

---

# Disentangling the Roles of Curation, Data-Augmentation and the Prior in the Cold Posterior Effect

---

Lorenzo Noci<sup>\*1</sup> Kevin Roth<sup>\*1</sup> Gregor Bachmann<sup>\*1</sup> Sebastian Nowozin<sup>2</sup> Thomas Hofmann<sup>1</sup>

## Abstract

The “cold posterior effect” (CPE) in Bayesian deep learning describes the disturbing observation that the predictive performance of Bayesian neural networks can be significantly improved if the Bayes posterior is artificially sharpened using a temperature parameter  $T < 1$ . The CPE is problematic in theory and practice and since the effect was identified many researchers have proposed hypotheses to explain the phenomenon. However, despite this intensive research effort the effect remains poorly understood. In this work we provide novel and nuanced evidence relevant to existing explanations for the cold posterior effect, disentangling three hypotheses: 1. The *dataset curation hypothesis* of (Aitchison, 2020): we show empirically that the CPE does not arise in a real curated data set but can be produced in a controlled experiment with varying curation strength. 2. The *data augmentation hypothesis* of (Izmailov et al., 2021) and (Fortuin et al., 2021): we show empirically that data augmentation is sufficient but not necessary for the CPE to be present. 3. The *bad prior hypothesis* of (Wenzel et al., 2020): we use a simple experiment evaluating the relative importance of the prior and the likelihood, strongly linking the CPE to the prior. Our results demonstrate how the CPE can arise in isolation from synthetic curation, data augmentation, and bad priors. Cold posteriors observed “in the wild” are therefore unlikely to arise from a single simple cause; as a result, we do not expect a simple “fix” for cold posteriors.

## 1. Introduction

For a general *introduction to Bayesian deep learning* see Section A in the Appendix. The “cold posterior effect” (CPE) states that among all tempered posteriors  $p(\boldsymbol{\theta}|\mathcal{D})^{1/T}$  the best posterior predictive performance on holdout data is achieved at temperature  $T < 1$  (Wenzel et al., 2020). Formally, tempering the posterior corresponds to a  $1/T$ -scaling of the potential energy function  $U(\boldsymbol{\theta})$ ,

$$p(\boldsymbol{\theta}|\mathcal{D})^{1/T} \propto \exp\left(-\frac{1}{T}U(\boldsymbol{\theta})\right), \quad (1)$$

where  $U(\boldsymbol{\theta}) = -\sum_{i=1}^n \log p(y_i|\mathbf{x}_i, \boldsymbol{\theta}) - \log p(\boldsymbol{\theta})$ , i.e., both the log-likelihood and the log-prior are scaled by  $1/T$ . For  $T = 1$  we have the Bayes posterior, whereas for  $T \rightarrow 0$  we obtain a sequence of distributions which have their mass more and more confined around the MAP mode of the distribution (Leimkuhler et al., 2019). We can thus think of the  $T \rightarrow 0$  limit of posterior inference as MAP estimation. See Section B in the Appendix for a discussion of the *difference between tempering the posterior and tempering the likelihood only*.

There are three main building blocks in Bayesian deep learning that may be at fault for the CPE to emerge: (i) *model misspecification*: the likelihood model could be misspecified (Wenzel et al., 2020; Adlam et al., 2020; Aitchison, 2020; Zeno et al., 2020), (ii) *bad priors*: the priors currently used in deep BNNs may be inadequate (Wenzel et al., 2020; Fortuin et al., 2021), or (iii) *inaccurate inference*: the inference method might not yield an accurate enough approximation to the true posterior (Wenzel et al., 2020; Adlam et al., 2020; Fortuin et al., 2021; Izmailov et al., 2021). Next we review some of the most prominent hypotheses for the emergence of the CPE. For an *extended discussion of these hypotheses* see Section B in the Appendix.

**Inaccurate Inference Hypothesis:** To rule out obvious problems with the inference mechanism, (Wenzel et al., 2020) proposed a set of diagnostics, based on comparing ensemble statistics to their theoretically known values. We have closely monitored these diagnostics in our experiments, however, despite our own extensive efforts to ensure accurate inference, we cannot exclude the possibility that our inference may be inaccurate although the diagnostics match.

---

<sup>\*</sup>Equal contribution <sup>1</sup>Dept of Computer Science, ETH Zürich <sup>2</sup>Microsoft Research, Cambridge, UK. Correspondence to: Lorenzo Noci <lorenzo.noci@inf.ethz.ch>, Kevin Roth <kevin.roth@inf.ethz.ch>, Gregor Bachmann <gregor.bachmann@inf.ethz.ch>.

To investigate the inference hypothesis (Izmailov et al., 2021) recently applied full-batch Hamiltonian Monte Carlo (HMC)<sup>1</sup>, which is considered to be the gold-standard in terms of inference accuracy, to the models of (Wenzel et al., 2020), showing that with HMC inference and when data augmentation is turned off the CPE disappears. While these results can easily be misconstrued as evidence that the approximate inference is at fault in the CPE, there are good reasons to think otherwise.

For instance, (Izmailov et al., 2021) show, using the code of (Wenzel et al., 2020), that turning off data augmentation alone is sufficient to remove the CPE (cf. Table 7 in (Izmailov et al., 2021) Appendix G). We too can confirm that the CPE does not arise with SG-MCMC based inference applied to (Wenzel et al., 2020)’s models when data augmentation is turned off, cf. Figure 1 in Section 3. From this we can already conclude that either SG-MCMC inference is accurate enough for this specific setting, or inaccurate inference is not necessary for the CPE to emerge<sup>2</sup>.

A more direct counter argument follows from (Adlam et al., 2020), which have demonstrated that there can be a CPE in Gaussian Processes (GP) regression, where the posterior has a closed form solution, provided that the aleatoric uncertainty is overestimated. Hence, the CPE can arise in a setting where exact inference is possible. Thus, while inaccurate inference may be sufficient, it does not appear to be necessary for the CPE to emerge.

**Curation Hypothesis (model misspecification):** (Aitchison, 2020) devises a theory that attributes the effectiveness of tempering to the fact that standard benchmark datasets such as CIFAR-10 are carefully curated and that we should take this curation into account by tempering BNNs. In (Aitchison, 2020)’s model of curation (actual dataset curation may differ), a datapoint  $\mathbf{x}$  is added to the dataset if and only if *all*  $S$  labellers independently agree on the label  $y_s$  to be assigned to  $\mathbf{x}$ , while the datapoint is discarded if at least one pair of labellers  $s, s'$  disagrees  $y_s \neq y_{s'}$ . The main argument put forward by (Aitchison, 2020) is that if we a priori know that the dataset is curated in the sense described above, we should take this into account in our likelihood model. The proposed likelihood to be used in the case of curation should then be of the following form

$$p(y, \mathbf{x} | \boldsymbol{\theta}) \propto p(\{y_s = y\}_{s=1}^S | \mathbf{x}, \boldsymbol{\theta}) = \prod_s p(y_s = y | \mathbf{x}, \boldsymbol{\theta}) \quad (2)$$

Thus, assuming that the labellers are i.i.d., the probability of consensus on label  $y$  is

$$p(y, \mathbf{x} | \boldsymbol{\theta}) \propto p(y | \mathbf{x}, \boldsymbol{\theta})^S, \quad (3)$$

<sup>1</sup>(Izmailov et al., 2021) parallelized the HMC computation over hundreds of Tensor Processing Units (TPUs)

<sup>2</sup>We can also conclude that data augmentation is sufficient for the CPE to arise (although it is not necessary, as we will see later).

which corresponds to a cold posterior where only the log-likelihood is re-scaled while the prior is not. (Aitchison, 2020) argues that we should observe  $S \approx 1/T$ , i.e., the optimal temperature should roughly be inversely proportional to the total number of labellers involved in the curation.

Note that (Aitchison, 2020)’s model of curation only includes a data point if *all* labellers agree, but in practice, datasets are often collected with some tolerance of labeller disagreement, e.g. in that a datapoint is included if a certain fraction of labellers agree. However, a more realistic (weaker) model of curation, that filters out fewer “hard” instances, would only give rise to a weaker, less pronounced CPE, see Section B in the Appendix for more details.

**Data Augmentation Hypothesis (model misspecification):** Current deep learning practices use a number of techniques that technically do not obey the likelihood principle (see Appendix K in (Wenzel et al., 2020) for a discussion of so-called “dirty likelihoods”). The data augmentation hypothesis specifically says that the CPE is largely an artifact of using data augmentation and that turning off data augmentation is sufficient to remove the CPE.

The hypothesis has recently gained traction with both (Izmailov et al., 2021) and (Fortuin et al., 2021) pointing out that turning off data augmentation is sufficient to remove the CPE in (Wenzel et al., 2020)’s ResNet CIFAR10 setting (cf. Table 7 in (Izmailov et al., 2021) Appendix G and Figure A.11 in (Fortuin et al., 2021) Appendix A). On the other hand, (Wenzel et al., 2020)’s CNN-LSTM IMDB model already had a clean likelihood function and still gave rise to a CPE. From these observations we can already conclude that data augmentation is sufficient but not necessary for the CPE to arise. It remains an interesting open problem how to properly account for data augmentation in a Bayesian sense.

