# ON UNIFORM, BAYESIAN, AND PAC-BAYESIAN DEEP ENSEMBLES

Anonymous authors

Paper under double-blind review

### Abstract

It is common practice to combine deep neural networks to ensembles. These deep ensembles can profit from the cancellation of errors effect: Errors by ensemble members may average out and the ensemble achieves better generalization performance than each individual network. Bayesian neural networks learn a posterior distribution over model parameters, and sampling and weighting networks according to this posterior yields an ensemble model referred to as Bayes ensemble. In this study, we stress that neither the sampling nor the weighting in a Bayes ensemble are particularly well-suited for increasing generalization performance, as they do not support the cancellation of errors effect, which is evident in the limit from the Bernstein-von Mises theorem for misspecified models. In contrast, a weighted average of models, where the weights are optimized by minimizing a PAC-Bayesian generalization bound, can improve generalization performance. This requires that the optimization takes correlations between models into account, which can be achieved by minimizing the tandem loss at the cost that hold-out data for estimating error correlations need to be available. The PAC-Bayesian weighting increases the robustness against correlated models and models with lower performance in an ensemble. This allows us to safely add several models from the same learning process to an ensemble, instead of using early-stopping for selecting a single weight configuration. Our study presents empirical results supporting these conceptual considerations on four different classification datasets. We show that state-of-the-art Bayes ensembles from the literature, despite being computationally demanding, do not improve over simple uniformly weighted deep ensembles and cannot match the performance of deep ensembles weighted by optimizing the tandem loss, which additionally come with non-vacuous generalization guarantees.

004

010 011

012

013

014

015

016

017

018

019

021

023

025

026

028

029

031

### 1 INTRODUCTION

Combining different deep neural networks to an ensemble model is a common way to increase 037 generalization performance (Bishop and Bishop, 2023; Goodfellow et al., 2016). Such deep ensembles can profit from the cancellation of errors effect (Eckhardt and Lee, 1985): When the individual networks perform better than random guessing and make independent errors, the errors in an additive 040 ensemble (or committee) average out and the ensemble tends to outperform the individual models 041 (Hansen and Salamon, 1990; Perrone and Cooper, 1993; Dietterich, 2000). Many strategies to create 042 the networks for an ensemble have been explored (e.g. Buschjäger et al., 2020; D'Angelo and Fortuin, 043 2021; Dziugaite and Roy, 2017; Garipov et al., 2018; Huang et al., 2017; Jiang et al., 2017; Lee 044 et al., 2015; Liu and Yao, 1999; Masegosa et al., 2020; Monteith et al., 2011; Ortega et al., 2022; Pérez-Ortiz et al., 2021), however, simple deep ensembles, which solely rely on repeating the neural network training with random weight initializations and stochastic training to generate ensemble 046 members, often already improve generalization (Lakshminarayanan et al., 2017). 047

One way to build a neural network ensemble is based on the predictive posterior of Bayesian neural networks (Adlam et al., 2020; Aitchison, 2021; Chen et al., 2014; Farquhar et al., 2020; Gal and Ghahramani, 2016; Grünwald, 2012; Izmailov et al., 2021; Kapoor et al., 2022; Kendall and Gal, 2017;
MacKay, 1992b; Maddox et al., 2019; Monteith et al., 2011; Nabarro et al., 2022; Neal, 1992; 1996; 2011; Pearce et al., 2020; Ritter et al., 2018; Welling and Teh, 2011; Wenzel et al., 2020a;b; Wilson and Izmailov, 2020; Zhang et al., 2020), and this approach has gained attention recently (Bachmann et al., 2022; Fortuin et al., 2022; Rudner et al., 2022; Wiese et al., 2023). A Bayesian neural network

054 learns a posterior distribution over the network weights (MacKay, 1992a; Neal, 1996). Averaging 055 networks sampled from this posterior yields an approximation of the Bayesian model average (BMA) 056 and is referred to as Bayes ensemble. This approach comes with great expectations, for example, Domingos (2000) stated 'Given the "correct" model space and prior distribution, Bayesian model averaging is the optimal method for making predictions; in other words, no other approach can 058 consistently achieve lower error rates than it does.' We argue that neither the sampling nor the weighting in a Bayes ensemble are particularly well-suited for increasing generalization performance 060 in theory and practice (see also Ortega et al., 2022). The Bernstein-von Mises theorem (assuming 061 identifiability of the model) shows that the BMA converges towards the maximum likelihood point 062 estimate with growing dataset size (Kleijn and Van der Vaart, 2012). In the limit, it will eventually 063 concentrate on a single model without exploiting ensemble diversity. Thus, the Bayes ensemble 064 typically does not leverage the cancellation of errors effect. This is particularly problematic because 065 the considered model space in the BMA setting is typically not 'correct': The Bayes ensemble itself 066 is in general outside the model space from which its members are sampled - and in the space of 067 ensemble models the Bayes ensemble cannot be expected to be optimal. 068

PAC-Bayesian methods (McAllester, 1998) provide an alternative approach to incorporate prior 069 information about models and to improve the weighting of ensemble members. We bring forward optimizing the weighting of deep ensemble members using a PAC-generalization bound based on 071 the tandem loss, which accounts for pairwise correlations between networks (Masegosa et al., 2020). This allows us to maintain diversity and increase generalization performance of the deep ensemble. 073 The PAC-Bayesian weight optimization is especially useful when intermediate weight configurations 074 (snapshots) from model training are considered, in which case the optimization performs model 075 selection taking both individual network performance and ensemble diversity into account. The price to pay for the improved weighting is additional hold-out data for the optimization of the weighs, 076 however, then the resulting PAC-Bayesian bound provides a rigorous performance guarantee. 077

Our study transfers existing, but scattered and often partly neglected knowledge on the (1) behavior
of Bayesian model averaging (BMA), in particular its suboptimality in terms of generalization
performance in the (typical) misspecified setting, (2) cancellation of errors, (3) and second order
PAC-Bayesian bounds to ensembling of deep neural networks (DNNs). The main contributions can
be summarized as:

- We discuss the conceptual differences between simple, Bayes and PAC-Bayesian deep ensembles, stressing that Bayes ensembles do not consider the cancellation of errors effect, which explain some observations in the current literature (e.g., on "cold posteriors").
- We support our conceptual arguments through unbiased experiments on DNN benchmark problems that were used to promote BMA. Our results on four datasets show that complex Bayesian approximate inference methods can often be surpassed by more efficient simple deep ensembles. Bagging decreased predictive performance, showing the trade-off between the degree of randomization and single-network performance given limited training data.
  - We show that the optimization of a PAC-Bayesian generalization bound using the tandem loss can improve the predictive performance of deep ensembles and provides non-vacuous generalization guarantees. This is in contrast to what could have been expected from the results by Ortega et al. (2022). Our results stress the importance of *second order* PAC-Bayesian analysis for ensembles.
  - We demonstrate that the inclusion of intermediate checkpoints from the same training run in a neural network ensemble can increase its predictive performance, especially when the PAC-Bayesian bound optimization weights their contributions.

In the next section, we will briefly summarize the background on Bayesian model averaging, the cancellation of errors effect, deep ensembles, and PAC-Bayesian majority voting. Section 3 will bring these topics together: It contrasts uniform, Bayes, and PAC-Bayesian deep ensembles and formulates our main hypotheses. These hypotheses will be empirically studied in Section 4 before we conclude.

104 105

084

085

090

092

093

094

095

096

098

### 2 BACKGROUND

We consider data  $\mathcal{D} = \{(x_i, y_i)\}_{i=1}^n$  drawn i.i.d. from a distribution p. The loss of a hypothesis/model h on (x, y) is denoted by  $\ell(h(x), y)$ . The generalization error (risk) of h is given by L(h) = 1

 $\mathbb{E}_{(x,y)\sim p}[\ell(h(x), y)]$  and the empirical risk by  $\hat{L}_{\mathcal{D}}(h) = \frac{1}{n} \sum_{i=1}^{n} \ell(h(x_i), y_i)$ . We write  $h_w$  to stress that a model is parameterized by  $w \in \mathbb{R}^d$ . We consider ensembles of M models  $h_i = h_{w_i} \in \mathcal{H}$ ,  $i = 1, \ldots, M$ , weighted by  $\rho \in \mathbb{R}^M$ . For regression, the ensemble regressor is  $h_{\rho}(x|h_1, \ldots, h_M) = M$ . 108 109 110 111  $\sum_{i=1}^{M} \rho_i h_i(x)$ . For classification, we distinguish two ways of combining the predictions of the 112 ensemble members. We assume a predictive distribution  $p_i(y|x) = p(y|x, w_i)$  associated with 113 each hypothesis  $h_{w_i}(x)$ . The ensemble prediction  $h_{\rho}(x|h_1,\ldots,h_M)$  can be either defined by the  $\rho$ -weighted average of these distributions  $\operatorname{AVG}_{\rho}(x) = \operatorname{arg} \max_{y} \sum_{i=1}^{M} \rho_{i} p_{i}(y|x)$  or by the  $\rho$ -weighted majority vote  $\operatorname{MV}_{\rho}(x) = \operatorname{arg} \max_{y} \sum_{i=1}^{M} \rho_{i} \mathbb{1}[y = \operatorname{arg} \max_{y'} p_{i}(y'|x)]$ , where  $\mathbb{1}[\cdot]$  is 1 if its argument is true and zero otherwise. Uniform weighted aggregations are denoted by  $\operatorname{AVG}_{u}$  and 114 115 116 117  $MV_u$ . 118

119 120 120 120 121 122 Bayesian model average. In Bayesian inference, we update our prior belief about the parameters w of a model predicting y given x according to p(y|x, w) to a posterior distribution p(w|D) based on training data D. Marginalizing p(y|x, w) over this posterior gives the posterior predictive distribution

$$p(y|x, \mathcal{D}) = \int p(y|x, w) p(w|\mathcal{D}) \mathrm{d}w, \qquad (1)$$

also known as Bayesian Model Average (BMA), Bayes ensemble or posterior predictive. It predicts by averaging over all possible models weighted by their posterior probability. The BMA explicitly models aleatoric uncertainty (given by p(y|x, w)) and epistemic uncertainty in the form of p(w|D)(see Caprio et al., 2024 for a recent discussion). As the integral in (1) is in general intractable, it is approximated in practice by an average of M models sampled from the posterior:

123 124

132

137

148

 $p(y|x, \mathcal{D}) \approx \frac{1}{M} \sum_{m=1}^{M} p(y|x, w_m), \quad w_m \sim p(w|\mathcal{D})$ (2)

133 It is in general difficult to sample from p(w|D), which can, for example, be addressed by Markov 134 chain Monte Carlo (MCMC) sampling (Chen et al., 2014; Neal, 1992; 2011; Welling and Teh, 135 2011; Wenzel et al., 2020a; Zhang et al., 2020), Monte Carlo Dropout (Folgoc et al., 2021; Gal and 136 Ghahramani, 2016), and Laplace approximation (MacKay, 1992b; Ritter et al., 2018).

