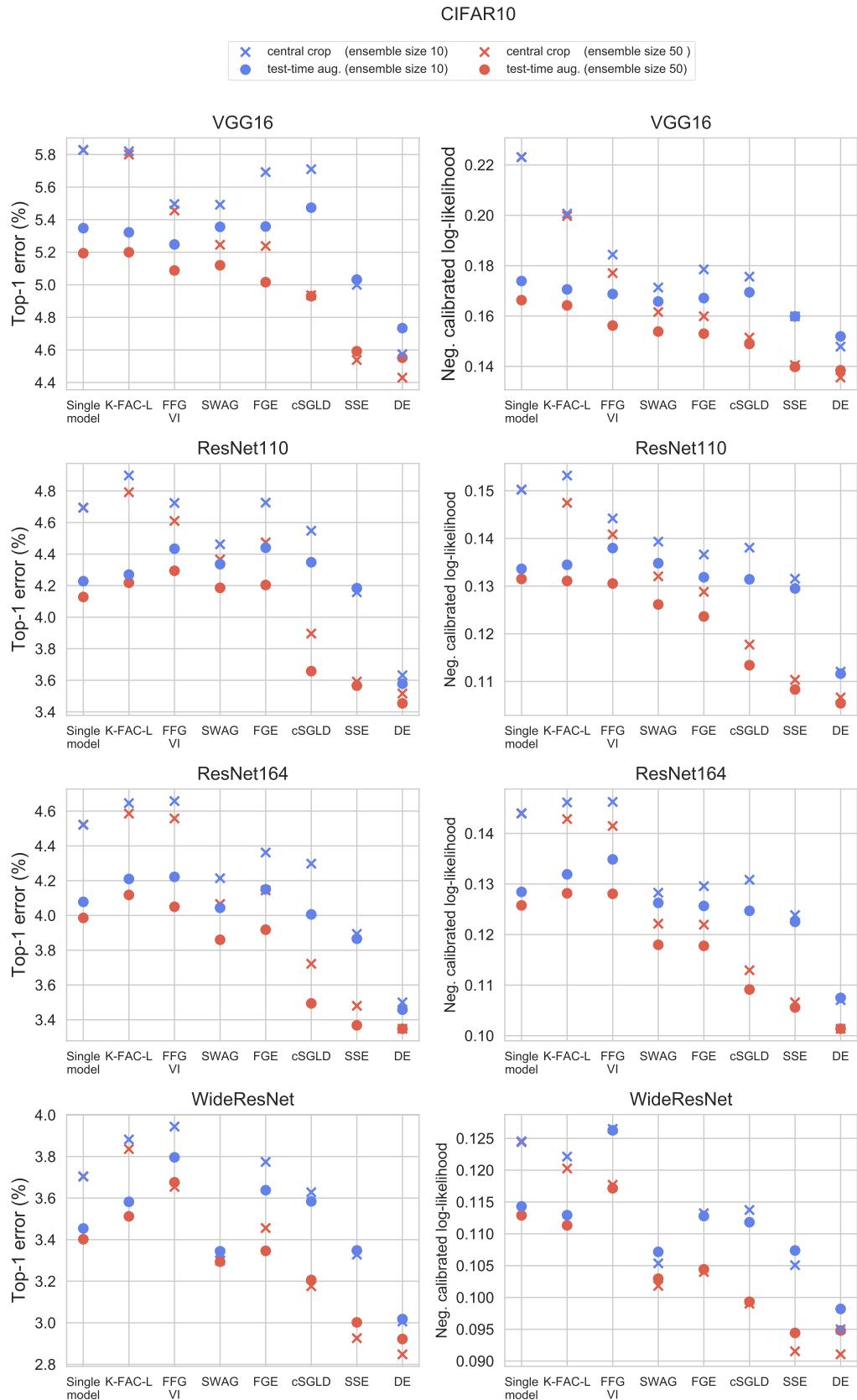


Figure 14: Classification error and negative calibrated log-likelihood before vs. after test-time augmentation on CIFAR-10

