

A LAPLACE PREDICTIVE POSTERIOR APPROXIMATIONS

In this section, we review several methods for obtaining the predictive posterior of Laplace approximation.

A.1 MONTE CARLO SAMPLING

As described in main text, we can obtain a closed-form Gaussian posterior on output logits,

$$f_{\theta}(\mathbf{x}_*) \sim \mathcal{N}(f_{\theta_{\text{MAP}}}(\mathbf{x}_*), \mathbf{\Lambda}), \quad (15)$$

where

$$\mathbf{\Lambda} = (\nabla_{\theta} f_{\theta}(\mathbf{x}_*)|_{\theta_{\text{MAP}}})^T \mathbf{\Sigma} (\nabla_{\theta} f_{\theta}(\mathbf{x}_*)|_{\theta_{\text{MAP}}}). \quad (16)$$

To obtain samples of $f_{\theta}(\mathbf{x}_*)$, we can decompose the covariance using the Cholesky factorization $\mathbf{\Lambda} = \mathbf{L}\mathbf{L}^T$ with

$$\tilde{f}_{\theta}(\mathbf{x}_*) = f_{\theta_{\text{MAP}}}(\mathbf{x}_*) + \mathbf{L}\boldsymbol{\xi}, \quad (17)$$

where $\boldsymbol{\xi}$ is a vector of IID standard normal random variables. We can compute the Bayesian model average by computing the average probabilities (passing the sampled logits through softmax function) under the Gaussian random noise from $\boldsymbol{\xi}$.

A.2 PROBIT APPROXIMATION

A closed-form approximation of the predictive posterior for classification can be obtained by integrating out the posterior over weights with a generalized probit approximation (Lu et al., 2020; Daxberger et al., 2021a) of the likelihood,

$$p(y_*|\mathbf{x}_*, \mathcal{D}) \approx \text{Categorical} \left(y_*, \text{softmax} \left(\frac{f_{\theta_{\text{MAP}}}(\mathbf{x}_*)}{\sqrt{1 + \frac{\pi}{8} \text{diag}(\mathbf{\Lambda})}} \right) \right), \quad (18)$$

where

$$\mathbf{\Lambda} = (\nabla_{\theta} f_{\theta}(\mathbf{x}_*)|_{\theta_{\text{MAP}}})^T \mathbf{\Sigma} (\nabla_{\theta} f_{\theta}(\mathbf{x}_*)|_{\theta_{\text{MAP}}}). \quad (19)$$

Although probit approximation provably preserves decision boundary in binary sigmoid classification (Kristiadi et al., 2020), it does hold for softmax multiclass classification.

A.3 LAPLACE BRIDGE

The Laplace bridge (MacKay, 1998; Daxberger et al., 2021a) maps a Gaussian $\mathcal{N}(\boldsymbol{\mu}, \mathbf{\Sigma})$ to a Dirichlet distribution $\mathcal{D}(\boldsymbol{\alpha})$ over classes with

$$\alpha_i = \frac{1}{\sum_{ii}} \left(1 - \frac{2}{C} + \frac{\exp(\boldsymbol{\mu}_i)}{C^2} \sum_j \exp(-\boldsymbol{\mu}_j) \right), \quad (20)$$

where C denotes the number of classes. Similar to the generalized probit approximation, it also ignores the covariance terms and only considers the diagonal of the covariance $\mathbf{\Sigma}_{ii}$.

A.4 EMPIRICAL COMPARISON

Figure 3 shows a comparison of these approximations in fine-tuning Llama2-7B and applying full Laplace-LoRA post-hoc. Specifically, we considered Monte Carlo sampling using the full covariance (MC joint), Monte Carlo sampling only using the diagonal covariance (MC indep), generalized probit approximation (probit), and Laplace bridge (bridge). MC joint consistently achieves highest accuracy and among the best NLL, while bridge is often the worst, probit and MC indep can sometimes give suboptimal performance. This is likely due to bridge, probit and MC indep are all approximations that ignored the covariances between logits, whereas MC joint faithfully approximates the true predictive posterior.

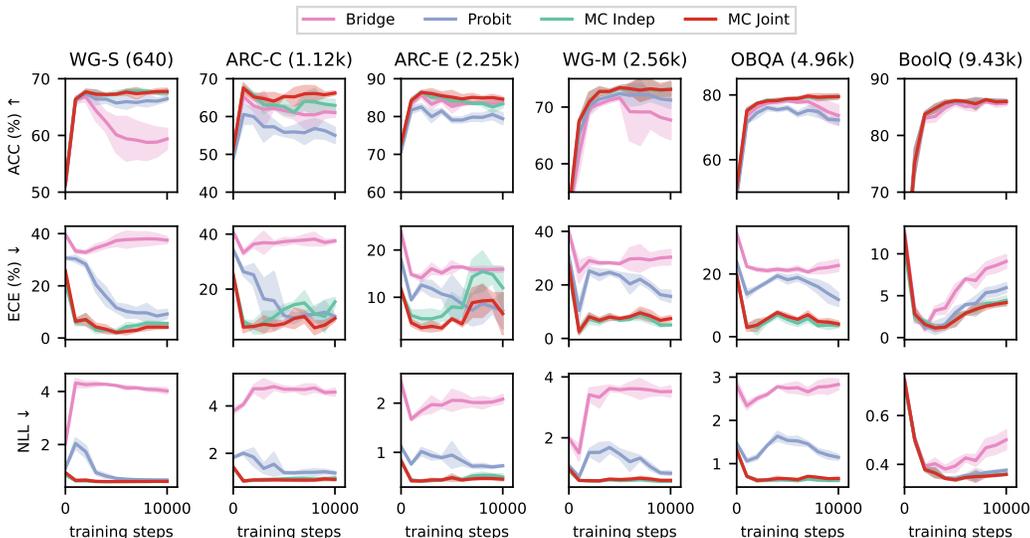


Figure 3: Fine-tuning of LLaMA2-7B across six common sense reasoning tasks, comparing different Laplace predictive posterior approximations: Laplace bridge approximation (bridge), generalized probit approximation (probit), Monte Carlo sampling using the diagonal covariance (MC indep), and Monte Carlo sampling using the full covariance (MC joint).

Task	Prompt
Winogrande (WG-S/WG-M)	Select one of the choices that answers the following question: {question} Choices: A. {option1}. B {option2}. Answer:
ARC (ARC-C/ARC-E)	Select one of the choices that answers the following question: {question} Choices: A. {choice1}. B. {choice2}. C. {choice2}. D. {choice2}. Answer:
Openbook QA (OBQA)	Select one of the choices that answers the following question: {question} Choices: A. {choice1}. B. {choice2}. C. {choice2}. D. {choice2}. Answer:
BoolQ	Answer the question with only True or False: {question} Context: {passage}.

Table 5: Prompt templates for fine-tuning LLaMA-7B on common sense reasoning tasks.

B PROMPT TEMPLATES FOR LLAMA2 COMMON SENSE REASONING TASKS

We present our prompt templates used to fine-tune LLaMA2 on common sense reasoning tasks in Table 5. For ARC datasets, although the majority of questions have four choices, there are a tiny amount of questions with three or five choices which we remove for consistency. In ARC-C, there are 1/1119 with three choices and 1/1119 with five choices in the training set; 3/299 with three choices and 1/299 with five choices in the evaluation set. In ARC-E, there are 6/2251 with three choices and 4/2251 with five choices in the training set; 1/570 with three choices and 2/570 with five choices in the evaluation set.

C HYPERPARAMETERS

We follow the default hyperparameters from Huggingface’s Transformer (Wolf et al., 2020) and PEFT (Mangrulkar et al., 2022) libraries. Hyperparameters used for fine-tuning RoBERTa-base and RoBERTa-large with LoRA are shown in Table 6, those for fine-tuning LLaMA-7B are shown in Table 7. Note the only differences between the two settings are a smaller batch size to fit into our GPU memory, and a longer max sequence length to account for prompt length.