**Cancellation of errors.** The cancellation of errors effect refers to the fact that, when individual 138 ensemble members perform better than random guessing and make independent errors, their errors 139 average out and the combined model outperforms the individual predictors (Eckhardt and Lee, 1985; 140 Hansen and Salamon, 1990). Let us assume binary classification and that each ensemble member 141 has an error probability lower bounded by  $p_{\text{max}} < \frac{1}{2}$ . Given M models, the generalization error of the ensemble is upper bounded by  $\exp\left(-2\left(\frac{M+1}{2} - Mp_{\text{max}}\right)^2 M^{-1}\right)$ , see Appendix A.3. That is, 142 143 144 in the idealized (and not realistic) setting of independent errors, the generalization error vanishes with increasing M if the classifiers are better than random guessing. The corresponding regression 145 case relating expected mean-squared error and correlation of the model outputs is discussed in deep 146 learning textbooks by (Goodfellow et al., 2016, sec. 7.11) and (Bishop and Bishop, 2023, sec. 9.6). 147

**Deep ensembles.** Deep ensembles (Lakshminarayanan et al., 2017) are ensembles with deep neural 149 networks as members. More specifically, we refer to a deep ensemble if the networks all share 150 the same structure. These networks are typically the result of independent training processes. The 151 diversity is a result of random initialization and the stochastic optimization. This randomness alone 152 can yield well-performing ensembles, while additional randomization of the training data by bootstrap 153 aggregation (bagging) bears the risk of the ensemble being worse than a single network trained on 154 all data (Lakshminarayanan et al., 2017; Lee et al., 2015). Other ways to create deep ensemble 155 members include the use of a cyclical learning rate schedule and taking the intermediate checkpoints 156 as snapshot ensembles (SSEs) (Huang et al., 2017) as well as searching for other ensemble members 157 in neighborhood of a single pre-trained network (fast geometric ensembling (Garipov et al., 2018)). While being introduced as an alternative to Bayesian approaches, Bayesian interpretations of deep 158 ensembles can be found in recent work (D'Angelo and Fortuin, 2021; Wilson and Izmailov, 2020). 159

- 160
- **PAC-Bayesian majority voting.** PAC-Bayesian analysis (McAllester, 1998; Valiant, 1984) provides bounds on the generalization error of Gibbs classifiers defined by a distribution  $\rho$  over a (subset of a)

162 hypothesis space  $\mathcal{H}$  given empirical risks for the hypotheses. A Gibbs classifier draws a hypothesis h 163 according to  $\rho$  at random for each input x and returns the prediction h(x). PAC-Bayesian bounds hold 164 for all distributions  $\rho$  over  $\mathcal{H}$  simultaneously. This allows us to directly optimize a PAC-Bayesian 165 bound in terms of  $\rho$ , e.g., by gradient-based methods (Dziugaite and Roy, 2017; Masegosa et al., 2020; 166 Masegosa, 2020; Ortega et al., 2022; Pérez-Ortiz et al., 2021; Thiemann et al., 2017). In a weighted majority voting classifier  $MV_{\rho}$  (e.g., a deep ensemble), the randomized prediction is replaced by a 167  $\rho$ -weighted vote by all hypotheses (Germain et al., 2009; 2015). PAC-Bayesian bounds can be applied 168 to  $MV_{\rho}$  by bounding  $L(MV_{\rho})$  by twice the risk of the corresponding Gibbs classifier (Germain et al., 169 2015; Langford and Shawe-Taylor, 2002; Masegosa et al., 2020). Gibbs classifiers ignore interactions 170 between the models (such as cancellation of errors), and optimization of the weighting using bounds 171 on the Gibbs classifier directly does typically not increase generalization performance of the ensemble 172 (e.g., Lorenzen et al., 2019). This can be addressed by second order PAC-Bayesian bounds (Germain 173 et al., 2015; Lacasse et al., 2006; Masegosa et al., 2020). In particular, Masegosa et al. derive a bound 174 in terms of the tandem loss  $L(h, h') = \mathbb{E}_{p(x,y)}[\mathbb{1}(h(x) \neq y \land h'(x) \neq y)]$ , which takes pairwise correlations into account: For any probability distribution  $\pi$  on  $\mathcal{H}$  that is independent of  $\mathcal{D}$  and 175 176 any  $\delta \in ]0,1[$ , with probability at least  $1-\delta$  over a random draw of  $\mathcal{D}$  with n elements, for all 177 distributions  $\rho$  on  $\mathcal{H}$  and all  $\lambda \in ]0, 2[$  simultaneously:

178 179

181

182

183

 $L(\mathrm{MV}_{\rho}) \le 4\left(\frac{\mathbb{E}_{\rho^2}[\hat{L}_{\mathcal{D}}(h, h'])}{1 - \lambda/2} + \frac{2\mathrm{KL}(\rho \| \pi) + \ln(2\sqrt{n}/\delta)}{\lambda(1 - \lambda/2)n}\right)$ (3)

Here  $\text{KL}(\rho \| \pi)$  denotes the Kullback–Leibler divergence between prior  $\pi$  and posterior  $\rho$ . The bound can be efficiently optimized w.r.t. to  $\rho$  and  $\lambda$  (Thiemann et al., 2017).

We consider the optimization of  $\rho$  (i.e., the weights in a weighted ensemble) given M models as proposed by Masegosa et al. (2020). One could instead also turn a second order PAC-Bayesian bound into an objective function for optimizing  $\rho$  and the parameters of the M models simultaneously. This has been done by Ortega et al. (2022) to optimize deep ensembles. However, because of the simultaneous optimization of the neural networks' parameters and their weighting, the formal guarantees by the PAC-Bayesian bounds get lost. Furthermore, Ortega et al. did not find improvements compared to uniformly weighted deep ensembles for larger networks (e.g., ResNet20, He et al., 2016) in their empirical evaluation.

192 193

194

199 200

### 3 UNIFORM, BAYESIAN, AND PAC-BAYESIAN DEEP ENSEMBLES

**Bayes ensembles and cancellation of errors.** As outlined above, general ensemble methods and Bayesian neural networks have different motivations. In a weighted additive ensemble, a finite set of hypotheses from  $\mathcal{H}$  are combined to a new model  $h_{\rho} \in \mathcal{H}_{\rm E}$  from the hypothesis class

$$\mathcal{H}_{\mathrm{E}} = \left\{ h_{\rho}(x|h_{1}, \dots, h_{M}) \, \middle| \, M \in \mathbb{N}, \rho \in \mathbb{R}^{M}, h_{i} \in \mathcal{H} \right\}.$$
(4)

As discussed in Section 2,  $h_{\rho}$  can generalize better than the individual  $h_1, \ldots, h_M$  if the ensemble members are diverse in the sense that their errors are not strongly correlated. Typically,  $\mathcal{H} \subset \mathcal{H}_E$ , and for most ensemble methods this is key. The ensemble members selected from  $\mathcal{H}$  are often referred to as weak learners and are not expected to have a very low risk (generalization performance) as long as they are better than random guessing: The desired model is expected to be in  $\mathcal{H}_E \setminus \mathcal{H}$ . Boosting with decision stumps is a prime example (Freund et al., 1996).

In Bayesian machine learning, the posterior distribution  $p(h|\mathcal{D})$  over  $\mathcal{H}$  models the uncertainty of 207 identifying a desired model  $h^* \in \mathcal{H}$  given finite training data  $\mathcal{D}$  and the prior over  $\mathcal{H}$ .<sup>1</sup> In Bayesian 208 deep learning with a fixed neural network architecture  $h_w$  parameterized by weights  $w \in \mathbb{R}^d$ , we have 209  $\mathcal{H} = \{h_w | w \in \mathbb{R}^d\}$  and the posterior over the hypotheses is modelled by a distribution  $p(w | \mathcal{D})$  over 210 the weights w. To simplify the discussion, let us assume that the models are parameterized in a way 211 that ensures identifiability, that is, different parameters generate different hypotheses. With increasing 212 training data set size  $n = |\mathcal{D}|$  the uncertainty about the desired model decreases, and  $p(w | \mathcal{D})$  should 213 concentrate on  $w^*$  (i.e., the weights of the desired model  $h^* = h_{w^*}$ ). This is formally shown by the 214

<sup>214</sup> 

<sup>&</sup>lt;sup>1</sup>The desired model in  $\mathcal{H}$  is the model in  $\mathcal{H}$  closest to the true model  $h_p$  underlying the data generating distribution p. We do neither assume  $h_p \in \mathcal{H}$  nor  $h_p \in \mathcal{H}_{E}$ .