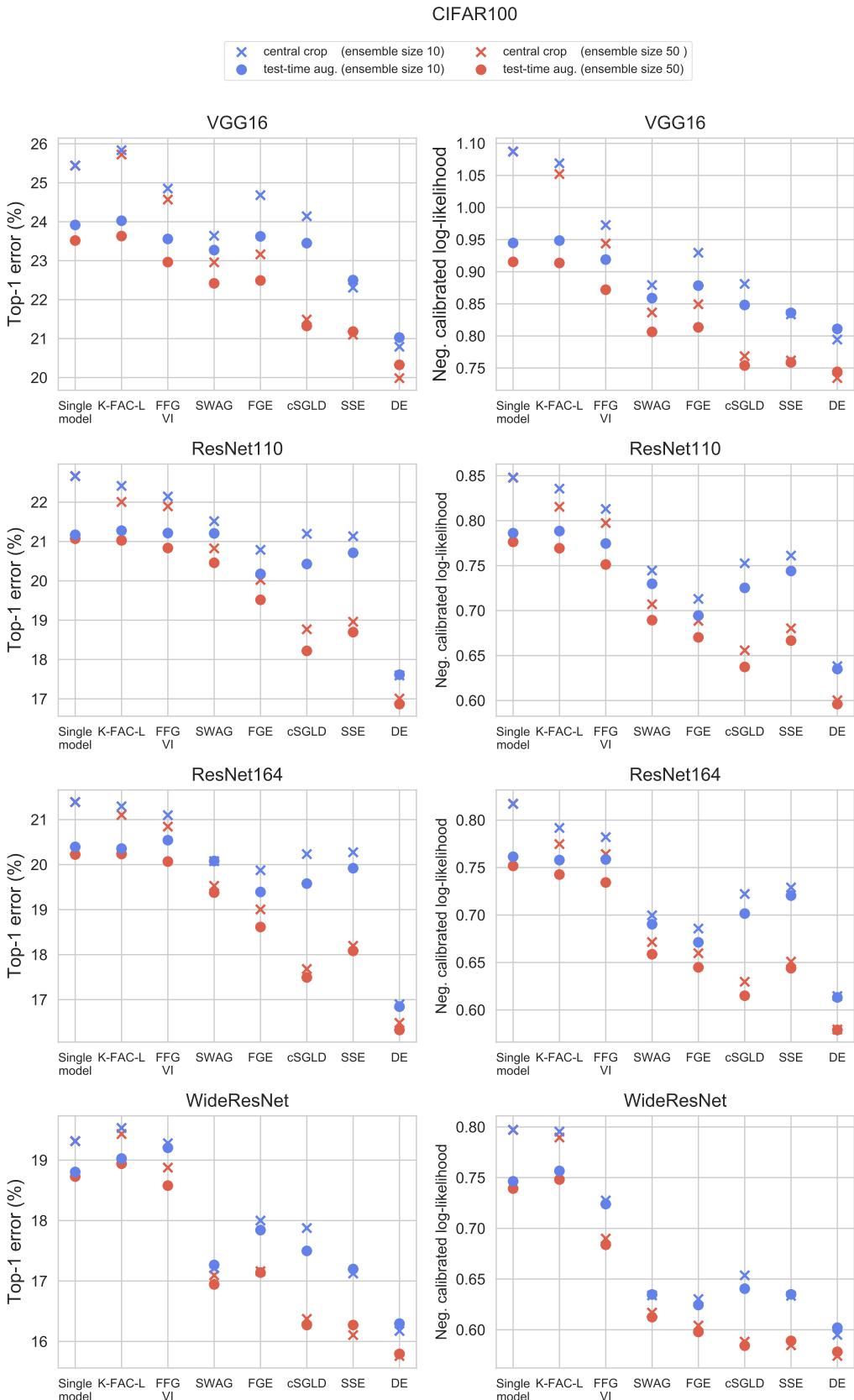


Figure 15: Classification error and negative calibrated log-likelihood before vs. after test-time augmentation on CIFAR-100