Hyperparameter	Value
LoRA r	8
LoRA α	16
Dropout Probability	0.1
Weight Decay	0
Learning Rate	5×10^{-5}
Learning Rate Scheduler	Linear
Batch Size	32
Max Sequence Length	256

Table 6: Hyperparameters used in fine-tuning RoBERTa-base and RoBERTa-large with LoRA.

Hyperparameter	Value
LoRA r	8
LoRA α	16
Dropout Probability	0.1
Weight Decay	0
Learning Rate	5×10^{-5}
Learning Rate Scheduler	Linear
Batch Size	4
Max Sequence Length	300

Table 7: Hyperparameters used in fine-tuning LLaMA-7B with LoRA.

D METRICS FOR UNCERTAINTY QUANTIFICATION

There are two commonly used metrics for measuring uncertainty quantification in neural networks: negative log-likelihood (NLL) and expected calibration error (ECE). NLL computes the sum of negative expected log probability of predicting the true label,

$$\text{NLL} = \sum_{i=1}^N -\log P(\hat{y}_i = y_i), \quad (21)$$

where $P(\hat{y}_i)$ is the model’s output distribution, y_i is the true label. NLL is also equivalent to cross-entropy between the one-hot data distribution and the model’s output distribution. NLL encourages the model to give high probability to correct answers. If the model is overconfident in a wrong answer, then the probability of giving the right answer would be low which raises NLL. On the other hand, ECE measures the alignment between model’s confidence and accuracy, by binning the highest predicted probabilities and compute a weighted average between the difference of average accuracy and confidence in each bin,

$$\text{ECE} = \sum_{m=1}^M \frac{|B_m|}{N} |\text{acc}(B_m) - \text{conf}(B_m)|, \quad (22)$$

where $\text{acc}(B_m)$ and $\text{conf}(B_m)$ are the average accuracy and average confidence in each bin,

$$\text{acc}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} \mathbf{1}(\hat{y}_i = y_i), \quad \text{conf}(B_m) = \frac{1}{|B_m|} \sum_{i \in B_m} P(\hat{y}_i), \quad (23)$$

and $|B_m|$ is the number of examples in bin m . However, expected calibration error cannot be optimized directly like negative log-likelihood, as a completely random model will have the same accuracy and confidence for each datapoint, thus achieving zero ECE (Ashukha et al., 2020).

E OPTIMIZING LAPLACE PRIOR PRECISION

In this section, we present how we optimize the Laplace prior precision. When there is only a training set available with no validation set (such as in Figure 1, 4, 5 and Table 1, 8, 9), we can use the

Algorithm 1 Optimize Laplace prior precision using the training set model evidence

```

Initialize prior precision  $\lambda$ , learning rate  $\eta$ , optimization steps  $M$ 
Obtain  $\theta_{\text{MAP}}$  from a fine-tuned checkpoint, pre-compute  $\log P(\mathbf{y}|\mathbf{X}, \theta)$  and Fisher  $\mathbf{F}$ 
for  $step = 1, \dots, M$  do
  Compute posterior covariance  $\Sigma = \mathbf{F} + \lambda \mathbf{I}$ 
  Calculate  $\mathcal{L}(\mathbf{y}, \mathbf{X}; \theta) = \log P(\mathbf{y}|\mathbf{X}, \theta) + \log P(\theta) = \log P(\mathbf{y}|\mathbf{X}, \theta) + \lambda \|\theta_{\text{MAP}}\|_2^2$ 
  Perform a gradient step with respect to  $\lambda$  to maximize log model evidence (Eq.13):
     $\lambda \leftarrow \lambda + \eta \nabla_{\lambda} (-\mathcal{L}(\mathbf{y}, \mathbf{X}; \theta_{\text{MAP}}) + \frac{1}{2} \log |\Sigma|)$ 
end for

```

Algorithm 2 Optimize Laplace prior precision using validation log-likelihood

```

Initialize prior precision  $\lambda$ , learning rate  $\eta$ , batch size  $b$ , validation set  $(\mathbf{X}, \mathbf{y})$ 
Obtain  $\theta_{\text{MAP}}$  from a fine-tuned checkpoint, pre-compute mean  $f_{\theta_{\text{MAP}}}$ , Jacobian  $\mathbf{J} = \nabla_{\theta} f_{\theta}(\mathbf{X})|_{\theta_{\text{MAP}}}$ ,
Fisher  $\mathbf{F}$ 
for  $step = 1, \dots, M$  do
  Randomly sample a batch of validation data  $\mathbf{X}_b, \mathbf{y}_b$  with corresponding Jacobian  $\mathbf{J}_b$ 
  Compute posterior covariance  $\Sigma = \mathbf{F} + \lambda \mathbf{I}$ 
  Calculate batch logits covariance  $\Lambda_b = \mathbf{J}_b^T \Sigma \mathbf{J}_b$  and Cholesky  $\mathbf{L}_b = \mathbf{L}_b \mathbf{L}_b^T$ 
  Obtain batch Bayesian model average  $\tilde{f}_{\theta}(\mathbf{X}_b) = f_{\theta_{\text{MAP}}}(\mathbf{X}_b) + \mathbf{L}_b \xi$ 
  Evaluate validation likelihood  $P(\mathbf{y}_b|\mathbf{X}_b, \theta) = \text{Categorical}(\mathbf{y}_b; \text{softmax}(\tilde{f}_{\theta}(\mathbf{X}_b)))$ 
  Perform a gradient step with respect to  $\lambda$  to maximize mini-batch validation log-likelihood:
     $\lambda \leftarrow \lambda + \eta \nabla_{\lambda} \log P(\mathbf{y}_b|\mathbf{X}_b, \theta)$ 
end for

```

Laplace model evidence in Equation 13 to optimize prior precision. Our algorithm is presented in Algorithm 1, and we chose $\eta = 0.1$, and $M = 100$.

When we introduce a validation set by splitting the training set (such as in Figure 8, 6, 7 and Table 2, 10, 11), we can use the validation log-likelihood to optimize the Laplace prior precision. For memory and computational efficiency, we precompute the mean $f_{\theta_{\text{MAP}}}$ and the Jacobian $\mathbf{J} = \nabla_{\theta} f_{\theta}(\mathbf{X})|_{\theta_{\text{MAP}}}$, then perform mini-batch gradient descent on λ (reparametrization in Bayesian model average allows gradient flowing through) as detailed in Algorithm 2. We chose $\eta = 0.1$, $M = 1000$, and $b = 4$.

F ADDITIONAL EXPERIMENTS

F.1 FINE-TUNING ROBERTA FOR TEXT CLASSIFICATION

In this section, we present additional results of fine-tuning RoBERTa-base (Fig. 4) and RoBERTa-large (Fig. 5) (Liu et al., 2019) with LoRA on text classification tasks from GLUE (Wang et al., 2019b) and SuperGLUE (Wang et al., 2019a). Results for RoBERTa-base are shown in Figure 4 and Table 8, and results for RoBERTa-large are shown in Figure 5 and Table 9. Surprisingly, checkpoint ensemble and Monte-Carlo (MC) dropout exhibit distinct behavior on RoBERTa models compared to LLaMA2-7B. Checkpoint ensemble often performs much worse than the Maximum a Posteriori (MAP) estimation in terms of Expected Calibration Error (ECE) and Negative Log-Likelihood (NLL), while MC dropout often offers much more improvement on RoBERTa models. We suspect this difference arises due to an extra hidden penultimate layer with an additional dropout layer in front in the default RoBERTa fine-tuning setup; whereas, in LLaMA2-7B fine-tuning, we have LoRA weights and dropout on LoRA layers at the end. However, the gain from MC dropout diminishes on RoBERTa-large compared to RoBERTa-base. On the other hand, Laplace-LoRA (LA) consistently delivers substantial gains on any models we have tested (RoBERTa-base, RoBERTa-large, and LLaMA2-7B), demonstrating the robustness of Laplace-LoRA. Moreover, the Last-Layer Laplace-LoRA (LLLA) offers modest improvements as in LLaMA2-7B when optimized by Laplace model evidence, underscoring the significance of performing Bayesian inference on all LoRA weights.