216 Bernstein-von Mises theorem (Kleijn and Van der Vaart, 2012; Van der Vaart, 2000), which was first 217 discovered by Laplace (1820), taken up by Bernstein (1917) and von Mises (1931) independently, and 218 proven in a more general form by Le Cam (1953), who coined the name. Let the parametric model 219  $h_w$  be twice differentiable with non-singular Fisher information matrix  $\mathcal{I}(w^*)$  at  $w^*$ , identifiable and 220 continuous on a compact parameter space. Further, let the prior measure be absolutely continuous in a neighbourhood of  $w^*$  with a continuous density non-zero at  $w^*$ . Let us consider the sequence 221 of posterior distributions  $p_n(w \mid \mathcal{D}_n)$  and maximum likelihood estimates  $\hat{w}_n$  given  $\mathcal{D}_n$  for data sets 222  $\mathcal{D}_n \subset \mathcal{D}_{n+1}$ . The Bernstein-von Mises theorem implies 223

224 225

$$\|p_n(w|\mathcal{D}_n) - \mathcal{N}(\hat{w}_n, n^{-1}\mathcal{I}^{-1}(w^*))\|_{\mathrm{TV}} \stackrel{p_{w^*}}{\to} 0.$$
(5)

That is, for increasing data set size the posterior converges in probability to a normal distribution  $\mathcal{N}(\hat{w}_n, n^{-1}\mathcal{I}^{-1}(w^*))$ , where  $\|\cdot\|_{\text{TV}}$  denotes the total variation norm (see Van der Vaart, 2000). The normal distribution is centered on the maximum likelihood estimate with covariance  $n^{-1}\mathcal{I}^{-1}(w^*)$ , that is, the distribution gets more and more concentrated with increasing data set size. Under the same assumptions, this implies that the expectation of the Bayes ensemble (2) converges to the maximum likelihood estimate  $h_{w^*} \in \mathcal{H}$  for  $n \to \infty$  (while the Bayes ensemble is not necessarily in  $\mathcal{H}$  for finite n).

Thus, as already argued by Masegosa (2020, sec. 4), from the perspective of general ensemble 233 learning and the goal of finding a model with minimum risk, a sample from the posterior  $p(w \mid \mathcal{D})$ 234 does not appear to be a particularly good choice for ensemble members. As the sample will be 235 distributed around the maximum likelihood estimate  $h_{w^*} \in \mathcal{H}$ , the ensemble members cannot be 236 expected to be very diverse, and the lack of diversity reduces the cancellation of errors effect. The 237 Bayes ensemble will tend to the model with the highest likelihood in  $\mathcal{H}$  – and not to the model with 238 the highest likelihood in  $\mathcal{H}_{\rm E}$ . One may well argue that the presented asymptotic argument is of little 239 practical relevance if  $\mathcal{H}$  already has a high capacity (e.g., consists of overparameterized deep neural 240 networks), because then  $h_{w^*} \in \mathcal{H}$  may be close to the optimal solution. However, the important point is that even for finite n, the M ensemble members can be expected to be quite similar in their 241 predictive behavior. Thus, we formulate the following hypothesis: 242

Hypothesis 1 The Bayes ensemble is not a particularly good way of selecting and weighting networks in a deep ensemble.

While this hypothesis appears to be evident given the theoretical considerations above, many studies evaluate the performance of the Bayes ensemble by measuring generalization performance on a test set and provide favorable comparisons to baselines suggesting the posterior predictive as a way to maximize generalization performance. This study presents unbiased comparisons of Bayes ensembles with simple deep ensemble approaches to test Hypothesis 1.

A Bayesian way to finding the best generalizing ensemble is to lift the Bayesian reasoning from  $\mathcal{H}$ to  $\mathcal{H}_{\rm E}$ . This has been demonstrated by Monteith et al. (2011), who apply BMA not to  $h_w \in \mathcal{H}$ , but to  $h_{\rho} \in \mathcal{H}_{\rm E}$ , where  $\mathcal{H}_{\rm E}$  are all ensembles of a constant size M. In such an approach, the posterior distribution contracts around the maximum likelihood estimate of the best weighting of ensemble members. Thus, we go from a distribution over the weights of a single neural network to a distribution over all parameters of the M networks and the ensemble weights. This growth in the number of dimensions renders the approach computationally infeasible in larger settings, however.

258

**PAC-Bayesian deep ensembles.** If Hypothesis 1 holds true, this does neither mean that we cannot 259 do better than uniform averaging given M models from  $\mathcal{H}$  (as already pointed out in early works 260 on ensembles, e.g., Krogh and Sollich, 1997) nor that we cannot make use of priors and Bayesian 261 reasoning to improve ensembles. For example, there is no need to put weight on an ensemble member 262 that is worse than random guessing or always errs when another ensemble member errs (i.e., it can 263 never correct an error). Lorenzen et al. (2019); Masegosa et al. (2020); Wu et al. (2021) Given the 264 results by Masegosa et al. (2020); Wu et al. (2021), it appears to be promising to weight ensemble 265 members based on PAC-Bayesian generalization bounds optimized using the tandem loss. It is 266 crucial to consider a second order bound that takes interactions between ensemble members – error 267 cancellation – into account, otherwise the approach suffers from the same problem as BMA and concentrates all weight on a single ensemble member in the limit. However, the bounds require for 268 each ensemble member  $h_i$  that there are data points  $\mathcal{D}_i$  not used for training  $h_i$  and that for each 269 pair of models  $h_i$  and  $h_j$  the overlap  $\mathcal{D}_i \cap \mathcal{D}_j$  is large enough for a good enough estimate of the

tandem loss: In (3), we compute  $\hat{L}_{\mathcal{D}_i \cap \mathcal{D}_j}(h_i, h_j)$  for each hypotheses pair  $h_i$  and  $h_j$  and n becomes  $n = \min_i |\mathcal{D}_i|$ . For random forests as considered by Lorenzen et al. (2019) and Masegosa et al. (2020), data sets  $\mathcal{D}_i \subset \mathcal{D}$  naturally arise from the bagging procedure and the proposed bounds can be computed and optimized for free. When applied to deep ensembles, we have to pay by leaving data out when training the ensemble members, which can be expected to reduce the performance of the individual networks. Taking this into account, we hypothesize:

Hypothesis 2 PAC-Bayesian weighting optimized using the tandem loss can improve the generalization performance of a deep ensemble, in particular if extra data for computing the bound are available, and provide non-vacuous performance guarantees.

280 In most deep ensemble approaches, each ensemble member is the result of an independent training 281 process, where the weight configuration is chosen by early-stopping (requiring some hold-out data) or 282 simply after a predefined number of learning iterations (which is then also a crucial hyperparameter). 283 In contrast, Wenzel et al. (2020a) sample from the posterior using a single stochastic gradient Markov 284 chain Monte Carlo (SG-MCMC) sequence and SSEs combine networks from a single learning process (Huang et al., 2017). That is, one process creates all ensemble members. If the latter part of 285 Hypothesis 2 is true, one need not worry about adding models to the ensemble that do not perform 286 particularly well and/or are correlated. The PAC-Bayesian weighting performs a model selection 287 taking both individual performance and diversity into account, which should provide protection 288 from networks not contributing positively to a ensemble. This makes it safe to add several weight 289 configurations (checkpoints) from a single learning process to an ensemble. It also renders using 290 a hold-out validation dataset for early-stopping the neural network training unnecessary: It should 291 neither be harmful to add underfitted nor overfitted models. The early-stopping data can instead be 292 used to optimize the PAC-Bayesian weighting: 293

**Hypothesis 3** Optimizing the weighting using the tandem loss allows inclusion of several models from a training run in a way that efficiently improves performance and makes early-stopping unnecessary.

295 296 297

298

294

### 4 EXPERIMENTS AND RESULTS

# 299 4.1 EXPERIMENTAL SETUP300

Four different datasets and four neural network architectures were considered.<sup>2</sup> We evaluated a CNN 301 LSTM (Yenter and Verma, 2017) on the IMDB binary classification benchmark (Maas et al., 2011) 302 and compared it to a state-of-the-art Bayesian approximate inference method from Wenzel et al. 303 (2020a) (referred to as cSGHMC-ap), who combine stochastic gradient Langevin dynamics (SGLD) 304 (Chen et al., 2014; Welling and Teh, 2011) with a cyclical learning rate schedule (Zhang et al., 2020), 305 and adaptive preconditioning (Li et al., 2016; Ma et al., 2015). Furthermore, we considered the 306 multi-class data sets CIFAR-10 and CIFAR-100 (Krizhevsky, 2009) with ResNet110 (He et al., 2016) 307 and WideResNet28-10 (WRN28-10) (Zagoruyko and Komodakis, 2016) architectures, for which 308 Ashukha et al. (2020) and Vadera et al. (2022) provide reference results for a variety of Bayesian 309 methods. Lastly, the ResNet50 was evaluated on the EyePACS dataset containing diabetic retinopathy 310 diagnoses distinguishing five degrees of severity. It was introduced by Band et al. (2021) as a 311 benchmark for Bayesian approximate inference methods. The authors provide results for a variety of Bayesian neural networks and ensembles of those, creating "deep ensembles of Bayes ensembles". 312 As we are concerned with a comparison to Bayes ensembles, we focus on the results of the individual 313 Bayesian methods as reference. The hyperparameters for each experiment are given in Appendix A.3. 314

To optimize the weighting while retaining the original training data set size, *test-time cross validation* as introduced by Ashukha et al. (2020) was employed. We used 50% of the hold-out data for bound optimization and the other half for testing. This was then repeated with the two subsets switching roles, and the results were averaged to decrease the variance for the unbiased performance estimate, see Appendix A.2 for details.