**Bad Prior Hypothesis:** Isotropic Gaussian priors are the de facto standard for modern Bayesian neural network inference (Fortuin et al., 2021). However, it is questionable whether such simplistic priors are optimal and whether they accurately reflect our true beliefs about the weight distributions. The bad prior hypothesis says that the CPE may only be an epi-phenomenon of a misspecified prior. The underlying argument is as follows: In classic Bayesian learning the number of parameters remains small and the prior is quickly dominated by the data. In contrast, in Bayesian deep learning the model dimensionality is typically on the same order if not larger than the dataset size (Kaplan et al., 2020). For such large models the prior will not be dominated by the data and will continue to exert an influence on the posterior. Hence the prior is critical.

The hypothesis recently got additional empirical support by (Fortuin et al., 2021), who found that the CPE can be partially alleviated by using heavy-tailed non-Gaussian priors.

## 2. Testing the Relative Influence of the Prior

A straightforward way to assess the bad prior hypothesis is to monitor the CPE while continuously trading off the relative influence between the prior term and the likelihood term: if the CPE becomes stronger as the relative influence of the prior increases, this would be an indication that the prior is poor. The relative weight of the prior vs. the likelihood in Bayesian inference is given by the dataset size  $n$ . To see this, recall the posterior energy function in Equation 1. Note how the log-likelihood is a sum over  $n$  datapoints, i.e. it scales with the dataset size  $n$ , while the log-prior is independent of  $n$ . This means that the prior will exert a stronger influence for smaller dataset sizes  $n$ .

In order to test the relative influence of the prior, we devise a simple experiment in which we train BNNs on random subsamples of different sizes, recording the optimal temperature for each value of the dataset size  $n$ . By varying the dataset size  $n$  we can test the following two hypotheses:

- If a bad prior causes the CPE, we would expect to see a stronger CPE for smaller dataset sizes  $n$ .
- From the theory of curation (Aitchison, 2020) we would expect that random subsamples of the dataset do not cause a change in the optimal temperature  $T^*$ , i.e. we would expect the same  $T^*$  regardless of  $n$ .

Note that care must be taken to ensure that the SG-MCMC inference has the same total number of parameter gradient updates across all data set sizes. In particular, fixing the batch size, we have to increase the cycle length (i.e. epochs per cycle) for smaller datasets. This ensures that the number of samples from the posterior and the overall number of gradient updates are the same across sample sizes.

## 3. CPE: A Symptom with Many Causes?

We use the SG-MCMC implementation of (Wenzel et al., 2020) for all our experiments, see Appendix F for a detailed description of the experimental setup. We generally report performance in terms of test cross-entropy. Shaded areas in the plots denote standard errors w.r.t. the number of random seeds (three in our case). We also define the *CPE-ratio* ‘‘CPE-ratio’’,

$$\text{CPE-ratio} = \ell_{T^*} / \ell_{T=1} \in (0, 1], \quad (4)$$

as the ratio between the cross-entropy loss at the optimal temperature  $T^*$  versus  $T = 1$  (Bayes posterior). A low CPE-ratio indicates that the performance of the tempered posterior is significantly better than the Bayes posterior. We perform experiments on SVHN (Netzer et al., 2011) and CIFAR-10 (Krizhevsky & Hinton, 2009). We use ResNet-20 neural networks (He et al., 2016) with Gaussian priors  $\mathcal{N}(0, 1)$ , unless stated otherwise. Further results can be found in the Appendix.

### Starting point: no CPE without data augmentation

We run SG-MCMC on SVHN and CIFAR-10 without data augmentation. As can be seen in Figure 1, we observe that SG-MCMC inference on the full dataset  $\mathcal{D}$  does not show any sign of a CPE, i.e.  $T = 1$  is close to optimal. From this we can conclude that *either SG-MCMC inference is accurate enough for this specific setting, or inaccurate inference is not necessary for the CPE to emerge*. We can also conclude that *the curation of SVHN and CIFAR10 does not give rise to a CPE*, which is somewhat surprising from (Aitchison, 2020)’s consensus theory standpoint.

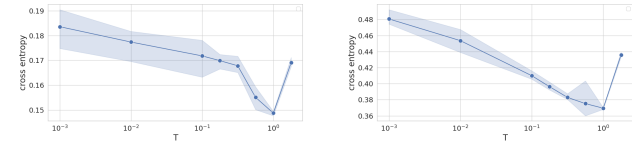
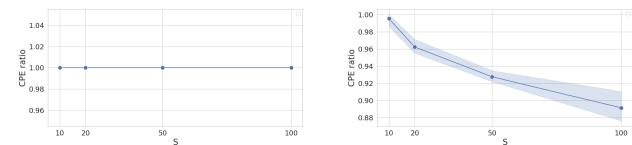


Figure 1: Test cross-entropy as a function of temperature  $T$  for (left) SVHN and (right) CIFAR-10. We can see that  $T = 1$  is optimal for SVHN and CIFAR-10.

### 3.1. CPE with curation, data augmentation, and random sub-sampling

**CPE and synthetic curation** We now test whether curation can cause the CPE in a simulated environment in which we control the number of labellers and can scale the amount of curation to large levels. We do so via synthetic curation by splitting the training set into a pre-training dataset  $\mathcal{D}_{pre}$  and a dataset  $\mathcal{D}_{tr}$ , by training a probabilistic classifier  $\hat{S}$  on  $\mathcal{D}_{pre}$  and by using  $S$  copies drawn from  $\hat{S}$  to independently re-label  $\mathcal{D}_{tr}$ , effectively simulating the behavior of  $S$  i.i.d. labellers. Using the labels induced by  $\hat{S}$ , we apply the curation procedure described in (Aitchison, 2020) to filter  $\mathcal{D}_{tr}$ , obtaining  $\mathcal{D}_{tr}^{cur}$ . Further details can be found in Section F.3.

Surprisingly, *we do not observe any cold posterior effect when only the training set is curated*, as shown in Figure 2a. However, as shown in Figure 2b, *curation of both the training and test set causes a CPE*.



(a) Only training set curated (b) Training & test set curated

Figure 2: CPE-ratio as a function of the number of labellers  $S$ .

**CPE and data augmentation** Recent works (Izmailov et al., 2021; Fortuin et al., 2021), have identified data augmentation as a cause for the CPE. We confirm that data augmentation can cause the CPE on SVHN and CIFAR-10,

as shown in Figure 3. On the other hand, the CPE can arise even without data augmentation. From these observations we can conclude that *data augmentation is sufficient but not necessary for the CPE to arise*. Note that data augmentation could hint at a problem with the prior (i.e. regularization).

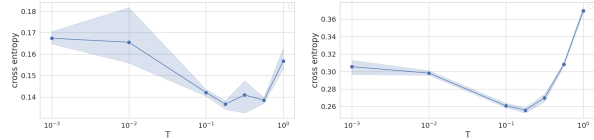


Figure 3: CPE on (left) SVHN and (right) CIFAR-10, both with data augmentation.

**CPE and random sub-sampling** We now discuss the random sub-sampling experiment to investigate the quality of the prior. We run SG-MCMC *without data augmentation* on subsets of SVHN with different sample sizes  $n$ . As explained, we ensure that the SG-MCMC inference has the same total number of parameter gradient updates across all data set sizes. The results in Figure 4 show that sub-sampling alone can cause CPE: the CPE is stronger for smaller  $n$ , as is evident from the CPE ratio in the plot on the (right). As the influence of the prior is larger on smaller datasets, *this is a strong indication that the prior is at fault*.

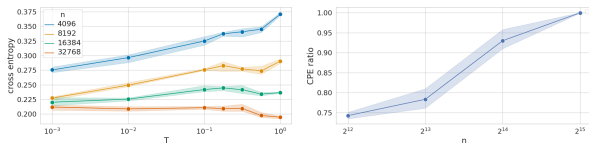


Figure 4: SVHN sub-sampling experiment. (left) Test cross-entropy for different temperatures. Each line represents a different subsample size. (right) CPER metric as a function of  $n$ . Note how the CPE is stronger for smaller  $n$ .

### 3.2. Comparing the relative influence of curation, data augmentation and sub-sampling

In Section F.4 we investigate the effect of sub-sampling on top of curation and data augmentation resp. the additional impact of curation and data augmentation on a sub-sampled dataset. The results are shown in Figure 7 in the Appendix.

### 3.3. Do standard Gaussian priors give too much weight to complicated hypotheses?

Here we investigate the influence of the prior in a synthetic toy dataset. To this end, we generate two clusters of datapoints from two 2D Gaussians with variance  $\sigma^2 = 1$  centered at  $(-1, -1)$  and  $(1, 1)$  respectively, in which the Bayes optimal classifier is given by the straight line  $y = -x$ .