Similarly, we also conduct experiments by splitting the training set into a 80% training set and a 20% validation set, then tune temperature and Laplace prior precision on the validation set. The results are

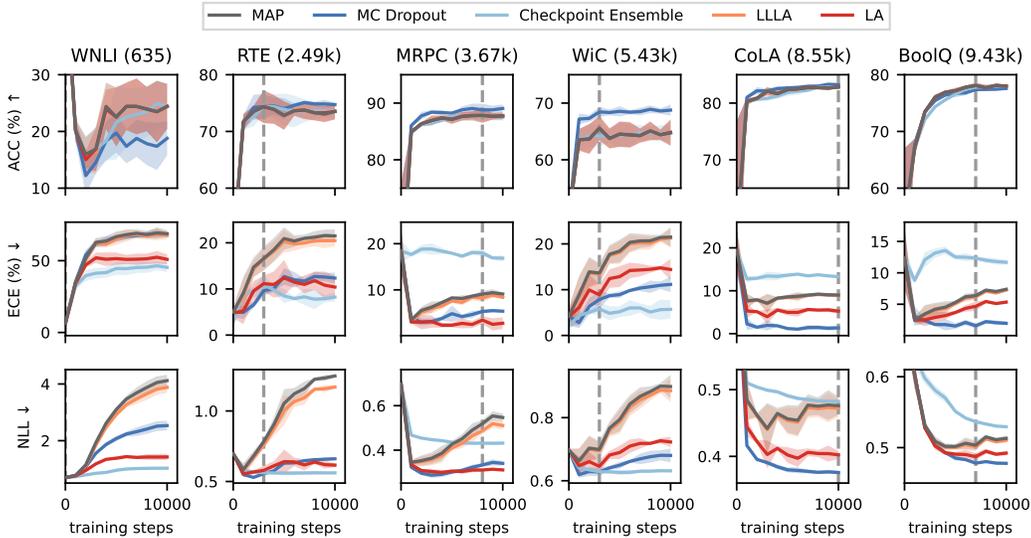


Figure 4: Fine-tuning of RoBERTa-base across six GLUE and SuperGLUE tasks (presented column-wise, with number of training examples in brackets), evaluated on the test set every 1000 gradient steps. The vertical dashed line gives the number of training steps with optimal MAP performance. Note that RoBERTa-base seems to fail on WNLI, but RoBERTa-large succeeds (Fig. 5).

Table 8: Comparison of different post-hoc methods applied to the fine-tuned RoBERTa-base across six GLUE and SuperGLUE tasks. Results are evaluated at the early stopping point of 5000 gradient steps.

Methods	Metrics	WNLI	RTE	MRPC	WiC	CoLA	BoolQ
ACC ↑	MAP	22.5 _{4.6}	72.8 _{2.2}	87.1 _{0.7}	64.5 _{2.0}	82.4 _{0.6}	77.2 _{0.4}
	MC Drop	19.7 _{3.0}	74.0 _{1.6}	88.2 _{0.0}	68.5 _{0.8}	82.7 _{0.1}	76.7 _{0.6}
	Ckpt Ens	21.6 _{5.8}	73.9 _{2.5}	87.7 _{0.6}	64.3 _{1.9}	81.8 _{0.5}	76.6 _{0.2}
	LLLA	22.5 _{4.6}	72.8 _{2.2}	87.1 _{0.7}	64.5 _{2.0}	82.4 _{0.6}	77.2 _{0.4}
	LA	22.5 _{4.6}	72.8 _{2.2}	87.2 _{0.7}	64.5 _{2.0}	82.4 _{0.6}	77.2 _{0.5}
ECE ↓	MAP	66.7 _{3.2}	20.9 _{2.1}	8.3 _{0.1}	18.4 _{2.2}	8.6 _{0.5}	5.3 _{0.2}
	MC Drop	65.6 _{1.8}	12.7 _{1.6}	4.8 _{1.6}	8.9 _{0.6}	1.2_{0.1}	1.5_{0.4}
	Ckpt Ens	44.6_{3.9}	8.5_{1.3}	17.9 _{0.1}	5.2_{1.0}	13.6 _{0.3}	12.5 _{0.3}
	LLLA	65.2 _{3.7}	20.3 _{2.2}	7.6 _{0.3}	18.1 _{2.2}	8.5 _{0.4}	5.1 _{0.1}
	LA	51.4 _{3.8}	12.4 _{3.4}	2.3_{0.4}	12.9 _{2.0}	5.0 _{0.6}	3.7 _{0.2}
NLL ↓	MAP	3.10 _{0.09}	1.05 _{0.09}	0.43 _{0.01}	0.79 _{0.03}	0.46 _{0.02}	0.50 _{0.00}
	MC Drop	2.07 _{0.08}	0.63 _{0.01}	0.30_{0.01}	0.65 _{0.01}	0.38_{0.00}	0.49_{0.00}
	Ckpt Ens	0.96_{0.02}	0.56_{0.01}	0.44 _{0.00}	0.63_{0.01}	0.49 _{0.00}	0.55 _{0.00}
	LLLA	2.92 _{0.08}	1.00 _{0.09}	0.41 _{0.01}	0.79 _{0.03}	0.45 _{0.02}	0.50 _{0.00}
	LA	1.38 _{0.06}	0.64 _{0.08}	0.30_{0.00}	0.69 _{0.02}	0.40 _{0.01}	0.49_{0.00}

shown in Figure 6 and Table 10 for RoBERTa-base, and Figure 7 and Table 11 for RoBERTa-large. Again, full Laplace-LoRA LA offers the most improvements most of the time, and is the most robust method overall.

F.2 FINE-TUNING LLAMA2-7B FOR COMMON SENSE REASONING

Figure 8 displays the results of tuning temperature scaling and Laplace prior precision on a held-out validation set split from the training set. LLLA is missing a few zero checkpoint evaluation results due to Cholesky errors. Comparing Figure 8 to Figure 1 in the main text, it is evident that splitting the training set into a smaller training set (80% data) and a validation set (20% data) slightly impact the accuracy of fine-tuned model.

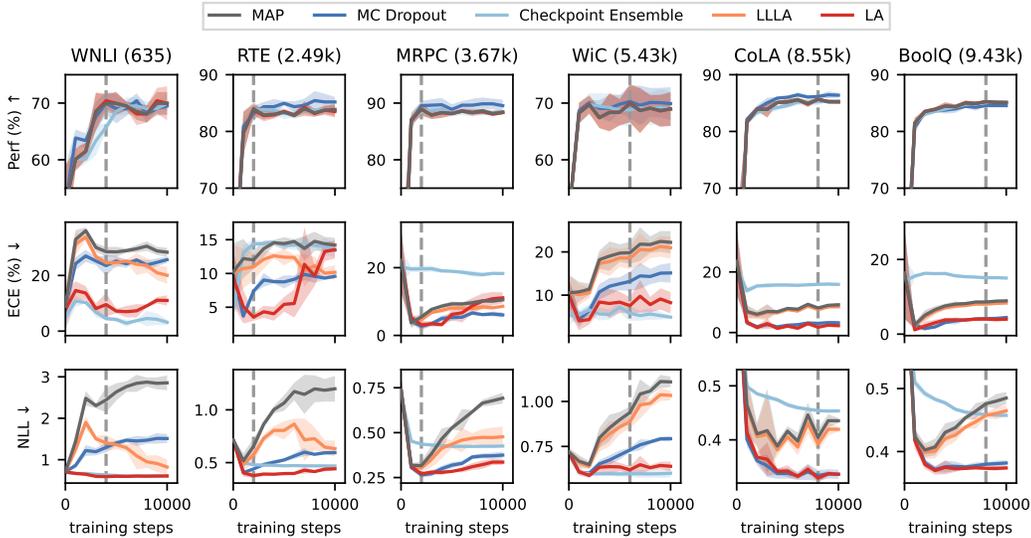


Figure 5: Fine-tuning RoBERTa-large across the six GLUE and SuperGLUE tasks in Fig. 4. The vertical dashed line gives the number of training steps with optimal MAP performance.