For each dataset and model architecture, a simple deep ensemble was constructed based on the reference paper's single-run hyperparameters. Intermediate checkpoints were stored and either included (referred to as *all* from now) or ignored in the final ensemble (*last*). Further, ensembles with

<sup>&</sup>lt;sup>2</sup>We consumed approximately 320 GPU days for all experiments on a local cluster.

sub-sampled training sets (bagging) as well as snapshot ensembles (SSE) using a cyclical learning
 rate were formed. Finally, the trade-off between ensemble size and number of training epochs was
 studied in a sequential setting in Appendix A.4.2.

We trained  $M_{\text{total}}$  networks for each setting, where  $M_{\text{total}}$  was 30, 50, and 10 for CIFAR, IMDB, and EyePACS, respectively. Ensembles of size M were created by sampling without replacement from these models. This was repeated five times for every setting, and we report mean and standard deviation  $\sigma$  over these five trials. Generalization bounds are reported for  $\delta = 0.05$ .



Figure 1: Mean test accuracy  $\hat{A} \pm \sigma$  vs. ensemble size M over five ensembles for uniformly and PAC-Bayesian weighted deep ensembles (AVG<sub>u</sub> and AVG<sub>p</sub>), using either only the last or all training checkpoints, and the best single model (SGD). References are Bayesian ensembles cSGHMC-ap (Wenzel et al., 2020a), cSGLD (Ashukha et al., 2020), *MC-Dr*opout and *MC-Dr*opout Deep Ensemble (i.e., ensemble of Bayesian ensembles), as well as simple Deep Ensembles for EyePACS from Band et al. (2021). Numbers in brackets indicate the ensemble sizes of the baselines. Additional results for other settings are shown in Figure A.4 in the appendix.

4.2 Results

364 365

366

**Bayesian vs. uniform deep ensembles (Hypothesis 1).** We compared the uniformly weighted 367 deep ensembles across all datasets and model architectures with Bayesian ensembles from the 368 literature. Our results shown in Table 1 and figures 1 and A.4 demonstrate that uniformly weighted 369 deep ensembles can match the performance of models based on BMA. On IMDB, a uniform deep 370 ensemble of size  $M \ge 20$  matched the reported performance from Wenzel et al. (2020a), and 371 equivalent results were observed for all other datasets and architectures. On CIFAR-100 with a 372 ResNet110, the uniform deep ensemble outperformed the Bayesian reference from Ashukha et al. 373 (2020) by 1.3ppt. Interestingly, the SGD and deep ensemble results from Ashukha et al. are notably 374 better than ours. As the authors base their Bayesian methods on pre-trained SGD solutions, we 375 can assume that their Bayesian reference performances could also profit from the better accuracies achieved in their training environment. Thus, it is remarkable that our uniform deep ensembles 376 nevertheless matched the Bayesian references, and we would expect an even better performance if 377 compared in the same environment.

In accordance with the literature, bagging reduced the performance, see Appendix A.4.1 for details.
The experiments on the trade-off between *M* and the number of training epochs showed that, even in a sequential setting with a strict limit on the number of overall training epochs, deep ensembles could match the BMA approach, see Appendix A.4.2 for details.

Table 1: Mean test accuracies  $\hat{A}$  over five ensembles, ensemble size is given in brackets (for other ensemble sizes see figures A.4, A.5, A.6, and A.7 in the appendix). The *Simple* results refer to simple deep ensembles with SGD hyperparameters from the reference papers. The subscript  $\rho$  indicates that the weighting of the ensemble members was based on minimizing a PAC-Bayesian tandem bound. The *Bayesian reference* results are taken from Ashukha et al. (2020) (CIFAR ResNet110 & WRN28-10), Band et al. (2021) (EyePACS) and Wenzel et al. (2020a) (IMDB, CIFAR-10 ResNet20).

Model (Dataset)	Experiment	$\hat{A}(AVG_u)$ last	$\begin{array}{l} \hat{A}(\mathrm{AVG}_{\rho}) \\ \mathrm{last} \end{array}$	$\hat{A}(AVG_u)$ all	$\hat{A}(AVG_{\rho})$ all	SGD	Bayesian reference
CNN LSTM (IMDB)	Simple Checkp. SSE	0.871[40] 0.872[40] 0.729[8]	0.872[40] 0.873[40] 0.741[8]	- 0.868[40·5] 0.858[8·10]	- 0.873[40·5] <b>0.874[8·10]</b>	0.853 0.849 0.861	0.870[7]
ResNet20 (CIFAR-10)	Simple SSE	<b>0.938[24]</b> 0.894[24]	<b>0.938[24]</b> 0.894[24]	0.928[24·5] 0.903[24·5]	<b>0.938[24·5]</b> 0.908[24·5]	0.918 0.896	0.936[27]
ResNet110 (CIFAR-10)	Simple SSE	<b>0.956[20]</b> 0.954[20]	<b>0.956[20]</b> 0.955[20]	0.944[20·5] 0.951[20·5]	<b>0.956[20·5]</b> 0.954[20·5]	0.943 0.941	0.955[10]
WRN28-10 (CIFAR-10)	Simple SSE	<b>0.967[24]</b> 0.964[24]	<b>0.967[24]</b> 0.963[24]	0.96[24·5] 0.964[24·5]	<b>0.967[24·5]</b> 0.964[24·5]	0.961 0.956	0.967[10]
ResNet110 (CIFAR-100)	Simple SSE	<b>0.794[24]</b> 0.79[24]	0.793[24] 0.79[24]	0.78[24·5] 0.788[24·5]	0.793[24·5] <b>0.794[24·5]</b>	0.745 0.733	0.781[10]
WRN28-10 (CIFAR-100)	Simple SSE	<b>0.83[24]</b> 0.818[24]	<b>0.83[24]</b> 0.818[24]	0.827[24·5] 0.822[24·5]	0.829[24·5] 0.826[24·5]	0.798 0.791	0.828[10]
ResNet50 (EyePACS)	Simple	0.91[8]	0.91[8]	0.912[8.6]	0.913[8.6]	0.895	0.909[5]

408 409

382

397

410 Bayesian vs. PAC-Bayesian deep ensembles (Hypothesis 2). Optimizing the weighting by min-411 imizing the second order PAC-Bayesian bound matched the uniform performance in all cases 412  $\pm 0.1$ ppt, while even improving the ensemble performance slightly on IMDB (Table 1:  $A(AVG_u)$ ) 413 last vs.  $A(AVG_{o})$  last). At the same time, the PAC generalization guarantees in Table 2 – which 414 still hold after optimization – tightened dramatically, most significantly on CIFAR-10 ResNet20 by 415 31.8ppt, from 0.425 to 0.743. For ResNet110 on CIFAR-100, the optimization was necessary to get 416 a non-trivial bound, increasing from 0.0 to 0.127. The PAC-Bayesian guarantees hold for  $MV_{\rho}$  as 417 aggregation method, and  $AVG_{\rho}$  usually performs slightly better in practice. However, Table A.4 in 418 the appendix shows how similar the two aggregation methods behave. That is, a slight decrease in 419 performance can give rigorous generalization bounds. Figures 2 and 3 exemplify the weighting from minimizing the tandem loss (and a first order) PAC-generalization bound. They show that the second 420 order objective function avoids putting all weight on a single hypothesis. 421

Computing and optimizing the PAC-Bayesian bound requires additional hold-out data, which biases
the comparison with uniform weighting. Instead of using extra data, one could use bagging (Lorenzen
et al., 2019; Masegosa et al., 2020). The results of our bagging experiments are presented in
Appendix A.4.1. For our neural networks ensembles and rather small data sets, bagging decreased
the performance, which is in line with the literature (Lakshminarayanan et al., 2017; Lee et al., 2015).

427

**Bayesian vs. PAC-Bayesian deep snapshot ensembles (Hypothesis 3).** We evaluated snapshot ensembling (Huang et al., 2017) as well as simply taking checkpoints from the original learning rate schedule. In SSEs, going from one to all intermediate checkpoints with optimized weighting improved predictive performance in all but one cases. On IMDB, this improvement was most pronounced with 13.3ppt (Table 1:  $\hat{A}(AVG_{\rho})$  last vs.  $\hat{A}(AVG_{\rho})$  all). An SSE with 8 members (and 10 snapshots per member) gave the best accuracy with a training budget ( $8 \cdot 50 = 400$  epochs) below the budget of the Bayesian reference method (500 epochs). SSEs include networks when the learning rate is lowest, just before re-starting the learning rate cycle. However, including checkpoints following the original schedule also matched the baseline of taking only the last checkpoint across all experiments. This straightforward approach improved performance on IMDB and EyePACS (Table 1:  $A(AVG_{a})$  last vs.  $\hat{A}(AVG_{o})$  all). This is in contrast to the uniform weighting results, where adding the snapshots decreased the performance compared to only taking the last network. As hypothesized, the tandem loss bound optimization allows to include worse performing models in the ensemble (at no additional training cost) while maintaining or improving performance. Figures 2 and 3 show that the second order PAC-Bayesian approach picks snapshots from the different training runs. This can be seen most clearly from the regular pattern in Figure 3, bottom left. Figure 3 shows the differences between the Simple and SSE training setup. In the former, there was a tendency to pick the last snapshots of each network training run, in the latter, which used a cyclic learning rate schedule, also intermediate snapshots were selected. 



Figure 2: IMDB accuracies ( $\hat{A}$ , top) and weight distribution per member for first order (middle, Lorenzen et al., 2019) and tandem bound weighting (bottom) for the *Simple* setting considering only a single network per training process (left) and when adding checkpoints (right). The number of training runs was 40. Four intermediate checkpoints were added, giving a total of five weight configurations per training process.



Figure 3: ResNet20 on CIFAR-10 accuracies (top) and weight distribution per member for first order (middle) and tandem bound weighting (bottom) for the *Simple* setting (left) and *SSE* (right), both including all checkpoints (24 training processes, five checkpoints per process).