To investigate the influence of the prior, we run full-batch MCMC for a 1 hidden layer MLP on subsets of the toy

dataset with varying sample sizes  $n$  and record the CPE across various temperatures  $T$ . The results, in Figure 5, show that the CPE becomes stronger for smaller datasets  $n$ .

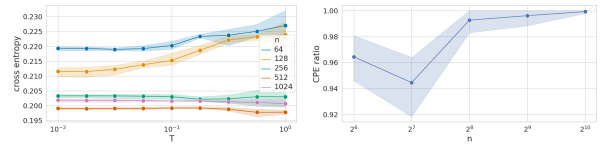


Figure 5: Testing the relative influence of the prior on a toy dataset. (left) Cross-entropy (CE) as a function of  $T$ , (right) CPE ratio as a function of dataset size  $n$ . We can see that the CPE becomes stronger for smaller dataset size  $n$ .

Interestingly, an analysis of the decision boundary, shown in Figure 6, suggests that the sharpened posterior obtained through tempered MCMC induces simpler functions than the Bayes posterior at  $T = 1$ .

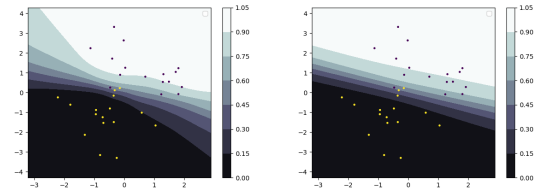


Figure 6: Decision boundary of 1 hidden layer MLP on toy dataset of size  $n = 32$ , (left)  $T = 1$ , (right)  $T = 10^{-2}$ . Color shading according to posterior mean prediction.

## 4. Discussion & Conclusion

Our results demonstrate how the CPE can arise in isolation from synthetic curation, data augmentation, and random sub-sampling. Specifically, we have confirmed that there is no CPE on SVHN and CIFAR-10 if data augmentation is turned off (Figure 1). On the other hand, the CPE can arise when both the training and test set are synthetically curated, i.e. when the role of the labellers is played by trained neural networks (Figure 2). Most importantly, we have shown that the CPE can also arise when a dataset that does not show any sign of CPE is randomly sub-sampled (Figure 4), providing a strong indication that the prior is at fault in the CPE. Our results suggest that priors *do* matter in Bayesian deep learning. We therefore consider it to be timely to study suitable priors for deep BNNs.

Since many of the recent deep learning advances, such as data augmentation, batch normalization and initialization distributions have been designed specifically for DNNs, it is not surprising that tempering, which gets BNNs closer to DNNs, improves their performance. The implication of this is that the Bayesian deep learning community has to find their own “advances” specifically tailored to BNNs.

## References

- Adlam, B., Snoek, J., and Smith, S. L. Cold posteriors and aleatoric uncertainty. *arXiv preprint arXiv:2008.00029*, 2020.
- Aitchison, L. A statistical theory of cold posteriors in deep neural networks. *arXiv preprint arXiv:2008.05912*, 2020.
- Barber, D. and Bishop, C. M. Ensemble learning for multi-layer networks. In *Advances in neural information processing systems*, pp. 395–401, 1998.
- Bhattacharya, A., Pati, D., Yang, Y., et al. Bayesian fractional posteriors. *The Annals of Statistics*, 47(1):39–66, 2019.
- Blundell, C., Cornebise, J., Kavukcuoglu, K., and Wierstra, D. Weight uncertainty in neural network. 37:1613–1622, 2015.
- Chen, C., Carlson, D., Gan, Z., Li, C., and Carin, L. Bridging the gap between stochastic gradient mcmc and stochastic optimization. In *Artificial Intelligence and Statistics*, pp. 1051–1060. PMLR, 2016.
- Chen, T., Fox, E., and Guestrin, C. Stochastic gradient Hamiltonian Monte Carlo. In *International Conference on Machine Learning*, pp. 1683–1691, 2014.
- Dayan, P., Hinton, G. E., Neal, R. M., and Zemel, R. S. The helmholtz machine. *Neural computation*, 7(5):889–904, 1995.
- Fortuin, V., Garriga-Alonso, A., Wenzel, F., Rätsch, G., Turner, R., van der Wilk, M., and Aitchison, L. Bayesian neural network priors revisited. *arXiv preprint arXiv:2102.06571*, 2021.
- Grünwald, P. Safe learning: bridging the gap between bayes, mdl and statistical learning theory via empirical convexity. In *Proceedings of the 24th Annual Conference on Learning Theory*, pp. 397–420. JMLR Workshop and Conference Proceedings, 2011.
- Grünwald, P. The safe bayesian. In *International Conference on Algorithmic Learning Theory*, pp. 169–183. Springer, 2012.
- Grünwald, P., Van Ommen, T., et al. Inconsistency of bayesian inference for misspecified linear models, and a proposal for repairing it. *Bayesian Analysis*, 12(4): 1069–1103, 2017.
- He, K., Zhang, X., Ren, S., and Sun, J. Deep residual learning for image recognition. In *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pp. 770–778, 2016.
- Hinton, G. and Van Camp, D. Keeping neural networks simple by minimizing the description length of the weights. In *in Proc. of the 6th Ann. ACM Conf. on Computational Learning Theory*, 1993.
- Izmailov, P., Vikram, S., Hoffman, M. D., and Wilson, A. G. What are bayesian neural network posteriors really like? *arXiv preprint arXiv:2104.14421*, 2021.
- Jansen, L. Robust Bayesian inference under model misspecification, 2013. Master thesis.
- Kaplan, J., McCandlish, S., Henighan, T., Brown, T. B., Chess, B., Child, R., Gray, S., Radford, A., Wu, J., and Amodei, D. Scaling laws for neural language models. *arXiv preprint arXiv:2001.08361*, 2020.
- Krizhevsky, A. and Hinton, G. Learning multiple layers of features from tiny images. 2009.
- Leimkuhler, B., Matthews, C., and Vlaar, T. Partitioned integrators for thermodynamic parameterization of neural networks. *arXiv preprint arXiv:1908.11843*, 2019.
- Li, C., Chen, C., Carlson, D., and Carin, L. Preconditioned stochastic gradient Langevin dynamics for deep neural networks. In *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- Ma, Y.-A., Chen, T., and Fox, E. A complete recipe for stochastic gradient MCMC. In *Advances in Neural Information Processing Systems*, pp. 2917–2925, 2015.
- MacKay, D. J. A practical bayesian framework for backpropagation networks. *Neural computation*, 4(3):448–472, 1992.
- MacKay, D. J. et al. Ensemble learning and evidence maximization. In *Proc. Nips*, volume 10, pp. 4083. Citeseer, 1995.
- Neal, R. M. *Bayesian learning for neural networks*. PhD thesis, University of Toronto, 1995.
- Netzer, Y., Wang, T., Coates, A., Bissacco, A., Wu, B., and Ng, A. Y. Reading digits in natural images with unsupervised feature learning. 2011.
- Ovadia, Y., Fertig, E., Ren, J., Nado, Z., Sculley, D., Nowozin, S., Dillon, J. V., Lakshminarayanan, B., and Snoek, J. Can you trust your model’s uncertainty? evaluating predictive uncertainty under dataset shift. In *Advances in Neural Information Processing Systems (NeurIPS 2019)*, 2019.
- Peterson, J. C., Battleday, R. M., Griffiths, T. L., and Rusakovsky, O. Human uncertainty makes classification more robust. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9617–9626, 2019.

- Quionero-Candela, J., Sugiyama, M., Schwaighofer, A., and Lawrence, N. D. *Dataset shift in machine learning*. The MIT Press, 2009.
- Sutskever, I., Martens, J., Dahl, G., and Hinton, G. On the importance of initialization and momentum in deep learning. In *International conference on machine learning*, pp. 1139–1147, 2013.
- Welling, M. and Teh, Y. W. Bayesian learning via stochastic gradient Langevin dynamics. In *Proceedings of the 28th international conference on machine learning (ICML-11)*, pp. 681–688, 2011.
- Wenzel, F., Roth, K., Veeling, B., Swiatkowski, J., Tran, L., Mandt, S., Snoek, J., Salimans, T., Jenatton, R., and Nowozin, S. How good is the bayes posterior in deep neural networks really? In *International Conference on Machine Learning*, pp. 10248–10259. PMLR, 2020.
- Wilk, M. v. d., Bauer, M., John, S., and Hensman, J. Learning invariances using the marginal likelihood. In *Proceedings of the 32nd International Conference on Neural Information Processing Systems*, pp. 9960–9970, 2018.
- Wilson, A. G. and Izmailov, P. Bayesian deep learning and a probabilistic perspective of generalization. *arXiv preprint arXiv:2002.08791*, 2020.
- Zeno, C., Golan, I., Pakman, A., and Soudry, D. Why cold posteriors? on the suboptimal generalization of optimal bayes estimates. *3rd Symposium on Advances in Approximate Bayesian Inference*, 2020. URL <https://openreview.net/pdf/7b7b21325b3db402632a8da3f607476b328671d7.pdf>.
- Zhang, R., Li, C., Zhang, J., Chen, C., and Wilson, A. G. Cyclical stochastic gradient MCMC for Bayesian deep learning. In *International Conference on Learning Representations (ICLR 2020)*, 2020.