Table 9: Comparison of different post-hoc methods applied to the fine-tuned RoBERTa-large across six GLUE and SuperGLUE tasks. Results are evaluated at the early stopping point of 5000 gradient steps.

Methods	Metrics	WNLI	RTE	MRPC	WiC	CoLA	BoolQ
ACC ↑	MAP	70.0 _{1.8}	83.4 _{0.3}	88.2 _{0.3}	68.4 _{2.0}	85.2 _{0.5}	84.4 _{0.4}
	MC Drop	69.0 _{2.3}	85.0 _{0.3}	88.9 _{0.8}	69.8 _{1.4}	86.0 _{0.1}	84.0 _{0.6}
	Ckpt Ens	68.5 _{0.7}	83.2 _{0.9}	88.2 _{0.5}	68.8 _{2.4}	84.6 _{0.4}	84.3 _{0.5}
	LLLA	70.0 _{1.8}	83.5 _{0.3}	88.2 _{0.3}	68.4 _{2.0}	85.2 _{0.5}	84.3 _{0.3}
	LA	70.0 _{1.8}	83.5 _{0.3}	88.2 _{0.3}	68.4 _{2.0}	85.2 _{0.5}	84.4 _{0.4}
ECE ↓	MAP	28.5 _{1.6}	14.4 _{0.8}	9.3 _{0.6}	19.7 _{1.8}	8.0 _{0.4}	8.0 _{0.6}
	MC Drop	25.2 _{1.6}	8.8 _{0.9}	5.3_{1.0}	12.6 _{1.3}	2.4 _{0.3}	3.6_{0.9}
	Ckpt Ens	4.2_{0.7}	14.2 _{0.5}	18.9 _{0.2}	7.2_{1.2}	15.7 _{0.5}	15.6 _{0.3}
	LLLA	24.6 _{1.7}	12.4 _{1.1}	8.1 _{0.6}	18.6 _{1.8}	7.6 _{0.3}	7.5 _{0.7}
	LA	7.2 _{0.7}	5.4_{1.0}	6.2 _{1.5}	8.4 _{2.1}	2.1_{0.3}	3.8 _{0.5}
NLL ↓	MAP	2.64 _{0.18}	1.06 _{0.06}	0.53 _{0.01}	0.90 _{0.04}	0.42 _{0.01}	0.44 _{0.01}
	MC Drop	1.40 _{0.06}	0.53 _{0.01}	0.32 _{0.02}	0.70 _{0.02}	0.34_{0.01}	0.38 _{0.01}
	Ckpt Ens	0.62 _{0.01}	0.47 _{0.00}	0.42 _{0.00}	0.60_{0.01}	0.46 _{0.00}	0.47 _{0.00}
	LLLA	1.36 _{0.10}	0.80 _{0.08}	0.44 _{0.02}	0.86 _{0.04}	0.41 _{0.01}	0.43 _{0.01}
	LA	0.60_{0.02}	0.40_{0.01}	0.28_{0.00}	0.64 _{0.02}	0.34_{0.01}	0.37_{0.00}

F.3 DIAGONAL LAPLACE APPROXIMATION

F.3.1 ROBERTA

In this section, we present the results for Laplace-LoRA utilizing a diagonal approximation of the Hessian. This approach is generally not found to be as effective as the Kronecker-factored Approximate Curvature (K-FAC) (Daxberger et al., 2021a) that approximates the block diagonal Hessian. The results on RoBERTa-base and RoBERTa-large are shown in Figure 9 Table 12, and Figure 10 Table 13. LLLA still offers a tiny advantage over Maximum a Posteriori (MAP) in ECE and NLL, however, LA show mixed performance across the datasets.

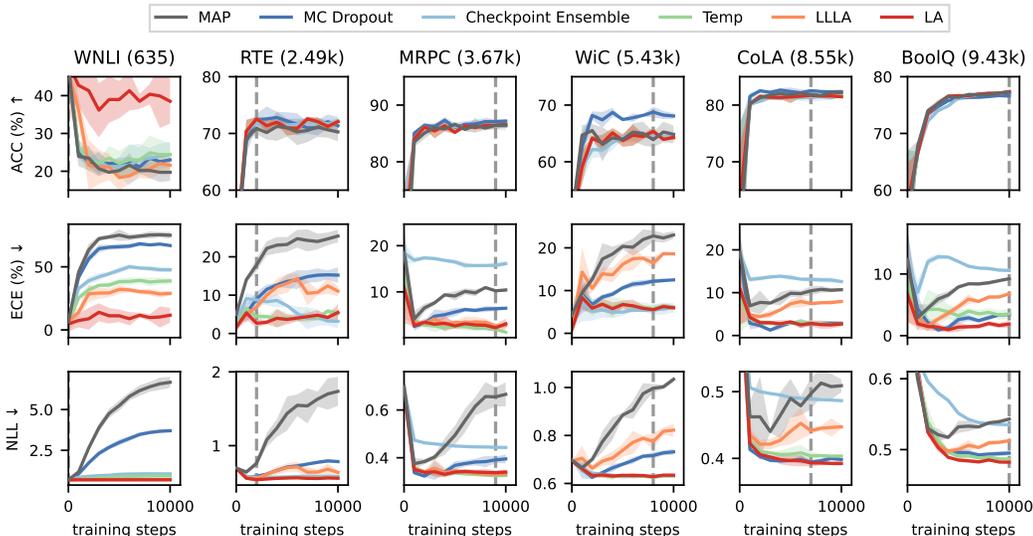


Figure 6: Fine-tuning RoBERTa-base across the six GLUE and SuperGLUE tasks. The vertical dashed line gives the checkpoint with optimal MAP performance on a held-out validation set.

Table 10: Comparison of different post-hoc methods applied to the fine-tuned RoBERTa-base across six GLUE and SuperGLUE tasks, with a validation set split from the training set used for tuning temperature and Laplace prior precision. Results are evaluated at the best MAP performance checkpoint observed on the validation set.

Methods	Metrics	WNLI	RTE	MRPC	WiC	CoLA	BoolQ
ACC ↑	MAP	47.9 _{6.0}	70.9 _{0.9}	86.4 _{0.6}	63.9 _{1.2}	81.8 _{0.3}	77.2 _{0.2}
	MC Drop	47.9 _{6.0}	72.4 _{0.9}	87.1 _{0.3}	68.8 _{0.9}	82.6 _{0.1}	76.6 _{0.4}
	Ckpt Ens	47.9 _{6.0}	71.8 _{0.9}	86.3 _{0.5}	64.7 _{0.3}	81.4 _{0.7}	77.2 _{0.1}
	Temp	47.9 _{6.0}	72.6 _{0.3}	86.5 _{0.3}	65.4 _{0.8}	81.8 _{0.8}	77.3 _{0.2}
	LLLA	46.0 _{10.0}	72.6 _{0.3}	86.4 _{0.4}	65.3 _{0.7}	81.8 _{0.7}	77.4 _{0.2}
	LA	48.4 _{7.7}	72.6 _{0.3}	86.4 _{0.3}	65.4 _{0.6}	81.7 _{0.7}	77.4 _{0.2}
ECE ↓	MAP	6.3 _{4.3}	17.8 _{1.4}	10.2 _{0.2}	22.8 _{1.1}	10.4 _{0.9}	9.2 _{0.1}
	MC Drop	6.4 _{4.4}	9.0 _{1.7}	6.3 _{0.4}	12.2 _{0.5}	2.8_{0.3}	3.4 _{0.1}
	Ckpt Ens	6.2 _{2.3}	8.8 _{0.9}	15.7 _{0.6}	5.5 _{1.0}	13.1 _{0.5}	10.6 _{0.1}
	Temp	5.5 _{1.7}	4.5 _{1.3}	1.9_{0.3}	5.8 _{0.9}	2.9 _{0.6}	3.5 _{0.7}
	LLLA	6.8 _{7.8}	6.9 _{0.7}	2.6 _{0.4}	16.4 _{1.4}	7.4 _{0.3}	6.8 _{0.2}
	LA	4.7_{5.2}	2.7_{1.0}	2.2 _{0.7}	5.5_{0.5}	2.8_{0.6}	1.9_{0.5}
NLL ↓	MAP	0.70 _{0.01}	0.76 _{0.03}	0.66 _{0.04}	1.00 _{0.02}	0.50 _{0.02}	0.54 _{0.00}
	MC Drop	0.70 _{0.01}	0.58 _{0.03}	0.39 _{0.01}	0.72 _{0.01}	0.39 _{0.00}	0.50 _{0.00}
	Ckpt Ens	0.69_{0.00}	0.57 _{0.00}	0.44 _{0.00}	0.63 _{0.00}	0.49 _{0.00}	0.53 _{0.00}
	Temp	0.69_{0.00}	0.54_{0.00}	0.32_{0.00}	0.62_{0.01}	0.40 _{0.01}	0.49 _{0.00}
	LLLA	0.69_{0.00}	0.56 _{0.01}	0.33 _{0.00}	0.78 _{0.03}	0.44 _{0.00}	0.51 _{0.00}
	LA	0.69_{0.00}	0.54_{0.00}	0.34 _{0.01}	0.62_{0.01}	0.39_{0.00}	0.48_{0.00}