## 5 DISCUSSION AND CONCLUSIONS

Bayesian model averaging (BMA) is not meant to maximize predictive performance but to get a
better idea about the maximum likelihood estimate of the parameters of a single model. Still, the
Bayes ensemble is often – implicitly and explicitly – brought forward as an approach to improve
generalization performance. Several lines of research study the generalization performance of the
Bayes ensemble and try to improve it, where performance is measured in terms of accuracy on a
test set. Here we argued conceptually in the line of Ortega et al. (2022) and showed empirically
that the predictive posterior is not a particularly good basis for selecting and weighting the networks

486 Table 2: Test accuracies for AVG<sub> $\rho$ </sub> and MV<sub> $\rho$ </sub> and PAC-Bayesian bounds ( $\delta = 0.05$ ) for all *Simple* 487 experiments (see Table 1). An overview over all experiments is shown in Figure A.4 in the appendix. 488 The bounds were computed for the *all* setting, except for IMDB, where checkpoints were taken in a separate experiment due to early-stopping in *Simple*;  $MV_u$  and  $MV_\rho$  refer to majority voting with 489 uniform and optimized  $\rho$ , respectively. 490

Model	Dataset	$\begin{array}{c} Bound \\ MV_u \end{array}$	Bound $MV_{\rho}$	$\hat{A}(AVG_{\rho})$ last	$\hat{A}(MV_{\rho})$ last	$\hat{A}(AVG_{\rho})$ all	$\hat{A}(MV_{\mu})$ all
CNN LSTM	IMDB	0.576	0.594	0.873	0.872	0.873	0.874
ResNet20	CIFAR-10	0.425	0.743	0.938	0.938	0.938	0.938
ResNet110	CIFAR-10	0.559	0.792	0.956	0.956	0.956	0.955
WRN28-10	CIFAR-10	0.794	0.833	0.967	0.966	0.967	0.966
ResNet110	CIFAR-100	0.0	0.127	0.793	0.791	0.793	0.791
WRN28-10	CIFAR-100	0.172	0.277	0.83	0.829	0.829	0.829
ResNet50	EyePACS	0.686	0.695	0.911	0.909	0.912	0.912

503

504

505 506

in a deep ensemble. A proper Bayesian way to build a deep ensemble would be to apply Bayesian 507 inference to the space of ensembles (Monteith et al., 2011), but this is computationally expensive and 508 - also for that reason - is not what is typically proposed in the research directions our study addresses. A simple uniformly weighted deep ensemble can be expected to perform on par with BMA based 510 approaches and it can be created more efficiently. It does not require sampling from the posterior and 511 training is embarrassingly parallel. 512

Our proposed optimization of the weighting of DNN ensemble members using the tandem loss 513 and additional data can improve uniform deep ensembles. Using a second order bound to derive 514 the optimization objective is crucial, because minimizing a first order bound will lead to a lack 515 of diversity similar to BMA (see also Lorenzen et al., 2019). To our knowledge, the only study 516 applying PAC-Bayesian methods to DNN ensembles is the work by Ortega et al. (2022), who present 517 a PAC-Bayesian generalization bound including a diversity term for different losses. They empirically 518 compared uniform DNN ensembles with ensembles of DNNs that are jointly trained with the newly 519 proposed loss function, directly optimizing the PAC-Bayesian bound during training. This procedure 520 implies that the generalization bounds are not valid anymore after training, as they are not based on independent data. For larger DNNs like ResNet20, the authors find that the direct optimization of 521 their bound leads to ensembles without higher diversity and with similar or even worse performance 522 than uniform deep ensembles with individually trained members – both is in contrast to our results. 523

524 The PAC-Bayesian weighting allows us to include several models from a training run, which can 525 improve performance efficiently and make early-stopping unnecessary. It has already been argued by Sollich and Krogh (1995) that having models overfitted to a subset of the training data in a properly 526 weighted ensemble need not be harmful. While no hold-out data for early-stopping is needed, it is a 527 limitation of the weighting approach that additional hold-out data is required. For the rather small data 528 sets and the models considered in this study, experiments with canonical bagging showed that leaving 529 out training data impaired the performance of the ensemble (as also observed by Lakshminarayanan 530 et al., 2017, and Lee et al., 2015), even though bagging increases diversity. However, the data used 531 for optimization of the weighting provide at the same time a rigorous, non-trivial upper bound on 532 the generalization performance – and getting such formal guarantees generally requires independent data. The PAC-Bayesian bounds still hold after optimization. In our case, this optimization was 534 necessary to elevate some of the bounds from being trivial to non-vacuous guarantees on the ensemble 535 performance. The bound presumes model combination by majority voting as prediction method, 536 however, in our experiments we did not observe a big difference in generalization performance 537 between majority voting and averaging predictive distributions. Thus, PAC-Bayesian optimization of the ensemble weights is highly recommended if not all data available for model construction 538 is needed for training so that a subset can be employed for optimization and getting performance 539 guarantees.

### 540 REPRODUCIBILITY STATEMENT 541

We provide a description of the experimental setup in the main text, as well as details about hyperparameters, augmentation, preprocessing etc. in Appendix A.3. All utilized code is either provided in
the supplementary material or referenced by pointing to public repositories.

### References

545 546

547

550

552

553

554

562

563

564

565

569

570

573

574

580

581

582

586

- B. Adlam, J. R. Snoek, and S. Smith. Cold posteriors and aleatoric uncertainty. In *International Conference on Machine Learning (ICML) Workshop on Uncertainty and Robustness in Deep Learning*, 2020.
  - L. Aitchison. A statistical theory of cold posteriors in deep neural networks. In *International Conference on Learning Representations (ICLR)*, 2021.
  - A. Ashukha, A. Lyzhov, D. Molchanov, and D. Vetrov. Pitfalls of in-domain uncertainty estimation and ensembling in deep learning. In *International Conference on Learning Representations*, 2020.
- G. Bachmann, L. Noci, and T. Hofmann. How tempering fixes data augmentation in Bayesian neural networks.
   In *International Conference on Machine Learning (ICML)*, volume 39, 2022.
- N. Band, T. G. J. Rudner, Q. Feng, A. Filos, Z. Nado, M. W. Dusenberry, G. Jerfel, D. Tran, and Y. Gal.
   Benchmarking Bayesian deep learning on diabetic retinopathy detection tasks. In Advances in Neural
   *Processing Systems (NeurIPS) Datasets and Benchmarks Track (Round 2)*, 2021.
- 560 S. Bernstein. Theory of probability, 1917. In Russian, cited after Kleijn and Van der Vaart (2012).
  - C. M. Bishop and H. Bishop. *Deep Neural Networks*. Springer, 2023.
  - S. Buschjäger, L. Pfahler, and K. Morik. Generalized negative correlation learning for deep ensembling. *arXiv* preprint arXiv:2011.02952, 2020.
- M. Caprio, S. Dutta, K. J. Jang, V. Lin, R. Ivanov, O. Sokolsky, and I. Lee. Credal Bayesian deep learning. arXiv preprint arXiv:2302.09656, 2024.
  - T. Chen, E. Fox, and C. Guestrin. Stochastic gradient Hamiltonian Monte Carlo. In *International Conference on Machine Learning (ICML)*, 2014.
- F. D'Angelo and V. Fortuin. Repulsive deep ensembles are Bayesian. In *Advances in Neural Processing Systems* (*NeurIPS*), volume 34, 2021.
  - T. G. Dietterich. Ensemble methods in machine learning. In *International Workshop on Multiple Classifier* Systems, pages 1–15. Springer, 2000.
- P. Domingos. Bayesian averaging of classifiers and the overfitting problem. In *International Conference on Machine Learning (ICML)*, pages 223–230, 2000.
- E. Dugas, Jared, Jorge, and W. Cukierski. Diabetic retinopathy detection. https://www.kaggle.com/
   competitions/diabetic-retinopathy-detection, 2015. Accessed: 15.04.2024.
  - G. K. Dziugaite and D. M. Roy. Computing nonvacuous generalization bounds for deep (stochastic) neural networks with many more parameters than training data. In *Conference on Uncertainty in Artificial Intelligence* (*UAI*), volume 33, 2017.
- D. E. Eckhardt and L. D. Lee. A theoretical basis for the analysis of multiversion software subject to coincident errors. *IEEE Transactions on Software Engineering*, SE-11(12), 1985.
  - S. Farquhar, M. A. Osborne, and Y. Gal. Radial Bayesian neural networks: Beyond discrete support in large-scale Bayesian deep learning. In *International Conference on Artificial Intelligence and Statistics (AISTATS)*, 2020.
- L. L. Folgoc, V. Baltatzis, S. Desai, A. Devaraj, S. Ellis, O. E. M. Manzanera, A. Nair, H. Qiu, J. Schnabel, and
   B. Glocker. Is MC Dropout Bayesian? *arXiv preprint arXiv:2110.04286*, 2021.
- V. Fortuin, A. Garriga-Alonso, S. W. Ober, F. Wenzel, G. Ratsch, R. E. Turner, M. van der Wilk, and L. Aitchison.
   Bayesian neural network priors revisited. In *International Conference on Learning Representations*, 2022.
- 593 Y. Freund, R. E. Schapire, et al. Experiments with a new boosting algorithm. In *International Conference on Machine Learning (ICML)*, volume 96, 1996.