## A. Introduction to Bayesian Deep Learning

Deep neural networks have achieved great success in predictive accuracy for supervised learning tasks. Unfortunately, however, they still fall short in giving useful estimates of their predictive *uncertainty*, i.e. meaningful confidence values for how certain the model is with its predictions (Ovadia et al., 2019). Quantifying uncertainty is especially crucial in real-world settings, which often involve data distributions that are shifted from the one seen during training (Quionero-Candela et al., 2009).

Bayesian deep learning combines deep learning with Bayesian probability theory. Instead of optimizing a single network, *Bayesian neural networks* (BNNs) learn a distribution over model parameters or equivalently sample an ensemble of likely models given the data, promising better generalization performance (no overfitting) and principled uncertainty quantification (robust predictions) (Neal, 1995; MacKay, 1992; Dayan et al., 1995).

In Bayesian deep learning we either learn a distribution  $q(\boldsymbol{\theta})$  over models compatible with the data, i.e.  $q(\boldsymbol{\theta}) \simeq p(\boldsymbol{\theta}|\mathcal{D})$ , or we *sample* an ensemble of models  $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K \sim p(\boldsymbol{\theta}|\mathcal{D})$  from the posterior over likely models

$$p(\boldsymbol{\theta}|\mathcal{D}) \propto p(\mathcal{D}|\boldsymbol{\theta}) p(\boldsymbol{\theta}), \quad (5)$$

where  $p(\mathcal{D}|\boldsymbol{\theta}) = \prod_{i=1}^n p(y_i|\mathbf{x}_i, \boldsymbol{\theta})$  is the likelihood, relating the model we want to learn to the observations  $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^n$ , and  $p(\boldsymbol{\theta})$  is a proper prior, e.g. a Gaussian density.

BNN predictions involve model averaging: rather than betting everything on a single model, we predict on a new instance  $\mathbf{x}$  by *averaging* over all likely models compatible with the data,

$$p(y|\mathbf{x}, \mathcal{D}) = \int p(y|\mathbf{x}, \boldsymbol{\theta}) p(\boldsymbol{\theta}|\mathcal{D}) d\boldsymbol{\theta} \simeq \sum_{k=1}^K p(y|\mathbf{x}, \boldsymbol{\theta}_k), \text{ where } \boldsymbol{\theta}_k \sim p(\boldsymbol{\theta}|\mathcal{D}). \quad (6)$$

Equation 6 is also known as the *posterior predictive* or *Bayesian model average*. Note that, in practice, solving the integral exactly is impossible. However, we can approximate it via Monte Carlo sampling using an ensemble of models  $\boldsymbol{\theta}_k \sim p(\boldsymbol{\theta}|\mathcal{D})$ .

The two main inference paradigms to learn a distribution over model parameters  $q(\boldsymbol{\theta}) \simeq p(\boldsymbol{\theta}|\mathcal{D})$  respectively to sample from the posterior  $\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_K \sim p(\boldsymbol{\theta}|\mathcal{D})$  are Variational Bayes (VB) (Hinton & Van Camp, 1993; MacKay et al., 1995; Barber & Bishop, 1998; Blundell et al., 2015) and Markov Chain Monte Carlo (MCMC) (Neal, 1995; Welling & Teh, 2011; Chen et al., 2014; Ma et al., 2015). We will focus on MCMC methods as they are simple to implement and can be scaled to large models and datasets when used with stochastic minibatch gradients (SG-MCMC) (Welling & Teh, 2011; Chen et al., 2014; Li et al., 2016).

In recent years, the Bayesian deep learning community has developed increasingly accurate and efficient approximate inference procedures for deep BNNs (cf. references in the paragraph above). Despite this algorithmic progress, however, important questions surrounding BNNs remain unanswered to this day. A recent and particularly prominent one concerns the “*cold posterior effect*” (CPE), which describes the disturbing observation that the predictive performance of BNNs can be significantly improved if the Bayes posterior is artificially sharpened  $p(\boldsymbol{\theta}|\mathcal{D})^{1/T}$  using a temperature parameter  $T < 1$  (Wenzel et al., 2020). Such cold posteriors sharply deviate from the Bayesian paradigm but are commonly used as heuristics in practice, see Section 2.3 in (Wenzel et al., 2020).

The CPE is problematic in theory and practice, and since the effect was identified many researchers have proposed hypotheses to explain the phenomenon. There has been an ongoing debate questioning the roles of isotropic Gaussian priors (Wenzel et al., 2020; Zeno et al., 2020; Fortuin et al., 2021), the likelihood model (Aitchison, 2020), inaccurate inference (Izmailov et al., 2021; Wenzel et al., 2020), and data augmentation (Izmailov et al., 2021; Fortuin et al., 2021). However, despite this intensive research effort the effect remains poorly understood.

## B. Cold Posteriors: Background & Related Work

**Tempering the Posterior vs. Tempering the Likelihood.** Note that, besides tempering the posterior as in Equation 1, one can also temper only the likelihood resp. scale only the log-likelihood, as is commonly done in VB, see Section 2.3 in (Wenzel et al., 2020). It is worth pointing out though that both variants are practically equivalent if the prior variance is multiplicative in log-prior pre-factors (such as for Gaussian priors) and if sufficiently many prior variances are grid-searched over as part of the inference pipeline, such that for the best performing posterior-tempered model there is a corresponding likelihood-tempered model with variance scaled by  $1/T$  in the grid and vice versa.

**Inaccurate Inference Hypothesis:** While recent works on the CPE have taken great care to ensure that their SG-MCMC based inference procedure yields as accurate an approximation to the true posterior as possible, it remains difficult to definitively assess the approximation accuracy without having access to the (intractable) true posterior. To rule out obvious problems with the inference mechanism, (Wenzel et al., 2020) proposed a set of diagnostics, based on comparing ensemble statistics to their theoretically known values. We have closely monitored these diagnostics in our experiments, however, despite our own extensive efforts to ensure accurate inference, we cannot exclude the possibility that our inference may be inaccurate even though the diagnostics match.

To investigate the inference hypothesis further, together with other foundational questions in Bayesian deep learning, (Izmailov et al., 2021) recently applied full-batch Hamiltonian Monte Carlo (HMC)<sup>3</sup>, which is considered to be the gold-standard in terms of inference accuracy, to the models of (Wenzel et al., 2020), showing that with HMC inference and when data augmentation is turned off the CPE disappears. While these results can easily be misconstrued as evidence that the approximate inference is at fault in the CPE, there are good reasons to think otherwise.

For instance, (Izmailov et al., 2021) show, using the code of (Wenzel et al., 2020), that turning off data augmentation alone is sufficient to remove the CPE (cf. Table 7 in (Izmailov et al., 2021) Appendix G). We too can confirm that the CPE does not arise with SG-MCMC based inference applied to (Wenzel et al., 2020)’s models when data augmentation is turned off, cf. Figure 1 in Section 3. From this we can already conclude that either SG-MCMC inference is accurate enough for this specific setting, or inaccurate inference is not necessary for the CPE to emerge<sup>4</sup>.

A more direct counter argument follows from (Adlam et al., 2020), which have demonstrated that there can be a CPE in Gaussian Processes (GP) regression, where the posterior has a closed form solution, provided that the aleatoric uncertainty is overestimated. Hence, the CPE can arise in a setting where exact inference is possible. They also provide experimental evidence that the CPE can arise in classification tasks in the infinite-width neural network Gaussian process (NNGP) limit. Thus, while inaccurate inference may be sufficient, it does not appear to be necessary for the CPE to emerge.

Note that a similar observation was made in (Grünwald, 2012; Grünwald et al., 2017), which demonstrated benefits of tempering, with  $T > 1$  in their setting, in the context of exact inference.