F.3.2 LLAMA2-7B

On LLaMA2-7B, both LLLA and LA show mixed performance and do not offer consistent improvements as shown in Figure 11 and Table 14. Specifically, the accuracy of diagonal LA is much worse than MAP on ARC datasets, and ECE is worse than MAP on ARC-E, OBQA and BoolQ.

On the other hand, when we split a validation set from the training set and tune the Laplace prior precision using the validation log-likelihood, the results are much better, as shown in Figure 12 and Table 15.

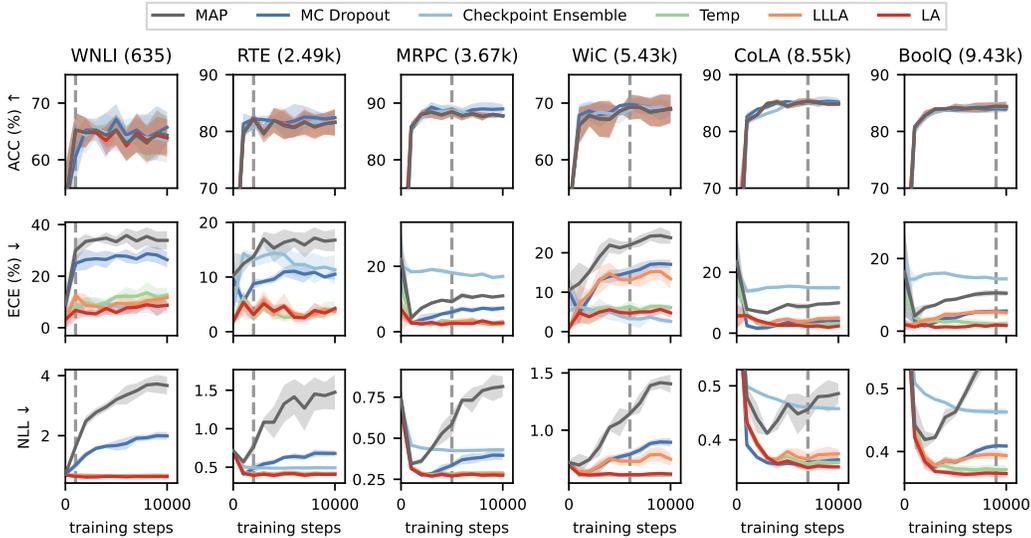


Figure 7: Fine-tuning RoBERTa-large across the six GLUE and SuperGLUE tasks. The vertical dashed line gives the checkpoint with optimal MAP performance on a held-out validation set.

Table 11: Comparison of different post-hoc methods applied to the fine-tuned RoBERTa-large across six GLUE and SuperGLUE tasks, with a validation set split from the training set used for tuning temperature and Laplace prior precision. Results are evaluated at the best MAP performance checkpoint observed on the validation set.

Methods	Metrics	WNLI	RTE	MRPC	WiC	CoLA	BoolQ
ACC ↑	MAP	65.3 _{2.7}	82.2 _{0.9}	88.5 _{0.7}	69.3 _{1.9}	85.3 _{0.6}	84.4 _{0.6}
	MC Drop	60.6 _{4.0}	82.2 _{0.3}	88.7 _{0.7}	69.8 _{0.8}	85.1 _{0.3}	84.0 _{0.6}
	Ckpt Ens	65.3 _{2.7}	80.7 _{0.3}	88.6 _{1.1}	68.8 _{1.0}	85.2 _{0.4}	84.3 _{0.8}
	Temp	65.3 _{2.7}	82.2 _{0.9}	88.5 _{0.7}	69.3 _{1.9}	85.3 _{0.6}	84.4 _{0.6}
	LLLA	65.3 _{2.7}	82.2 _{0.9}	88.5 _{0.7}	69.3 _{1.9}	85.3 _{0.6}	84.4 _{0.6}
	LA	65.3 _{2.7}	82.3 _{0.8}	88.5 _{0.7}	69.3 _{2.0}	85.3 _{0.5}	84.4 _{0.6}
ECE ↓	MAP	30.0 _{3.5}	13.7 _{0.4}	9.2 _{0.7}	22.0 _{0.8}	8.8 _{0.8}	10.5 _{0.6}
	MC Drop	24.9 _{2.6}	8.8 _{0.8}	6.0 _{0.9}	15.1 _{1.1}	3.3 _{0.8}	5.5 _{0.4}
	Ckpt Ens	7.7 _{3.0}	13.1 _{1.2}	18.0 _{0.5}	3.3 _{1.6}	15.0 _{0.3}	14.4 _{0.1}
	Temp	8.1 _{2.2}	2.6 _{0.7}	3.1 _{0.4}	6.4 _{1.3}	2.7 _{0.8}	1.9 _{0.7}
	LLLA	12.7 _{0.9}	3.3 _{1.0}	3.0 _{0.5}	13.1 _{1.2}	4.0 _{0.4}	5.3 _{0.8}
	LA	6.8 _{3.3}	3.1 _{1.1}	2.6 _{0.7}	4.7 _{1.7}	2.2 _{0.2}	1.7 _{0.4}
NLL ↓	MAP	1.63 _{0.09}	0.77 _{0.07}	0.59 _{0.06}	1.14 _{0.02}	0.46 _{0.02}	0.56 _{0.01}
	MC Drop	1.02 _{0.13}	0.48 _{0.06}	0.33 _{0.01}	0.79 _{0.02}	0.36 _{0.01}	0.41 _{0.00}
	Ckpt Ens	0.65 _{0.02}	0.49 _{0.01}	0.42 _{0.00}	0.61 _{0.01}	0.46 _{0.00}	0.46 _{0.00}
	Temp	0.65 _{0.01}	0.41 _{0.01}	0.28 _{0.02}	0.61 _{0.01} <i>s</i>	0.36 _{0.00}	0.37 _{0.01}
	LLLA	0.69 _{0.02}	0.40 _{0.02}	0.27 _{0.01}	0.73 _{0.03}	0.36 _{0.01}	0.40 _{0.01}
	LA	0.64 _{0.00}	0.39 _{0.01}	0.27 _{0.01}	0.61 _{0.01}	0.35 _{0.01}	0.37 _{0.01}

G OUT-OF-DISTRIBUTION EVALUATION

Here we present the specific far OOD MMLU datasets we used for evaluations in Table 3 and Table 4. The specific task splits we selected for each subject are shown in Table 16 assigned by Hendrycks et al. (2020).