594 Y. Gal and Z. Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty in deep 595 learning. In International Conference on Machine Learning (ICML), 2016. 596 T. Garipov, P. Izmailov, D. Podoprikhin, D. P. Vetrov, and A. G. Wilson. Loss surfaces, mode connectivity, and 597 fast ensembling of DNNs. In Advances in Neural Processing Systems (NeurIPS), volume 31, 2018. P. Germain, A. Lacasse, F. Laviolette, and M. Marchand. PAC-Bayesian learning of linear classifiers. In International Conference on Machine Learning (ICML), volume 26, 2009. 600 601 P. Germain, A. Lacasse, F. Laviolette, M. March, and J.-F. Roy. Risk bounds for the majority vote: From a 602 PAC-Bayesian analysis to a learning algorithm. Journal of Machine Learning Research, 16(26):787–860, 603 2015. 604 I. Goodfellow, Y. Bengio, and A. Courville. *Deep Learning*. MIT press, 2016. 605 606 P. Grünwald. The safe Bayesian: Learning the learning rate via the mixability gap. In International Conference 607 on Algorithmic Learning Theory, 2012. 608 L. K. Hansen and P. Salamon. Neural network ensembles. IEEE Transactions on Pattern Analysis and Machine 609 Intelligence, 12(10):993-1001, 1990. 610 K. He, X. Zhang, S. Ren, and J. Sun. Deep residual learning for image recognition. In IEEE/CVF Conference on 611 Computer Vision and Pattern Recognition (CVPR), 2016. 612 613 G. Huang, Y. Li, G. Pleiss, Z. Liu, J. E. Hopcroft, and K. Q. Weinberger. Snapshot ensembles: Train 1, get m for 614 free. In International Conference on Learning Representations (ICLR), 2017. 615 P. Izmailov, S. Vikram, M. D. Hoffman, and A. G. Wilson. What are Bayesian neural network posteriors really 616 like? In International Conference on Machine Learning (ICML), volume 38, 2021. 617 Z. Jiang, H. Liu, B. Fu, and Z. Wu. Generalized ambiguity decompositions for classification with applications in 618 active learning and unsupervised ensemble pruning. In AAAI Conference on Artificial Intelligence, volume 31, 619 2017. 620 621 S. Kapoor, W. J. Maddox, P. Izmailov, and A. G. Wilson. On uncertainty, tempering, and data augmentation in Bayesian classification. In Advances in Neural Processing Systems (NeurIPS), volume 35, 2022. 622 623 A. Kendall and Y. Gal. What uncertainties do we need in Bayesian deep learning for computer vision? In 624 Advances in Neural Processing Systems (NeurIPS), volume 30, 2017. 625 B. J. Kleijn and A. W. Van der Vaart. The Bernstein-von-Mises theorem under misspecification. Electronic journal 626 of Statistics, 6:354–381, 2012. 627 628 A. Krizhevsky. Learning multiple layers of features from tiny images. https://www.cs. toronto. edu/ kriz/learningfeatures-2009-TR. pdf, 2009. 629 630 A. Krogh and P. Sollich. Statistical mechanics of ensemble learning. *Physical Review E*, 55(1):811, 1997. 631 A. Lacasse, F. Laviolette, M. Marchand, P. Germain, and N. Usunier. PAC-Bayes bounds for the risk of the 632 majority vote and the variance of the gibbs classifier. In Advances in Neural Processing Systems (NeurIPS), 633 volume 19, 2006. 634 635 B. Lakshminarayanan, A. Pritzel, and C. Blundell. Simple and scalable predictive uncertainty estimation using deep ensembles. In Advances in Neural Processing Systems (NeurIPS), volume 30, 2017. 636 637 J. Langford and J. Shawe-Taylor. PAC-Bayes & margins. In Advances in Neural Processing Systems (NeurIPS), 638 volume 15, 2002. 639 P. S. Laplace. Théorie analytique des probabilités, volume 7. Courcier, 1820. 640 641 L. Le Cam. On some asymptotic properties of maximum likelihood estimates and related Bayes' estimates. 642 Univ. Calif. Publ. in Statist., 1:277-330, 1953. 643 S. Lee, S. Purushwalkam, M. Cogswell, D. Crandall, and D. Batra. Why m heads are better than one: Training a 644 diverse ensemble of deep networks. arXiv preprint arXiv:1511.06314, 2015. 645 646 C. Li, C. Chen, D. Carlson, and L. Carin. Preconditioned stochastic gradient Langevin dynamics for deep 647 neural networks. In Association for the Advancement of Artificial Intelligence (AAAI) Conference on Artificial

Intelligence, volume 30, 2016.

648 649	Y. Liu and X. Yao. Ensemble learning via negative correlation. Neural Networks, 12(10):1399–1404, 1999.
650 651	S. S. Lorenzen, C. Igel, and Y. Seldin. On PAC-Bayesian bounds for random forests. <i>Machine Learning</i> , 108(8): 1503–1522, 2019.
652 653	YA. Ma, T. Chen, and E. Fox. A complete recipe for stochastic gradient mcmc. In <i>Advances in Neural Processing Systems (NeurIPS)</i> , volume 28, 2015.
655 656 657	A. L. Maas, R. E. Daly, P. T. Pham, D. Huang, A. Y. Ng, and C. Potts. Learning word vectors for sentiment analysis. In <i>Annual Meeting of the Association for Computational Linguistics: Human Language Technologies</i> , volume 49, 2011.
658	D. J. MacKay. Bayesian interpolation. Neural Computation, 4(3):415-447, 1992a.
659 660 661	D. J. MacKay. A practical Bayesian framework for backpropagation networks. <i>Neural Computation</i> , 4(3): 448–472, 1992b.
662 663	W. J. Maddox, P. Izmailov, T. Garipov, D. P. Vetrov, and A. G. Wilson. A simple baseline for Bayesian uncertainty in deep learning. In <i>Advances in Neural Processing Systems (NeurIPS)</i> , volume 32, 2019.
664 665	A. R. Masegosa. Learning under model misspecification: Applications to variational and ensemble methods. In <i>Advances in Neural Processing Systems (NeurIPS)</i> , volume 33, 2020.
666 667 668	A. R. Masegosa, S. S. Lorenzen, C. Igel, and Y. Seldin. Second order PAC-Bayesian bounds for the weighted majority vote. In <i>Advances in Neural Processing Systems (NeurIPS)</i> , 2020.
669 670	D. A. McAllester. Some PAC-Bayesian theorems. In <i>Conference on Computational Learning Theory (COLT)</i> , volume 11, 1998.
671 672	K. Monteith, J. L. Carroll, K. Seppi, and T. Martinez. Turning Bayesian model averaging into Bayesian model combination. In <i>International Joint Conference on Neural Networks (IJCNN)</i> , pages 2657–2663, 2011.
673 674 675 676	S. Nabarro, S. Ganev, A. Garriga-Alonso, V. Fortuin, M. van der Wilk, and L. Aitchison. Data augmentation in Bayesian neural networks and the cold posterior effect. In <i>Conference on Uncertainty in Artificial Intelligence</i> , 2022.
677 678 679 680	Z. Nado, N. Band, M. Collier, J. Djolonga, M. Dusenberry, S. Farquhar, A. Filos, M. Havasi, R. Jenatton, G. Jerfel, J. Liu, Z. Mariet, J. Nixon, S. Padhy, J. Ren, T. Rudner, Y. Wen, F. Wenzel, K. Murphy, D. Sculley, B. Lakshminarayanan, J. Snoek, Y. Gal, and D. Tran. Uncertainty Baselines: Benchmarks for uncertainty & robustness in deep learning. <i>arXiv preprint arXiv:2106.04015</i> , 2021.
681 682	R. M. Neal. Bayesian learning via stochastic dynamics. In Advances in Neural Processing Systems (NeurIPS), volume 5, 1992.
683 684	R. M. Neal. Bayesian Learning for Neural Networks. Springer, 1996.
685 686	R. M. Neal. MCMC using Hamiltonian dynamics. Handbook of Markov Chain Monte Carlo, 2(11):113–162, 2011.
687 688 689	L. A. Ortega, R. Cabañas, and A. Masegosa. Diversity and generalization in neural network ensembles. In <i>International Conference on Artificial Intelligence and Statistics (AISTATS)</i> , volume 25, pages 11720–11743, 2022.
690 691 692	T. Pearce, F. Leibfried, and A. Brintrup. Uncertainty in neural networks: Approximately Bayesian ensembling. In <i>International Conference on Artificial Intelligence and Statistics</i> , 2020.
693 694	M. Pérez-Ortiz, O. Rivasplata, J. Shawe-Taylor, and C. Szepesvári. Tighter risk certificates for neural networks. <i>Journal of Machine Learning Research</i> , 22(227):1–40, 2021.
695 696 697	M. P. Perrone and L. N. Cooper. When networks disagree: Ensemble methods for hybrid neural networks. In R. J. Mammone, editor, <i>Neural networks for speech and image processing</i> , pages 81–99. Chapman-Hall, 1993.
698 699 700	H. Ritter, A. Botev, and D. Barber. A scalable Laplace approximation for neural networks. In <i>International Conference on Learning Representations (ICLR)</i> , volume 6, 2018.
701	T. G. Rudner, Z. Chen, Y. W. Teh, and Y. Gal. Tractable function-space variational inference in Bayesian neural networks. In <i>Advances in Neural Processing Systems (NeurIPS)</i> , volume 35, 2022.