**Curation Hypothesis (*model misspecification*):** (Aitchison, 2020) devises a theory that attributes the effectiveness of tempering to the fact that standard benchmark datasets such as CIFAR-10 are carefully curated and that we should take this curation into account by tempering BNNs. In (Aitchison, 2020)’s model of curation (actual dataset curation may differ), a datapoint  $\mathbf{x}$  is added to the dataset if and only if *all*  $S$  labellers independently agree on the label  $y_s$  to be assigned to  $\mathbf{x}$ , while the datapoint is discarded if at least one pair of labellers  $s, s'$  disagrees  $y_s \neq y_{s'}$ . The main argument put forward by (Aitchison, 2020) is that if we a priori know that the dataset is curated in the sense described above, we should take this into account in our likelihood model. The proposed likelihood to be used in the case of curation should then be of the following form

$$p(y, \mathbf{x}|\boldsymbol{\theta}) \propto p(\{y_s = y\}_{s=1}^S | \mathbf{x}, \boldsymbol{\theta}) = \prod_s p(y_s = y | \mathbf{x}, \boldsymbol{\theta}) . \quad (7)$$

Thus, assuming that the labellers are i.i.d., the probability of consensus on label  $y$  is

$$p(y, \mathbf{x}|\boldsymbol{\theta}) \propto p(y|\mathbf{x}, \boldsymbol{\theta})^S , \quad (8)$$

which corresponds to a cold posterior where only the log-likelihood is re-scaled while the prior is not (cf. discussion at the beginning of this Section). (Aitchison, 2020) argues that we should observe  $S \approx 1/T$ , i.e., the optimal temperature should roughly<sup>5</sup> be inversely proportional to the total number of labellers involved in the curation of a datapoint.

Note that (Aitchison, 2020)’s model of curation only includes a data point if *all* labellers agree, but in practice, datasets are often collected with some tolerance of labeller disagreement, e.g. in that a datapoint is included if a certain fraction of labellers agree. However, a more realistic (weaker) model of curation, that filters out fewer “hard” instances from the pool of uncurated datapoints, would only give rise to a weaker, less pronounced CPE: if we do not observe the CPE for (Aitchison, 2020)’s simplistic model of curation, which is the most extreme form of curation imaginable, we would not

<sup>3</sup>(Izmailov et al., 2021) parallelized the HMC computation over *hundreds* of Tensor Processing Units (TPUs)

<sup>4</sup>We can also conclude that data augmentation is sufficient for the CPE to arise (although it is not necessary, as we will see later).

<sup>5</sup>To obtain an exact correspondence, we would need access to the discarded datapoints in order to be able to marginalize them out, cf. “marginalise over unknown latents” in Section 3 in (Aitchison, 2020) for how to account for the discarded images in a proper Bayesian sense

expect to observe it under more realistic models of curation either. See Section D in the Appendix for additional details on how popular datasets like CIFAR-10 or SVHN were collected.

Finally, we would like to point out that (Adlam et al., 2020) made a similar, albeit somewhat more general argument regarding curation resp. mismatch of aleatoric uncertainty. Specifically, they show that the CPE can arise if the model overestimates the aleatoric uncertainty, which is naturally reduced when the dataset is curated.

**Data Augmentation Hypothesis (*model misspecification*):** Current deep learning practices use a number of techniques, including data augmentation, that technically do not obey the likelihood principle (see Appendix K in (Wenzel et al., 2020) for an in-depth discussion of so-called “dirty likelihoods”). The data augmentation hypothesis specifically says that the CPE is largely an artifact of using data augmentation and that turning off data augmentation is sufficient to remove the CPE.

The hypothesis has recently gained traction with both (Izmailov et al., 2021) and (Fortuin et al., 2021) pointing out that turning off data augmentation is sufficient to remove the CPE in (Wenzel et al., 2020)’s ResNet CIFAR10 setting (cf. Table 7 in (Izmailov et al., 2021) Appendix G and Figure A.11 in (Fortuin et al., 2021) Appendix A). On the other hand, (Wenzel et al., 2020)’s CNN-LSTM IMDB model already had a clean likelihood function and still gave rise to a CPE. From these observations we can already conclude that data augmentation is sufficient but not necessary for the CPE to arise.

As the performance of deep neural networks is often significantly better when using some form of data augmentation, it is not really an option to just turn it off in BNNs, while properly accounting for it in Bayesian inference does not seem trivial either. On the one hand, data augmentation affects the data points that enter the likelihood function. However, while data augmentation may increase the amount of data seen by the model, that increase is certainly not equal to the number of times each data point is augmented (after all, augmented data is not independent from the original data). On the other hand, considering data augmentation as a form of regularization (constraining the classification functions to be invariant to certain transformations), one can argue that it should be represented in the prior (Wilk et al., 2018). It remains an interesting open problem how to properly account for data augmentation in a Bayesian sense.

**Bad Prior Hypothesis:** Isotropic Gaussian priors are the de facto standard for modern Bayesian neural network inference (Fortuin et al., 2021). However, it is questionable whether such simplistic priors are optimal and whether they accurately reflect our true beliefs about the weight distributions. The bad prior hypothesis says that the CPE may only be an epiphenomenon of a misspecified prior. The underlying argument is as follows: In classic Bayesian learning the number of parameters remains small and the prior is quickly dominated by the data. In contrast, in Bayesian deep learning the model dimensionality is typically on the same order if not larger than the dataset size (Kaplan et al., 2020). For such large models the prior will not be dominated by the data and will continue to exert an influence on the posterior. Hence the prior is critical.

The hypothesis has already been put forward by (Wenzel et al., 2020), who reported that the CPE becomes stronger with increasing model dimensionality. It recently got additional empirical support by (Fortuin et al., 2021), who found that the CPE can be partially alleviated by using heavy-tailed non-Gaussian priors. More specifically, they find that for fully connected neural networks (FCNNs), heavy-tailed priors can both improve predictive performance and alleviate the CPE. For convolutional neural networks (CNNs), the CPE can also be removed with heavy-tailed priors, however, the resulting performance gains are less striking. On the other hand, the performance of CNNs can be improved with correlated priors, although they no longer appear to alleviate the CPE.

Finally, we note that the prior that ultimately matters is the prior over functions that is induced when a prior over parameters is combined with the functional form of a neural network architecture (Wilson & Izmailov, 2020; Izmailov et al., 2021). Still, this does not render the prior over parameters irrelevant, as innocent-looking priors may inadvertently be highly informative, for instance placing large prior mass on undesirable functions.

## C. Further Related Work

**“Warm” Posteriors:** Motivated by the behavior of Bayesian inference in *misspecified* models (Grünwald et al., 2017; Jansen, 2013) extensively studied the so called “generalized” Bayesian inference, i.e, the Bayes posterior in which only the likelihood is tempered. In particular, the “Safe Bayes” framework (Grünwald, 2012; 2011; Grünwald et al., 2017) was developed to tune the temperature parameter. However, these works consider only “warm posteriors”  $T > 1$  (the inverse temperature is called “learning rate” in the relevant literature), as a way to learn under model misspecification. Warm posteriors can arise in a context where the model is misspecified, for instance by assuming homoscedastic noise where the

data-generating noise is heteroscedastic. Under this misspecification, the model overfits the datapoints, despite the fact that the prior is well-specified (for instance in their case it is centered around the best performing and non-overfitting solution). We hypothesize that in (Grünwald et al., 2017) the prior favours simple models, hence it is beneficial to put more weight onto the prior and use warm posterior. The opposite might happen in BNNs: cold posteriors counteract the effect of a bad prior that tends to prefer overcomplicated solutions. Finally, we mention the work of (Bhattacharya et al., 2019), in which the authors develop *fractional posteriors* with the goal of decreasing posterior concentration.

## D. Dataset Collection & Curation

Here we review the way that the two datasets that are mainly used in our experiments, SVHN and CIFAR-10, have been collected and curated.

**SVHN** The Street View House Numbers dataset (Netzer et al., 2011), which is divided into a training corpus  $\mathcal{D}$  of around 73257 training images and a test set  $\mathcal{D}_{te}$  of around 26k images. Although we do not know the exact number of labellers, the dataset has undergone a curation procedure in the sense of (Aitchison, 2020). In particular, AMT was adopted<sup>6</sup> (quoting from (Netzer et al., 2011), "The SVHN dataset was obtained from a large number of Street View images using a combination of automated algorithms and the Amazon Mechanical Turk (AMT) framework").

**CIFAR-10** In CIFAR-10 (Krizhevsky & Hinton, 2009), labellers followed strict guidelines to ensure high quality labelling of the images. In particular, labellers were instructed that "it's worse to include one that shouldn't be included than to exclude one. False positives are worse than false negatives", and "If there is more than one object that is roughly equally prominent, reject". The reader is invited to review Appendix C in (Krizhevsky & Hinton, 2009) .

Unfortunately, in the relevant papers there are no details on the specific curation process that was applied, e.g. the number of labellers per image, or whether *all* labellers have to agree on a label or only a subset of them.

## E. MCMC Inference

In this section we review the basics of (SG)-MCMC inference. The description of the implementation and adaptations for deep learning can be found in Section F.2 below.