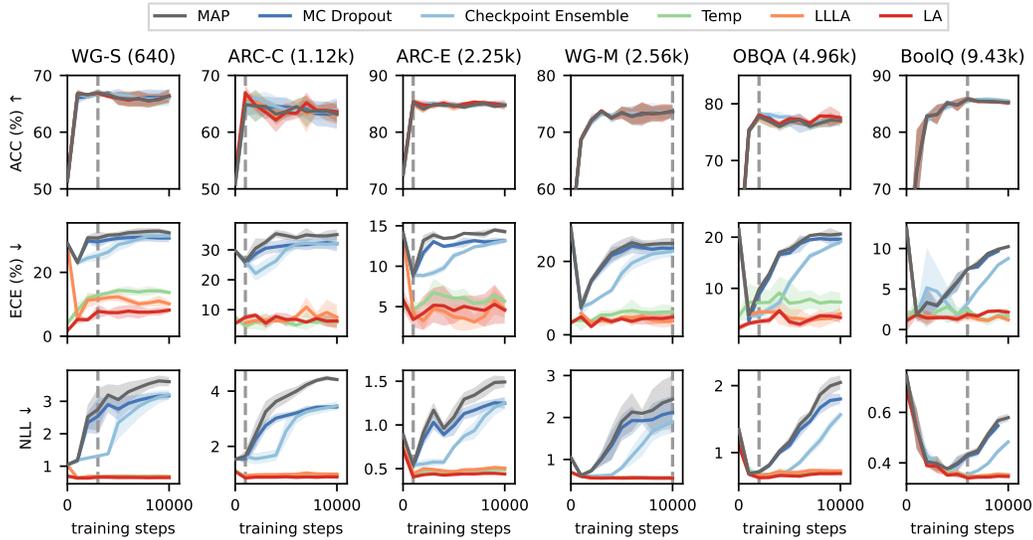


Figure 8: Fine-tuning of LLaMA2-7B across six common sense reasoning tasks (presented column-wise, with number of training examples in brackets), evaluated on the test set every 1000 gradient steps. The vertical dashed line gives the checkpoint with optimal MAP performance on a held-out validation set.

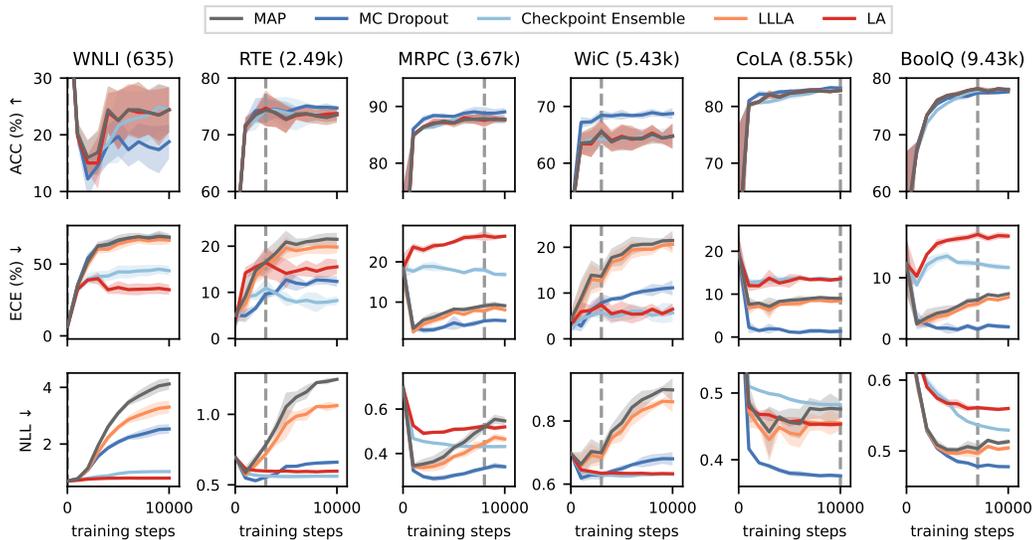


Figure 9: Fine-tuning of RoBERTa-base across six GLUE and SuperGLUE tasks (presented column-wise, with number of training examples in brackets), evaluated on the test set every 1000 gradient steps. The vertical dashed line gives the number of training steps with optimal MAP performance. LA and LLLA using diagonal Fisher approximation.

Table 12: Comparison of different post-hoc methods applied to the fine-tuned RoBERTa-base across six common GLUE and SuperGLUE tasks. Results are evaluated at the early stopping point of 5000 gradient steps. LA and LLLA using diagonal Fisher approximation.

Methods	Metrics	WNLI	RTE	MRPC	WiC	CoLA	BoolQ
ACC \uparrow	MAP	22.5 _{4.6}	72.8 _{2.2}	87.1 _{0.7}	64.5 _{2.0}	82.4 _{0.6}	77.2 _{0.4}
	MC Drop	19.7 _{3.0}	74.0 _{1.6}	88.2 _{0.0}	68.5 _{0.8}	82.7 _{0.1}	76.7 _{0.6}
	Ckpt Ens	21.6 _{5.8}	73.9 _{2.5}	87.7 _{0.6}	64.3 _{1.9}	81.8 _{0.5}	76.6 _{0.2}
	LLLA	22.5 _{4.6}	72.8 _{2.2}	87.1 _{0.7}	64.5 _{2.0}	82.4 _{0.6}	77.2 _{0.4}
	LA	22.5 _{4.6}	72.7 _{2.1}	87.1 _{0.5}	64.7 _{1.9}	82.3 _{0.7}	77.2 _{0.5}
ECE \downarrow	MAP	66.7 _{3.2}	20.9 _{2.1}	8.3 _{0.1}	18.4 _{2.2}	8.6 _{0.5}	5.3 _{0.2}
	MC Drop	65.6 _{1.8}	12.7 _{1.6}	4.8 _{1.6}	8.9 _{0.6}	1.2 _{0.1}	1.5 _{0.4}
	Ckpt Ens	44.6 _{3.9}	8.5 _{1.3}	17.9 _{0.1}	5.2 _{1.0}	13.6 _{0.3}	12.5 _{0.3}
	LLLA	64.0 _{3.7}	19.6 _{2.1}	7.5 _{0.2}	17.4 _{2.3}	7.7 _{0.5}	4.6 _{0.1}
	LA	33.7 _{4.0}	14.1 _{2.8}	25.0 _{0.5}	6.0 _{1.9}	13.1 _{1.4}	16.4 _{0.5}
NLL \downarrow	MAP	3.10 _{0.09}	1.05 _{0.09}	0.43 _{0.01}	0.79 _{0.03}	0.46 _{0.02}	0.50 _{0.00}
	MC Drop	2.07 _{0.08}	0.63 _{0.01}	0.30 _{0.01}	0.65 _{0.01}	0.38 _{0.00}	0.49 _{0.00}
	Ckpt Ens	0.96 _{0.02}	0.56 _{0.01}	0.44 _{0.00}	0.63 _{0.01}	0.49 _{0.00}	0.55 _{0.00}
	LLLA	2.59 _{0.10}	0.93 _{0.09}	0.39 _{0.01}	0.77 _{0.03}	0.44 _{0.02}	0.50 _{0.00}
	LA	0.80 _{0.00}	0.60 _{0.01}	0.51 _{0.01}	0.63 _{0.00}	0.46 _{0.01}	0.57 _{0.01}