- P. Sollich and A. Krogh. Learning with ensembles: How overfitting can be useful. In *Advances in Neural Information Processing Systems (NeurIPS)*, volume 8, 1995.
  - N. Thiemann, C. Igel, O. Wintenberger, and Y. Seldin. A strongly quasiconvex PAC-Bayesian bound. In International Conference on Algorithmic Learning Theory, 2017.
- M. Vadera, J. Li, A. Cobb, B. Jalaian, T. Abdelzaher, and B. Marlin. URSABench: A system for comprehensive benchmarking of Bayesian deep neural network models and inference methods. In *Machine Learning and Systems (MLSys)*, volume 4, 2022.
- L. G. Valiant. A theory of the learnable. *Communications of the ACM*, 27(11):1134–1142, 1984.
- A. W. Van der Vaart. *Asymptotic Statistics*, volume 3. Cambridge University Press, 2000.
- 713 R. von Mises. Wahrscheinlichkeitsrechnung. Springer, 1931.
- M. Welling and Y. W. Teh. Bayesian learning via stochastic gradient Langevin dynamics. In *International Conference on Machine Learning (ICML)*, volume 28, 2011.
- F. Wenzel, K. Roth, B. S. Veeling, J. Swiatkowski, L. Tran, S. Mandt, J. Snoek, T. Salimans, R. Jenatton, and
  S. Nowozin. How good is the Bayes posterior in deep neural networks really? In *International Conference on Machine Learning (ICML)*, volume 37, 2020a.
- F. Wenzel, J. Snoek, D. Tran, and R. Jenatton. Hyperparameter ensembles for robustness and uncertainty quantification. In *Advances in Neural Processing Systems (NeurIPS)*, volume 33, 2020b.
  - J. G. Wiese, L. Wimmer, T. Papamarkou, B. Bischl, S. Günnemann, and D. Rügamer. Towards efficient MCMC sampling in Bayesian neural networks by exploiting symmetry. In *Machine Learning and Knowledge Discovery in Databases: Research Track*, 2023.
  - A. G. Wilson and P. Izmailov. Bayesian deep learning and a probabilistic perspective of generalization. In *Advances in Neural Processing Systems (NeurIPS)*, volume 33, pages 4697–4708, 2020.
- Y.-S. Wu, A. R. Masegosa, S. S. Lorenzen, C. Igel, and Y. Seldin. Chebyshev-Cantelli PAC-Bayes-Bennett inequality for the weighted majority vote. In *Advances in Neural Processing Systems (NeurIPS)*, 2021.
- A. Yenter and A. Verma. Deep CNN-LSTM with combined kernels from multiple branches for IMDb review sentiment analysis. In *IEEE Annual Ubiquitous Computing, Electronics and Mobile Communication Conference (UEMCON)*, volume 8, 2017.
- S. Zagoruyko and N. Komodakis. Wide residual networks. In *British Machine Vision Conference (BMVC)*, 2016.
- R. Zhang, C. Li, J. Zhang, C. Chen, and A. G. Wilson. Cyclical stochastic gradient MCMC for Bayesian deep learning. In *International Conference on Learning Representations (ICLR)*, 2020.
- 737

704

705

706

711

714

722

723

724

725

726

727

740

741 742

743

744 745

- 746
- 747 748

749

750

- 752
- 753
- 754
- 755

## 756 A APPENDIX / SUPPLEMENTAL MATERIAL

# 758 A.1 EXPERIMENTAL DETAILS 759

The experimental setups including hyperparameters were taken from the respective cited references
and corresponding source code. Thus, no hyperparameter tuning was performed, except for the
EyePACS dataset, where the hyperparameters from the literature did not lead to convergence of the
models and the learning rate was decreased by a factor of 100.

764

785 786

787 788 789

**IMDB.** For IMDB, the described setup from Wenzel et al. (2020a) was copied as closely as possible. 765 The dataset from the tensorflow.keras.datasets API was used with 20,000 word features 766 and a maximum sequence length of 100. 20,000 training samples with 5,000 random validation 767 samples being used for training with early-stopping due to overfitting. The test set included all 25,000 768 samples from the original test set. The CNN LSTM model and code were taken from the Keras 769 example<sup>3</sup> and extended with a Gaussian prior  $\mathcal{N}(0, I)$  for regularization, as done by the authors in 770 their code repository<sup>4</sup>. As optimizer, SGD with Nesterov momentum of 0.98 and a batch size of 32 771 with a constant learning rate was utilized. For the cyclic learning rate in all snapshot ensemble (SSE) 772 experiments, the original schedule from Huang et al. (2017) was used. 773

774 CIFAR. The CIFAR datasets were taken from the tensorflow.keras.datasets API with 775 the original train and test split. Similar to Wenzel et al. (2020a), no early-stopping with a validation 776 set was employed. The ResNet20 model was taken from the Keras example  $^{5}$ . As in the example, all experiments used data augmentation during training with random left/right flipping and random 777 cropping with 4px of shift horizontally and vertically. The hyperparameters for the ResNet20 were 778 adopted from Wenzel et al. (2020a), while they were taken from Ashukha et al. (2020) for the 779 ResNet110 and Wide ResNet28-10. As optimizer, SGD with Nesterov momentum of 0.9 and a 780 batch size of 128 with a step-wise decreasing learning rate ( $\eta$ ) was employed (epoch,  $\eta$ -multiplier): 781 (80, 0.1), (120, 0.01), (160, 0.001), (180, 0.0005). For the ResNet110 and WRN28-10, a linearly 782 decreasing learning rate schedule starting at half of the total number of epochs was utilized, as 783 reported by Ashukha et al. (2020) and introduced by Garipov et al. (2018): 784

 $\eta(i) = \begin{cases} \eta_{\text{init}}, & i \in [0, 0.5 \cdot \text{epochs}] \\ \eta_{\text{init}} \cdot (1.0 - 0.99 \cdot (i/\text{epochs} - 0.5)/0.4), & i \in [0.5 \cdot \text{epochs}, 0.9 \cdot \text{epochs}] \\ \eta_{\text{init}} \cdot 0.01, & \text{otherwise} \end{cases}$ (A.6)

EyePACS. The EyePACS dataset was included in the 2015 Kaggle diabetic retinopathy detection 790 competition (Dugas et al., 2015). Diabetic retinopathy is the leading cause of blindness in the 791 working-age population of the developed world and estimated to affect over 92 million people. The 792 dataset contains high-resolution labeled RGB images of human retinas with varying degrees of 793 diabetic retinopathy on a five-grade scale, from none (0), to mild (1), moderate (2), severe (3) and 794 proliferating (4) development of the disease. It consists of 35126 training, 10906 validation and 795 42670 test images, each labeled by a medical expert. In order to binarize the labels, Band et al. (2021) 796 follow previous work and classify moderate or worse manifestation as sight-threatening (grades 2-4), 797 and remaining grades 0-1 as non sight-threatening. The dataset with binary labels is unbalanced, 798 such that 19.6% of the training and 19.2% of the test set have a positive label, which is why the 799 cross-entropy objective is weighted by the inverse of the global class distribution. Furthermore, the 800 images exhibit different kinds of noise (artifacts, focus, exposure) and are expected to show label noise due to misjudgement of the medical personnel (Dugas et al., 2015). The preprocessing of 801 images by Band et al. (2021) follows the winning entry of the original Kaggle challenge (Dugas et al., 802 2015). The images are first rescaled such that the retinas have a radius of 300 pixels, smoothed using 803 local Gaussian blur with a kernel standard deviation of 100 pixels and clipped to 90% to remove 804 boundary effects. Finally, they are resized to 512x512 pixels and stored. For our preprocessing, the 805 implementation from Nado et al. (2021) was used<sup>6</sup>. Except the preprocessing, no data augmentation 806

808 <sup>4</sup> https://github.com/google-research/google-research/tree/master/cold\_posterior\_bnn

<sup>5</sup>https://github.com/keras-team/keras/blob/1a3ee8441933fc007be6b2beb47af67998d50737/examples/cifar10\_resnet.py

<sup>&</sup>lt;sup>3</sup>https://github.com/keras-team/keras/blob/1a3ee8441933fc007be6b2beb47af67998d50737/examples/imdb\_cnn\_lstm.py

<sup>&</sup>lt;sup>6</sup>https://github.com/google/uncertainty-baselines/blob/main/uncertainty\_baselines/datasets/diabetic\_retinopathy\_dataset\_utils.py

Model (Dataset)	Experiment	Val. set	$lr_{ m init}$	СР	epochs	L2-reg.	LR Sched
	Simple	20%		1	50		constant
IMDB	Checkp.	20%	0.1	5	5	$\mathcal{N}(0,I)$	constant
(CNN LSTM)	Bagging	Bagging		1	50		constant
	SSE	20%			50		cyclic
	Ep.b.	20%		1	variable		constant
	Simple	0%	0.1	5	200		step
ResNet20 (CIFAR-10)	Bagging	Bagging	0.1	5	200	0.002	step
	SSE	0%	0.2	5	200		cyclic
	Ep.b.	0%	0.1	1	variable		step
ResNet110,	Simple	0%		5	300	0.0003	linear
WRN28-10 (CIFAR-10,	Bagging	Bagging	0.1	5	300		linear
	SSE	0%		5	300		cyclic
CIFAR-100)	Ep.b.	0%		1	variable		linear
ResNet50 (EyePACS)	Simple Bagging	0% Bagging	$2.3\cdot 10^{-4}$	6	90	$1.07\cdot 10^{-4}$	step

Table A.3: Hyperparameters for all experiments. *Ep.b.* (Epoch budget) refers to the experiments in section A.4.2.

835

836 837

827

was used. With the optimal SGD hyperparameters from Band et al. (2021), the resulting models
failed to converge in our case, which is why we reduced the learning rate and switched the optimizer
to Adam. The batch size was kept at 32. For the learning rate schedule, a step-wise decrease was
employed with a reduction by factor 0.2 at 30 and 60 epochs.