The two main inference paradigms to learn a distribution over model parameters  $q(\theta) \simeq p(\theta|\mathcal{D})$  respectively to sample from the posterior  $\theta_1, \dots, \theta_K \sim p(\theta|\mathcal{D})$  are Variational Bayes (VB) (Hinton & Van Camp, 1993; MacKay et al., 1995; Barber & Bishop, 1998; Blundell et al., 2015) and Markov Chain Monte Carlo (MCMC) (Neal, 1995; Welling & Teh, 2011; Chen et al., 2014; Ma et al., 2015; Li et al., 2016). We will focus on MCMC methods as they are simple to implement and can be scaled to large models and datasets when used with stochastic minibatch gradients (SG-MCMC) (Welling & Teh, 2011; Chen et al., 2014).

Markov Chain Monte Carlo (MCMC) methods allow to sample an ensemble of models  $\theta_1, \dots, \theta_K \sim p(\theta|\mathcal{D})^{1/T}$  from the (tempered) posterior  $p(\theta|\mathcal{D})^{1/T}$ , by performing a guided random walk in parameter space in which artificial noise is injected into the updates  $\theta_k \rightarrow \theta_{k+1}$  in such a way that the ensemble distribution converges to the desired posterior  $p(\theta|\mathcal{D})^{1/T}$  (in the limit of small step sizes and long enough run times) (Neal, 1995; Welling & Teh, 2011; Ma et al., 2015; Li et al., 2016). By injecting artificial noise, the algorithm explores the loss landscape instead of approaching a single point estimate  $\hat{\theta}$ .

**SG-MCMC** Recent advances in stochastic inference through Markov Chain Monte Carlo (MCMC) methods have made the task of sampling from the posterior distribution of deep neural networks more efficient (Welling & Teh, 2011; Zhang et al., 2020; Wenzel et al., 2020). In particular, the usage of mini-batches gave rise to stochastic gradient MCMC methods (SG-MCMC) (Welling & Teh, 2011), which is further improved through various techniques such as momentum variables (Chen et al., 2014), preconditioning (Li et al., 2016), and cyclical stepsize (Zhang et al., 2020). All these methods perform stochastic updates in parameter space that come from the discretization of a stochastic process (Ma et al., 2015). For the purpose of exposition, here we mention SG-MCMC in its simplest form, given by SGLD (stochastic gradient Langevin

<sup>6</sup><https://www.mturk.com/>

Table 1: Deep learning features and SG-MCMC hyper-parameters. The double horizontal line splits the experiments of the paper from those in the appendix.

Experiment	Hyper-parameter						
	data augm	batch norm	precond.	learning rate	burn-in	cycle length	epochs
SVHN & CIFAR-10 (Fig. 1)	✗	✓	✗	0.1	150	50	1500
SVHN (Fig.2)	✓	✓	✗	0.1	100	25	500
SVHN & CIFAR-10 (Fig. 3)	✓	✓	✗	0.1	150	50	1500
SVHN (Fig. 4)	✗	✓	✗	0.1	100	25	500
Toy data (Fig. 5)	✗	✗	✗	0.1	500	75	2000
SVHN & CIFAR-10 (Fig. 11)	✗	✓	✗	0.1	100	25	500
MNIST (Fig. 9)	✗	✗	✗	0.1	150	50	1500

dynamics), in which the updates have the form

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \frac{\epsilon_t}{2} \left( \frac{n}{B} \sum_{i=1}^B \nabla \log p(y|\mathbf{x}_i, \boldsymbol{\theta}) + \nabla \log p(\boldsymbol{\theta}) \right) + \boldsymbol{\eta}_t, \quad (9)$$

where  $B$  indicates the size of a mini-batch, and  $\boldsymbol{\eta}_t \sim \mathcal{N}(0, \epsilon_t I)$ . We define the gradient of the mini-match posterior energy as  $\nabla \tilde{U}(\boldsymbol{\theta}) := \frac{n}{B} \sum_{i=1}^B \nabla \log p(y|\mathbf{x}_i, \boldsymbol{\theta}) + \nabla \log p(\boldsymbol{\theta})$ . (Welling & Teh, 2011) show that if  $\epsilon_t$  is such that:

$$\sum_{i=1}^{\infty} \epsilon_t = \infty \quad \sum_{i=1}^{\infty} \epsilon_t^2 < \infty, \quad (10)$$

then convergence to a local maximum is guaranteed. We will use the SG-MCMC implementation of (Wenzel et al., 2020) throughout our experiments<sup>7</sup>, which combines the aforementioned techniques, further discussed in Section F.2 below.

## F. Experimental Setup

The experimental details, including the SG-MCMC hyperparameters, are included in Table 1. Note that for the subsampling experiments on SVHN (Figure 4), the table entries in the last three columns, that have the number of epochs as units (i.e. burn-in period, cycle length, epochs), refer to the *full dataset size*. When subsampling is applied, the number of epochs are adjusted such that the number of gradient steps is kept fixed. For instance, if half of the dataset is used, the number of epochs, cycle length epochs and burn-in epochs doubles.

### F.1. Neural Network Architectures

For the SG-MCMC experiments, we use a 20-layer architecture with residual layers (He et al., 2016) and batch normalization. For the SG-MCMC experiments on the toy dataset, we use a single hidden layer fully connected net with 20 units and ReLU activation function. The SG-MCMC method that we adopt is the one in (Wenzel et al., 2020) and summarized in Sections F.2 and E. In particular, no preconditioning is used.

<sup>7</sup>[https://github.com/google-research/google-research/tree/master/cold\\_posterior\\_bnn](https://github.com/google-research/google-research/tree/master/cold_posterior_bnn)

## F.2. Inference Method / Training Procedure

In this work, we will mainly use the inference method proposed in (Wenzel et al., 2020), which adapts recent advances in optimization for deep learning and stochastic inference to SG-MCMC. See also (Chen et al., 2016) for some interesting connections between SG-MCMC and stochastic optimization.

**Momentum Variables** Adding momentum to SGD is an optimization technique to accelerate gradient based optimization methods that is widely used in deep learning (Sutskever et al., 2013). Momentum variables were added to SG-MCMC methods in (Chen et al., 2014), giving raise to the stochastic-gradient version of Hamiltonian dynamics (SG-HMC). SGLD can be modified as follows to include them:

$$\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t + \epsilon_t \mathbf{m}_{t+1} \tag{11}$$

$$\mathbf{m}_{t+1} = (1 - \epsilon_t \alpha) \mathbf{m}_t - \epsilon_t \nabla \tilde{U}(\boldsymbol{\theta}_t) + \sqrt{2} \boldsymbol{\eta}_t, \tag{12}$$

where  $\alpha$  is the momentum weight.

**Layerwise Preconditioning** A subset of our experiments were performed both with and without preconditioning, which did not make a big difference. The reported results are without preconditioning.

**Cyclical step size** Cyclical step size was introduced by (Zhang et al., 2020) to guarantee better exploration, given the fact that posterior exploration is somewhat limited in standard SGLD due to the fact that the learning rate  $\epsilon_t$  must be small enough to avoid bias in estimation and MH acceptance/rejection steps. It consists in alternating updates with large learning rate, which allows to overshoot the local minima and therefore having a better *exploration* of the posterior landscape, and updates with very small learning rate, during which samples from the posterior are collected. All the details of the algorithm can be found in (Wenzel et al., 2020), Section 3.

## F.3. Training the synthetic labellers

For the curation experiment, the role of the labeller is played by a neural network.

The following experiment is designed to test whether curation can cause the CPE in a simulated environment in which we control the number of labellers and can scale the amount of curation to large levels. We do so by performing synthetic curation as follows:

- We first split  $\mathcal{D}$  into two non-intersecting sets: a pre-training dataset  $\mathcal{D}_{pre}$  and a dataset  $\mathcal{D}_{tr}$ .
- We train a probabilistic classifier  $\hat{\mathcal{S}}$  on  $\mathcal{D}_{pre}$ , that learns a categorical distribution over the labels (given by the output of the softmax activation of a neural network).
- We use  $S$  copies drawn from  $\hat{\mathcal{S}}$  to independently re-label  $\mathcal{D}_{tr}$ , effectively simulating the behavior of  $S$  i.i.d. labellers. The labels are obtained by sampling from the categorical distribution learned by the network. The more uncertain the model is about a prediction for some input, the more disagreement between labellers we expect for that input.
- Using the labels induced by  $\hat{\mathcal{S}}$ , we apply the curation procedure described in (Aitchison, 2020) and summarized earlier in Section 1, to filter  $\mathcal{D}_{tr}$  and obtain  $\mathcal{D}_{tr}^{cur}$ . Note that the consensus label does not necessarily match with the original label (which can happen for images where the model is confident on the wrong label).

As the labeller classifier  $\hat{\mathcal{S}}$ , we train an 8-layer ResNet on  $\mathcal{D}_{pre}$  with the Adam optimizer.

We perform SG-MCMC inference on the curated dataset  $\mathcal{D}_{tr}^{cur}$  using  $S$  labellers, for various values of  $S$ . The cross-entropy is evaluated both on the original test set  $\mathcal{D}_{test}$  - which is simply re-labeled according to the trained model  $\hat{\mathcal{S}}$  - and the curated one  $\mathcal{D}_{test}^{cur}$  using the same number of  $S$  labellers as in the training set.