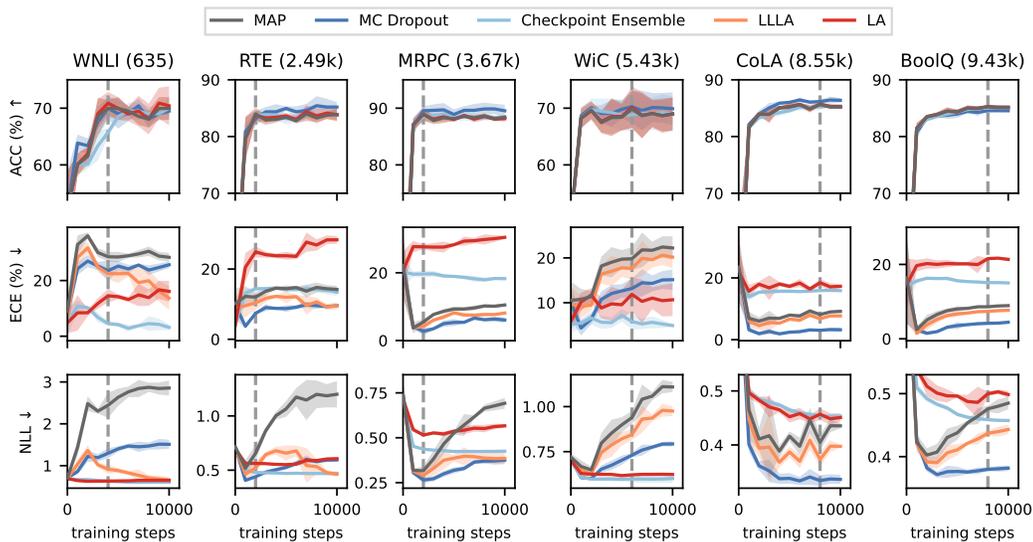


Figure 10: Fine-tuning of RoBERTa-large across six GLUE and SuperGLUE tasks (presented column-wise, with number of training examples in brackets), evaluated on the test set every 1000 gradient steps. The vertical dashed line gives the number of training steps with optimal MAP performance. LA and LLLA using diagonal Fisher approximation.

Table 13: Comparison of different post-hoc methods applied to the fine-tuned RoBERTa-large across six common GLUE and SuperGLUE tasks. Results are evaluated at the early stopping point of 5000 gradient steps. LA and LLLA using diagonal Fisher approximation.

Methods	Metrics	WNLI	RTE	MRPC	WiC	CoLA	BoolQ
ACC \uparrow	MAP	70.0 _{1.8}	83.4 _{0.3}	88.2 _{0.3}	68.4 _{2.0}	85.2 _{0.5}	84.4 _{0.4}
	MC Drop	69.0 _{2.3}	85.0 _{0.3}	88.9 _{0.8}	69.8 _{1.4}	86.0 _{0.1}	84.0 _{0.6}
	Ckpt Ens	68.5 _{0.7}	83.2 _{0.9}	88.2 _{0.5}	68.8 _{2.4}	84.6 _{0.4}	84.3 _{0.5}
	LLLA	70.0 _{1.8}	83.5 _{0.3}	88.2 _{0.3}	68.4 _{2.0}	85.2 _{0.5}	84.4 _{0.4}
	LA	70.0 _{1.8}	83.6 _{0.5}	88.0 _{0.3}	68.5 _{2.2}	85.3 _{0.6}	84.4 _{0.4}
ECE \downarrow	MAP	28.5 _{1.6}	14.4 _{0.8}	9.3 _{0.6}	19.7 _{1.8}	8.0 _{0.4}	8.0 _{0.6}
	MC Drop	25.2 _{1.6}	8.8_{0.9}	5.3_{1.0}	12.6 _{1.3}	2.4_{0.3}	3.6_{0.9}
	Ckpt Ens	4.2_{0.7}	14.2 _{0.5}	18.9 _{0.2}	7.2_{1.2}	15.7 _{0.5}	15.6 _{0.3}
	LLLA	22.4 _{3.0}	12.0 _{1.5}	8.2 _{0.3}	17.8 _{1.8}	6.7 _{0.4}	6.7 _{0.7}
	LA	14.2 _{0.8}	23.8 _{1.0}	28.3 _{0.2}	9.5 _{2.8}	17.3 _{0.6}	20.2 _{0.5}
NLL \downarrow	MAP	2.64 _{0.18}	1.06 _{0.06}	0.53 _{0.01}	0.90 _{0.04}	0.42 _{0.01}	0.44 _{0.01}
	MC Drop	1.40 _{0.06}	0.53 _{0.01}	0.32_{0.02}	0.70 _{0.02}	0.34_{0.01}	0.38_{0.01}
	Ckpt Ens	0.62_{0.01}	0.47_{0.00}	0.42 _{0.00}	0.60_{0.01}	0.46 _{0.00}	0.47 _{0.00}
	LLLA	0.88 _{0.05}	0.66 _{0.06}	0.39 _{0.02}	0.82 _{0.03}	0.39 _{0.01}	0.41 _{0.01}
	LA	0.64 _{0.01}	0.55 _{0.01}	0.54 _{0.01}	0.62 _{0.00}	0.45 _{0.01}	0.50 _{0.00}

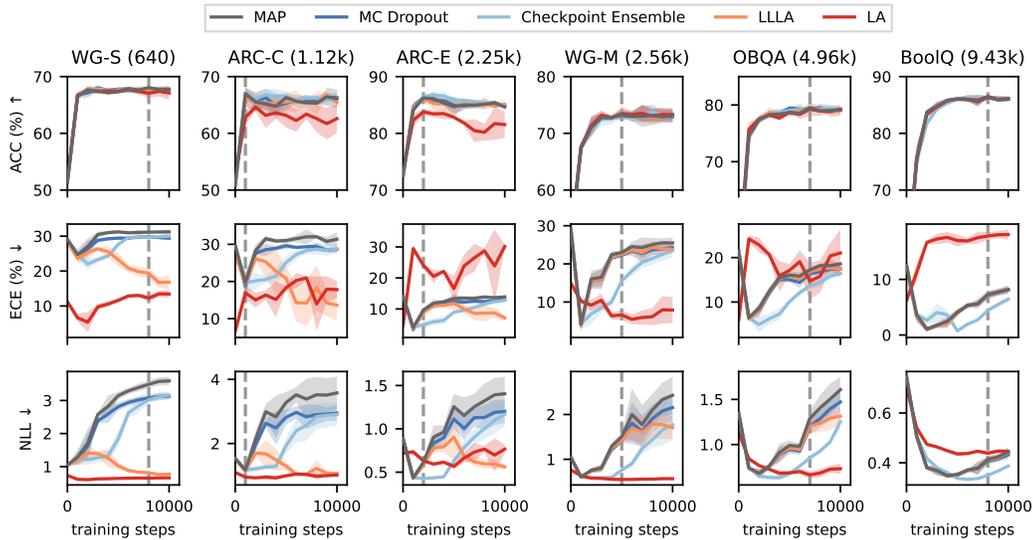


Figure 11: Fine-tuning of LLaMA2-7B across six common sense reasoning tasks (presented column-wise, with number of training examples in brackets), evaluated on the test set every 1000 gradient steps. Laplace (diagonal) prior precision is tuned using model evidence.

Table 14: Comparison of different post-hoc methods applied to the fine-tuned LLaMA2-7B across six common sense reasoning tasks. Results are evaluated at the early stopping point of 5000 gradient steps.