### A.2 TEST-TIME CROSS-VALIDATION

Ashukha et al. (2020) describe the problem of requiring a validation set (in their case for scaling the 838 logit outputs of neural networks with a temperature parameter), while the benchmark datasets only 839 feature a training and test dataset (e.g. CIFAR). In that case, when splitting the training set into a 840 training and validation subset, the performance may drop compared to methods that make use of the 841 full training data because the reduced training data set is a worse description of the task. In contrast, 842 when splitting off a validation set from the test data, one still obtains an unbiased estimate of the 843 generalization error, but with higher variance due to smaller test dataset size. In order to reduce the 844 variance, Ashukha et al. (2020) propose test-time cross-validation: splitting the test dataset randomly 845 and averaging the results of the generalization estimates to reduce variance. We follow this approach and divide the test dataset randomly in half. One half is used for the tandem loss bound optimization, 846 while the other serves as a generalization estimate. Afterwards, the datasets' roles are switched, and 847 the risk estimates are averaged. In our case, this setting does not result in a fair comparison. All 848 methods should make use of the same amount of data, and algorithms that do not need a validation 849 set could use the additional data for training, which can improve performance, in particular if data 850 are scarce. However, one can envision a scenario where the additional data are not available during 851 training but only later when deploying the model (e.g., local fine-tuning of centrally trained models). 852

### 853 854

855

856

858 859

860 861

### A.3 CANCELLATION OF INDEPENDENT ERRORS

Because this known result is often stated without proof, we provide an upper bound for the error probability of the majority vote for an ensemble of binary classifiers with independent errors. We consider M binary classifiers  $h_i, \ldots, h_M$  mapping to  $\{0, 1\}$  and the majority vote classifier given by:

$$h_{\rm MV}(x) = \begin{cases} 1 & \text{if } \sum_{i=1}^{M} h_i(x) \ge n/2\\ 0 & \text{otherwise} \end{cases}$$
(A.7)

We define the random variables  $A_i = \mathbb{1}[h_i(x) \neq y]$  indicating a mistake by  $h_i$  and  $S_M = \sum_{i=1}^M A_i$ . The probability  $P(h_{MV}(x) \neq y)$  that the majority vote classifier makes a mistake is the probability that  $P(S_M > M/2)$  or alternatively  $P(S_M \ge (M+1)/2)$ . We assume for all i = 1, ..., M that the risk is bounded by a constant

$$P_{(x,y)\sim p}(h_i(x) \neq y) = \mathbb{E}\{A_i\} \le p_{\max} < \frac{1}{2}$$
 (A.8)

and that the  $A_i$  are independent. Then

$$\mathbb{E}\left\{\sum_{i=1}^{M} A_i\right\} = \mathbb{E}\{S_M\} \le Mp_{\max} \tag{A.9}$$

and we have

866 867 868

870 871 872

877 878 879

882 883

885

$$P(S_M \ge (M+1)/2) = P(S_M - \mathbb{E}\{S_M\} \ge (M+1)/2 - \mathbb{E}\{S_M\})$$
(A.10)

$$\leq P(S_M - \mathbb{E}\{S_M\} \geq (M+1)/2 - Mp_{\max}).$$
 (A.11)

With Hoeffding's inequality and  $\varepsilon = (M+1)/2 - Mp_{max}$  we get the desired bound

$$P(S_M \ge (M+1)/2) \le P\{S_M - \mathbb{E}\{S_M\} \ge \varepsilon\} \le \exp\left(-\frac{2\varepsilon^2}{M}\right)$$

$$= \exp\left(-\frac{2\left(\frac{M+1}{2} - Mp_{\max}\right)^2}{M}\right).$$
(A.12)

### A.4 FURTHER EXPERIMENTAL RESULTS

### 886 A.4.1 BAGGING

In the bagging experiments, we used different training data for each ensemble member to increase diversity. We drew bootstrap samples from the training data uniformly at random with replacement. The sample size was equal to the size of the training data. The experiments illustrate the effect of decreasing the training data volume (e.g., to use the data for bound optimization).

The individual network performance suffered from omitting unique training data points, and although 892 the networks are assumed to be more diverse, the resulting ensembles performed worse performance 893 across all experiments. The results are visualized in Figure A.6, which shows that none of the 894 ensembles, neither uniform nor optimized in their weighting, could match the Bayesian reference. 895 This is in line with the literature (Lakshminarayanan et al., 2017; Lee et al., 2015), which finds the 896 same decrease in performance for neural network ensembles. PAC-Bayesian optimization is therefore 897 best suited for large datasets with additional data that is not essential for training. On a positive note, random initialization and stochastic training appears to be enough randomization for creating diverse ensembles. 899

900 901

### A.4.2 TRAINING TIME COMPARISON

It is not straight-forward to compare the computational resources required by the different approaches due to their different training and inference procedures. The BMA based approaches require training of the neural network and sampling the posterior, where the latter can be time consuming. In contrast, minimization of the PAC-Bayesian bound is highly efficient and its computational costs can be neglected. Simple deep ensembles are embarrassingly parallel, while methods that consider ensemble members from one process cannot be fully parallelized.

908 But even for simple ensembles, there can be a trade-off between number of networks and the number 909 of training iterations spent on each network. To explore this trade-off, we considered the worst 910 case scenario for deep ensemble training and assume that there is no parallelization. Furthermore, 911 we assumed that the sampling from the posterior does not take any time. That is, for methods 912 based on BMA only the training epochs for optimizing the network are counted, not the sampling. 913 In this sequential setting, we asked, given an overall budget of B training epochs, how does the 914 performance of our deep ensembles depend on the number M of networks when each network is 915 trained for |B/M| epochs. For the experiments taken from Wenzel et al. (2020a), we set B = 1500for CIFAR-10 and B = 500 for IMDB, corresponding to the number of training epochs used for the 916 Bayes ensemble. The results are shown in Figure A.7. Even in this biased comparison, there were 917 choices of M for which the simple deep ensemble outperformed the Bayesian approach.

918	Table A.4: Test accuracies for AVG <sub><math>\rho</math></sub> and MV <sub><math>\rho</math></sub> and PAC-Bayesian bounds for all experiments. The
919	bounds were computed for the all setting, except for IMDB, where checkpoints were taken in a
920	separate experiment due to early-stopping in Simple; $MV_u$ and $MV_\rho$ refer majority voting with
921	uniform and optimized $\rho$ , respectively.

Model (Dataset)	Experiment	Uniform Bound	TND Bound	$\hat{A}(AVG_{\rho})$ last	$\hat{A}(\mathrm{MV}_{\rho})$ last	$\hat{A}(AVG_{\rho})$ all	$\hat{A}(MV_{\rho})$ all
	Simple	0.585	0.59	0.872	0.873	-	-
CNN LSTM	Checkpointing	0.576	0.594	0.873	0.872	0.873	0.874
(IMDB)	Bagging	0.579	0.585	0.869	0.867	-	-
	SSE	0.459	0.575	0.741	0.74	0.874	0.872
	Ep.b. (500)	0.571	0.574	0.873	0.872	-	-
	Simple	0.425	0.743	0.938	0.938	0.938	0.938
ResNet20	Bagging	0.387	0.691	0.929	0.928	0.929	0.928
(CIFAR-10)	SSE	0.593	0.612	0.894	0.893	0.908	0.909
	Ep.b. (1500)	0.718	0.718	0.937	0.934	-	-
	Simple	0.559	0.792	0.956	0.956	0.956	0.955
D N - +110	Bagging	0.555	0.763	0.949	0.948	0.949	0.948
(CIEAD 10)	SSE	0.785	0.802	0.955	0.954	0.954	0.953
(CIFAK-10)	Ep.b. (1500)	0.799	0.8	0.954	0.953	-	-
	Ep.b. (300)	0.735	0.735	0.942	0.939	-	-
	Simple	0.794	0.833	0.967	0.966	0.967	0.966
Wide ResNet28-10 (CIFAR-10)	Bagging	0.764	0.808	0.959	0.958	0.959	0.958
	SSE	0.843	0.848	0.963	0.964	0.964	0.964
	Ep.b. (1500)	0.832	0.834	0.966	0.965	-	-
	Ep.b. (300)	0.805	0.806	0.961	0.957	-	-
	Simple	0.0	0.127	0.793	0.791	0.793	0.791
ResNet110 (CIFAR-100)	Bagging	0.0	0.0	0.771	0.769	0.771	0.768
	SSE	0.083	0.125	0.79	0.787	0.794	0.793
	Ep.b. (1500)	0.091	0.099	0.783	0.777	-	-
	Ep.b. (300)	0.0	0.0	0.749	0.728	-	-
Wide ResNet28-10	Simple	0.172	0.277	0.83	0.829	0.829	0.829
	Bagging	0.041	0.173	0.812	0.811	0.812	0.811
	SSE	0.276	0.289	0.818	0.818	0.826	0.827
(CIFAR-100)	Ep.b. (1500)	0.271	0.276	0.829	0.822	-	-
	Ep.b. (300)	0.233	0.233	0.815	0.811	-	-
ResNet50	Simple	0.686	0.695	0.911	0.909	0.912	0.912
(EyePACS)	Raooino	0.613	0.638	0.894	0.893	0.896	0.895

### A.5 LICENSES

Our work built on several, publicly available code repositories and datasets. The IMDB training code including the CNN LSTM model was based on the Keras example, similarly to the CIFAR training code from the CIFAR-10 Keras example, including all ResNets. Both examples are distributed under the Apache 2.0 license. The Gaussian prior for the CNN LSTM on IMDB was based on the code from Wenzel et al. (2020a), also under Apache 2.0 license. The preprocessing routine of the EyePACS dataset was taken and adapted from Google's uncertainty-baselines repository, again under Apache 2.0 license. The Wide ResNet (Zagoruyko and Komodakis, 2016) implementation was taken from a publicly available GitHub repository. Finally, the implementation of the tandem loss PAC-Bayesian bound optimization (Masegosa et al., 2020) was taken and extended from the official implementation by the authors. 



Figure A.4: Mean test accuracy  $\dot{A} \pm \sigma$  over five ensembles for uniformly and PAC-Bayesian weighted deep ensembles (AVG<sub>u</sub> and AVG<sub>ρ</sub>), using only the last or all training checkpoints, and best single model (SGD). References are Bayesian ensembles cSGHMC-ap (Wenzel et al., 2020a), cSGLD (Ashukha et al., 2020), *MC-Dr*opout and *MC-Dr*opout *Deep Ensemble* (i.e., ensemble of Bayesian ensembles), as well as simple *Deep Ensembles* for EyePACS from Band et al. (2021). Numbers in brackets indicate the ensemble sizes.

1021

1022

1023

1024









Figure A.7: Epoch budget in a sequential training setting: Mean ensemble accuracy (5 ensembles,  $\pm \sigma$ ). The number of epochs is given by  $\lfloor B/M \rfloor$ , where b = 1500.