## F.4. Relative influence of curation, data augmentation and subsampling

In the experiments Section, we have seen that sub-sampling alone can cause the CPE to arise. Here we investigate the effect of sub-sampling on top of curation and data augmentation resp. the additional impact of curation and data augmentation on

a sub-sampled dataset. The results are shown in Figure 7. The plot can be read in two ways: either one looks at a fixed dataset size  $n$  and compares the impact of curation (blue) and data augmentation (green) over the original test set (orange), or one looks at the relative change in CPER induced by sub-sampling for a given curve (e.g. how much each curve drops when going from  $n = 16384$  to  $n = 8192$ ).

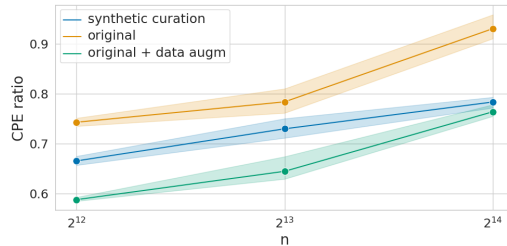


Figure 7: Comparing the CPER for random sub-sampling with curation and data augmentation on datasets of size 4096, 8192 and 16384.

Here we give details regarding the experimental setup. If one inspects the core of curation, one may argue that as it consists of *removing* datapoints, the CPE is caused by using smaller datasets rather than the curation procedure itself. Therefore, we want to compare the added contribution to the CPE of curation (and data augmentation) with respect to random sub-sampling. While adding data augmentation to random sub-sampling is straightforward, in the case of curation it is not practical to find the exact number of labellers  $S$  to match a desired dataset size. Instead, we proceed as follows: we choose the number of labellers such that the curated dataset contains around 17k-18k datapoints (the exact number may vary depending on the random seed).<sup>8</sup> Then we perform random sub-sampling to match the desired dataset size (16384, 8192, 4096 in our case) on the synthetically curated dataset. The only difference with the previous experiments is that we use the original labels, instead of the consensus labels. In this way, if one image  $x$  belongs to both the synthetically curated dataset and the randomly sampled one, it has the same label  $y$ . We run SG-MCMC with the same parameters settings as in Section 3, paragraph "CPE and random sub-sampling".

## G. Further Experimental Results

### G.1. Sharpening the Labeller’s distribution

What might cause the cold posterior effect to arise in the curation theory is the removal of ambiguous datapoints, and this does not have to necessarily correlate with the number of labellers. To test this, we use the same labelling scheme as in the previous section; the only difference is that this time we artificially sharpen (or flatten) the labelling distribution by raising the probabilities to the power of  $\alpha$  and then re-normalizing. If  $\alpha > 1$  we are sharpening the distribution, and  $\alpha < 1$  we are flattening it. Results are shown in Figure 8 in a summary plot. Note that in Figure 8  $\alpha$  correlates with an increase in optimal temperature: for a higher values of  $\alpha$ , the labeller is more confident, therefore the resulting dataset is *less* curated and so it follows that there is a weaker cold posterior effect. This confirms that even if the curation argument is correct, it is hard to find an exact correspondence between the number of labellers and optimal inverse temperature.

### G.2. Subsampling on MNIST

We repeat the subsampling experiment on MNIST. Results are shown in Figure 9. Note how for small sample sizes the CPER metric decreases significantly, indicating the presence of the cold posterior effect.

### G.3. CIFAR-10 and SVHN

For the subsampling experiment on CIFAR-10, see Figure 10. For other experiments on full SVHN and CIFAR-10, in which we use less number of epochs per cycle, see Figure 11.

<sup>8</sup>As we observe that the classifier  $\tilde{S}$  is over confident on many datapoints, we decide to smooth its output distribution, keeping the rank of the probabilities unchanged. See Appendix G.1

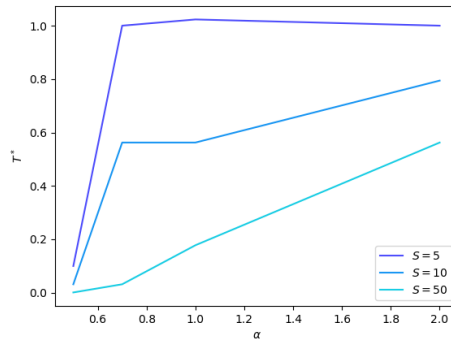


Figure 8: Optimal temperature as a function of  $\alpha$ . Lighter color indicates more labellers, 5, 10, 50, respectively. Note that the optimal temperature correlates with the confidence in prediction of the labeller.

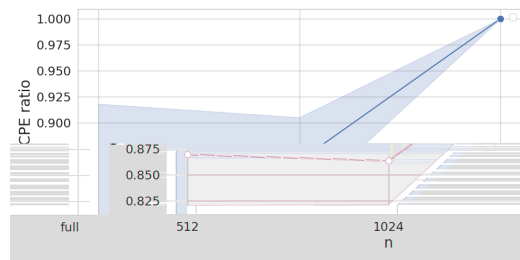


Figure 9: MNIST subsampling experiment

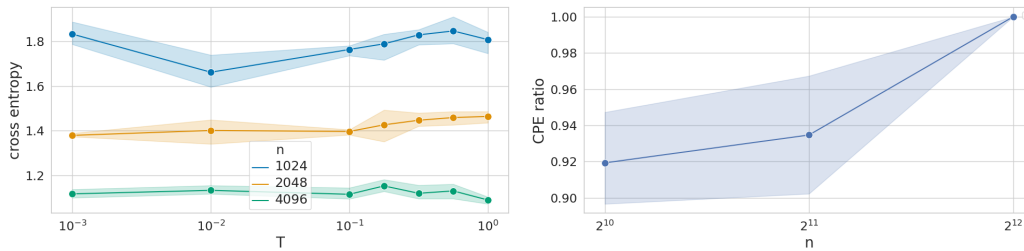


Figure 10: CIFAR subsampling experiment: we can see that even for low dataset sizes there is (almost) no CPE on CIFAR-10

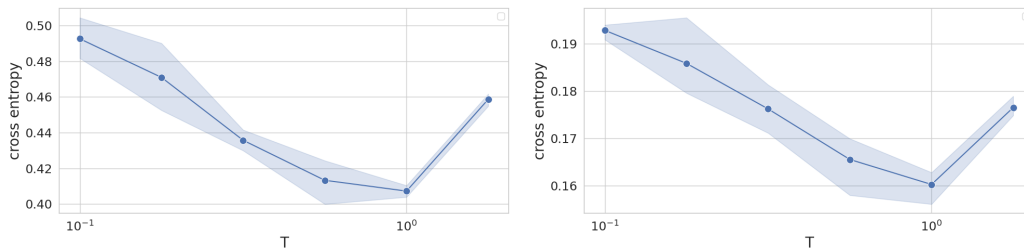


Figure 11: CIFAR experiment (left) and SVHN (right) on full dataset with smaller cycles and less samples (no data augmentation)

#### G.4. SVHN curated

In Figure 12 and 13, we show the plots underlying the curation experiment on SVHN, summarized in Figure 2. In Figure 14 we show the underlying plots of the curation + subsampling experiment of Figure 7.

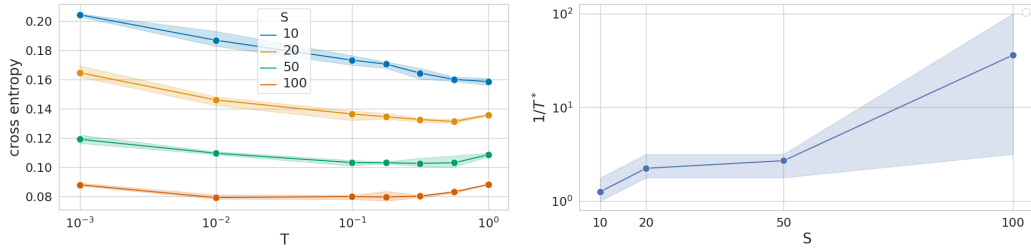


Figure 12: Plots underlying experiments when both train and test set are curated

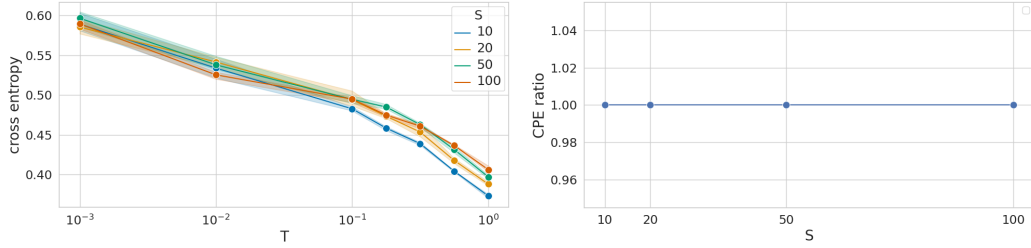


Figure 13: Plots underlying experiments when only the training set is curated: they all perform quite similarly.