Methods	Metrics	WG-S	ARC-C	ARC-E	WG-M	OBQA	BoolQ
ACC ↑	MAP	67.4 _{0.3}	66.3 _{0.6}	84.7 _{1.5}	73.4 _{0.4}	78.7 _{0.4}	86.1 _{0.2}
	MC Drop	67.8 _{0.1}	65.3 _{1.0}	85.0 _{1.3}	73.2 _{0.5}	79.5 _{0.2}	86.0 _{0.3}
	Ckpt Ens	67.4 _{0.3}	66.6 _{0.8}	85.9 _{0.2}	73.5 _{0.7}	79.1 _{0.2}	86.0 _{0.3}
	LLLA	67.7 _{0.3}	65.5 _{1.4}	84.6 _{1.2}	73.6 _{0.7}	78.7 _{0.7}	86.0 _{0.3}
	LA	67.4 _{0.4}	63.2 _{1.9}	82.8 _{0.2}	73.5 _{1.1}	78.7 _{0.4}	86.1 _{0.3}
ECE ↓	MAP	31.2 _{0.3}	31.0 _{0.5}	13.4 _{1.3}	23.0 _{0.1}	16.1 _{0.6}	4.0 _{0.5}
	MC Drop	29.4 _{0.3}	29.6 _{0.8}	12.4 _{1.2}	22.2 _{0.2}	15.0 _{0.4}	4.1 _{0.4}
	Ckpt Ens	27.6 _{0.6}	25.1 _{1.6}	9.0_{0.3}	15.5 _{0.2}	10.1_{0.3}	0.7_{0.2}
	LLLA	23.3 _{1.7}	20.9 _{2.6}	11.8 _{1.9}	22.5 _{0.6}	15.9 _{0.3}	4.1 _{0.6}
	LA	11.6_{0.5}	18.3_{0.5}	16.6 _{3.7}	6.6_{1.2}	17.2 _{1.2}	17.0 _{0.9}
NLL ↓	MAP	3.15 _{0.10}	3.28 _{0.29}	1.26 _{0.13}	1.51 _{0.05}	0.99 _{0.05}	0.35 _{0.01}
	MC Drop	2.81 _{0.11}	2.82 _{0.21}	1.11 _{0.10}	1.41 _{0.03}	0.95 _{0.04}	0.35 _{0.01}
	Ckpt Ens	1.84 _{0.22}	1.93 _{0.22}	0.63 _{0.01}	0.76 _{0.01}	0.68_{0.03}	0.34_{0.00}
	LLLA	1.02 _{0.10}	1.30 _{0.11}	0.90 _{0.24}	1.42 _{0.05}	0.96 _{0.05}	0.35 _{0.01}
	LA	0.64_{0.01}	1.00_{0.04}	0.57_{0.03}	0.55_{0.01}	0.69 _{0.01}	0.44 _{0.01}

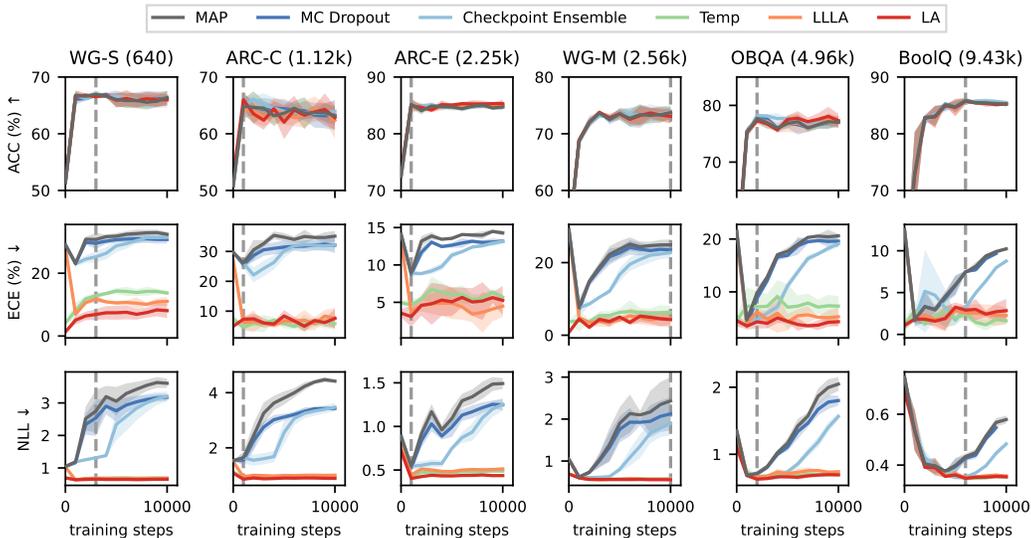


Figure 12: Fine-tuning of LLaMA2-7B across six common sense reasoning tasks (presented column-wise, with number of training examples in brackets), evaluated on the test set every 1000 gradient steps. The vertical dashed line gives the checkpoint with optimal MAP performance on a held-out validation set. Temperature scaling and Laplace (diagonal) prior precision are tuned on the validation set.

Table 15: Comparison of different post-hoc methods applied to the fine-tuned LLaMA2-7B across six common sense reasoning tasks, with a validation set split from the training set used for tuning temperature and Laplace prior precision. Results are evaluated at the best MAP performance checkpoint observed on the validation set.

Methods	Metrics	WG-S	ARC-C	ARC-E	WG-M	OBQA	BoolQ
ACC ↑	MAP	67.0 _{0.6}	64.9 _{1.1}	85.2 _{0.6}	73.7 _{0.9}	77.7 _{0.8}	85.8 _{0.4}
	MC Drop	66.7 _{0.3}	64.9 _{1.9}	85.1 _{0.5}	73.5 _{0.9}	77.7 _{0.2}	85.9 _{0.4}
	Ckpt Ens	66.7 _{0.3}	64.9 _{1.1}	85.2 _{0.6}	73.8 _{1.0}	78.2 _{0.2}	85.4 _{0.3}
	Temp	67.0 _{0.6}	64.9 _{1.1}	85.2 _{0.6}	73.7 _{0.9}	77.7 _{0.8}	85.8 _{0.4}
	LLLA	66.7 _{0.3}	64.4 _{1.0}	85.1 _{0.8}	73.1 _{1.2}	77.2 _{0.3}	85.7 _{0.4}
	LA	66.5 _{0.3}	66.0 _{1.2}	85.0 _{1.2}	73.1 _{1.1}	77.3 _{0.3}	85.7 _{0.5}
ECE ↓	MAP	30.8 _{1.8}	26.1 _{1.4}	8.9 _{0.3}	24.9 _{1.3}	9.8 _{1.0}	7.4 _{0.1}
	MC Drop	29.5 _{1.6}	25.6 _{0.7}	8.8 _{0.6}	23.5 _{1.2}	8.8 _{0.8}	7.5 _{0.1}
	Ckpt Ens	25.2 _{1.6}	26.1 _{1.4}	8.9 _{0.3}	22.8 _{1.4}	4.7 _{0.5}	3.2 _{0.5}
	Temp	12.8 _{0.9}	4.6 _{1.0}	4.7 _{0.8}	6.3 _{1.6}	7.2 _{2.6}	2.5 _{0.3}
	LLLA	11.8 _{1.1}	6.0 _{1.8}	3.7 _{0.5}	4.5 _{2.2}	6.2 _{0.5}	2.6 _{0.9}
	LA	6.9 _{1.5}	7.3 _{0.6}	3.1 _{1.2}	4.3 _{1.3}	4.3 _{1.1}	2.9 _{0.5}
NLL ↓	MAP	2.75 _{0.57}	1.64 _{0.19}	0.54 _{0.03}	2.43 _{0.50}	0.71 _{0.03}	0.43 _{0.01}
	MC Drop	2.54 _{0.49}	1.55 _{0.16}	0.52 _{0.04}	2.12 _{0.35}	0.71 _{0.04}	0.43 _{0.01}
	Ckpt Ens	1.31 _{0.04}	1.64 _{0.18}	0.54 _{0.03}	1.89 _{0.24}	0.65 _{0.02}	0.35 _{0.01}
	Temp	0.68 _{0.01}	0.90 _{0.01}	0.43 _{0.02}	0.58 _{0.01}	0.67 _{0.02}	0.35 _{0.00}
	LLLA	0.68 _{0.02}	0.95 _{0.03}	0.44 _{0.01}	0.57 _{0.01}	0.66 _{0.02}	0.35 _{0.00}
	LA	0.66 _{0.02}	0.86 _{0.03}	0.40 _{0.02}	0.55 _{0.01}	0.63 _{0.01}	0.35 _{0.01}

Subject	Task
Computer Science (CS)	college computer science
	computer security
	high school computer science
	machine learning
Engineering (Eng)	electrical engineering
Law	international law
	jurisprudence
	professional law
Health	anatomy
	clinical knowledge
	college medicine
	human aging
	nutrition
	professional medicine
	virology

Table 16: Far OOD MMLU subjects and tasks.