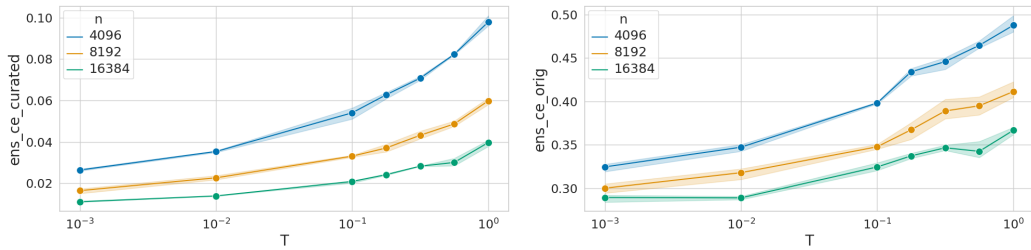


Figure 14: Curation + subsample: underlying plots. Left: curated test set. Right: original test set

### G.5. Data Augmentation

In Figure 15, we show the plots underlying Figure 7 for the data augmentation part.

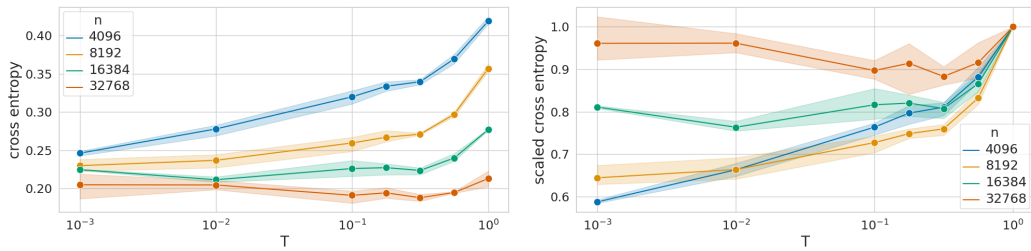


Figure 15: Data augmentation and sub-sampling. The right plot is the scaled version of the left plot.

### G.6. Experiments on CIFAR-10H

**Summary:** We repeat some of the experiment on CIFAR-10H by (Aitchison, 2020). Our results show that over-weighting the likelihood with respect to the prior is sufficient for the removal of the CPE. Furthermore, this weight does not equal the number of labellers, suggesting that the curation theory of (Aitchison, 2020) is not accurate enough to explain the CPE alone (as was already acknowledged by the author).

In (Aitchison, 2020), the authors state that the cold posterior effect might be due to the labelling process, based on consensus: each image is added to the dataset only if all  $S$  labellers agree on one class. In particular, they say that if we had access to the  $S$  original labels for each image, then we should not observe any cold posterior effect. They devise an experiment on CIFAR-10H (Peterson et al., 2019), in which all the (approximately 50) labels are given for each image. Our experiments show that over-weighting the likelihood with respect to the prior is sufficient to eliminate the CPE, and that the weight does not correspond to the number of labellers.

More formally, the dataset can be defined as  $\mathcal{D} := \{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^n$ , where  $\mathbf{y}_i$  is a vector containing the counts for each class, i.e. its  $j$ -th element is the number of labellers that chose class  $j$ . Given a mini-batch of size  $B$ , the log-likelihood used in an SG-MCMC method, is the following:

$$\mathcal{L}_c := \frac{n}{B} \sum_{i=1}^B \mathbf{y}_i^T \log f_i. \tag{13}$$

This is very similar to the use of label smoothing in which the label vector is  $\mathbf{y}_{ls} := \frac{1}{S} \mathbf{y}$ . The corresponding mini-batch likelihood has the form:

$$\mathcal{L}_{ls} := \frac{n}{B} \sum_{i=1}^B (\mathbf{y}_{ls})_i^T \log f_i = \frac{n}{BS} \sum_{i=1}^B \mathbf{y}_i^T \log f_i. \tag{14}$$

Therefore the only difference between  $\mathcal{L}_c$  and the "human-aware" label-smoothing loss  $\mathcal{L}_{ls}$  is that the former is  $S$  times stronger, and they are conceptually very similar. We will also consider "standard" label smoothing with parameter  $\alpha \in (0, 1)$ , :

$$\mathcal{L}_\alpha := \frac{n}{B} \sum_{i=1}^B \hat{\mathbf{y}}_i^T \log f_i, \tag{15}$$

where  $\hat{y}_i$  is  $1 - \alpha$  in the position corresponding to the correct label and  $\frac{\alpha}{C}$  otherwise.

For each of the likelihoods proposed above, we train a ResNet-20 on CIFAR-10H and evaluate on CIFAR-10 training set, as in (Aitchison, 2020). Regarding the SG-MCMC method, we use cyclical step size and layer preconditioning, adapting the code from (Wenzel et al., 2020). We leave unchanged all the hyperparameters. The only exception is in the case we are using  $\mathcal{L}_c$ , where we reduce the learning rate by a factor 50, the approximate number of labellers per image, as in (Aitchison, 2020). The reason for this reduction is that the likelihood gets  $\approx 50$  times stronger due to the fact that each datapoint is labelled by  $\approx 50$  labellers. We use 200 epochs as burn-in period and a cycle length of 100 epochs. We use data augmentation as in (Aitchison, 2020).

**Is it all about over-weighting the likelihood?** We apply standard label smoothing (loss  $\mathcal{L}_\alpha$ ) with  $\alpha = 0.1$  to the one hot encoded labels. Then, we overweight the likelihood by a factor of 50 and reduce the learning by the same factor (i.e. the approximate number of labellers). Results are in Figure 16.

Note that making the likelihood  $S$  times stronger helps to eliminate the cold posterior effect. Note also that label smoothing plus likelihood over-weighting is equivalent to assuming that the labellers have assigned different labels.

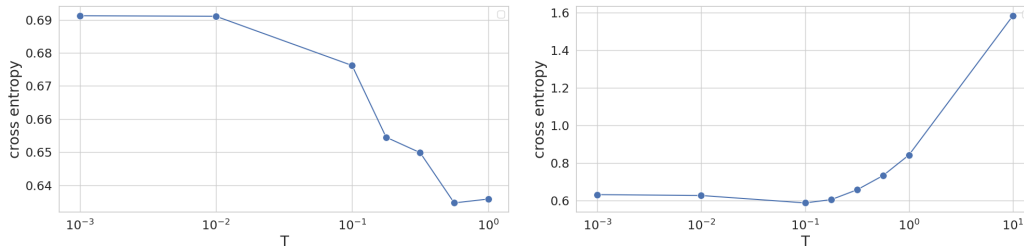


Figure 16: **left**: model with over-weighted likelihood with smoothed labels. **right**: model trained with standard label smoothing

**A smaller number of labellers is enough to alleviate the cold posterior effect** The previous experiments have shown that over-weighting the likelihood alleviates the cold posterior effect. Here, we test whether weighting the likelihood by a smaller amount is sufficient to alleviate the cold posterior effect. In particular, we subsample a different number of  $\tilde{S}$  labels for each image from the counts that are available in CIFAR-10H. We do this by considering the probabilities implied by the counts, and sampling  $\tilde{S}$  times from this categorical distribution, for each image. Results are shown in Figure 17. Note that  $T \geq 1$  is optimal for any value of  $\tilde{S} \geq 10$  that we consider, and the cold posterior effect appears only for  $\tilde{S} \leq 5$ . This is strong evidence that there is not an exact correspondence between the number of labellers and the optimal temperature. What matters is to over-weight the likelihood, but even a significantly smaller weight can fix the cold posterior problem.

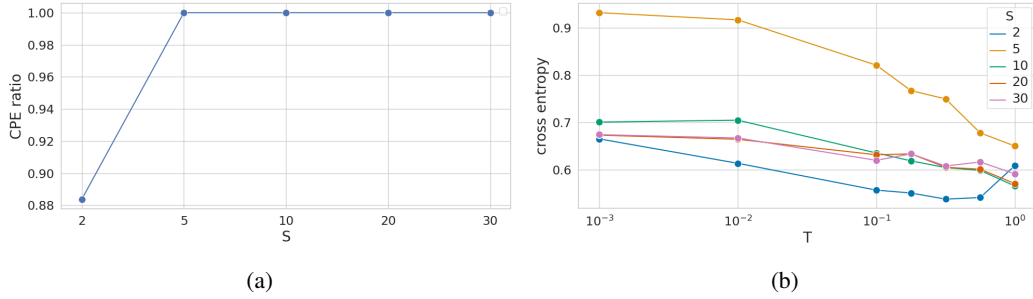


Figure 17: CPER and test CE for subsample of  $\tilde{S}$  labellers from Cifar10H out of a total of  $S \approx 50$  available labellers.