Minimal Ranks, Maximum Confidence: Parameter-efficient Uncertainty Quantification for LoRA

Patryk Marszałek* Klaudia Bałazy Jacek Tabor Tomasz Kuśmierczyk*†

Jagiellonian University

Abstract

Adaptation (LoRA) enables Low-Rank parameter-efficient fine-tuning of large language models by decomposing weight updates into low-rank matrices, significantly reducing storage and computational overhead. While effective, standard LoRA lacks mechanisms for uncertainty quantification, leading to overconfident and poorly calibrated models. Bayesian variants of LoRA address this limitation, but at the cost of a significantly increased number of trainable parameters, partially offsetting the original efficiency gains. Additionally, these models are harder to train and may suffer from unstable convergence. In this work, we propose a novel parameter-efficient Bayesian LoRA via subspace inference, demonstrating that effective uncertainty quantification can be achieved in very low-dimensional parameter spaces. The proposed method achieves strong performance with improved calibration and generalization while maintaining computational efficiency. Our empirical findings show that, with the appropriate projection of the weight space: (1) uncertainty can be effectively modeled in a low-dimensional space, and (2) weight covariances exhibit low ranks.

1 Introduction

LoRA (Low-Rank Adaptation) (Hu et al., 2022) reduces computational overhead by decomposing the update weights of pre-trained models into low-rank matrices, enabling efficient adaptation to downstream tasks. Minimizing the number of trainable parameters reduces memory and storage requirements, making large-scale model adaptation feasible. Reducing computational overhead speeds up training time and makes adaptation possible in resource-constrained settings.

Unlike pre-trained models, which are relatively well-calibrated (OpenAI, 2023), fine-tuned large

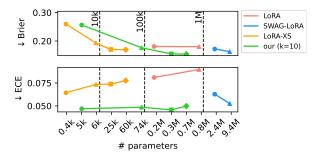


Figure 1: Performance averaged over multiple GLUE datasets (individual results in Fig. 3). Our method achieves superior calibration (ECE) and competitive predictive performance (Brier) while maintaining computational efficiency. For example, at r=8 (\clubsuit), we reduce ECE by half with only 1/10th LoRA parameters.

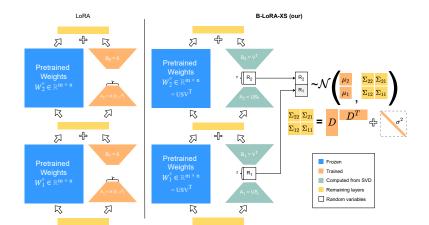
models (e.g., LLMs) often become overconfident and poorly calibrated (Jiang et al., 2021; Tian et al., 2023; Xiao et al., 2022; He et al., 2023), especially when trained on limited data. This hinders their usability for applications where uncertainty-aware decisions are performed.

Bayesian treatment is then frequently proposed to address overconfidence in neural networks (Blundell et al., 2015; Kristiadi et al., 2020; Aitchison et al., 2021; Izmailov et al., 2021). Consequently, recently proposed Bayesian variants of LoRA (Onal et al., 2024; Yang et al., 2024; Doan et al., 2025) address the aforementioned challenges by introducing uncertainty estimation directly into the fine-tuning process. During training, these models continuously adjust both the mean and covariance of fine-tuned parameters to achieve better generalization and uncertainty quantification.

Learning the posterior covariance matrix is necessary for modeling epistemic uncertainty. However, its size grows quadratically with the number of parameters, which can easily cancel out the benefits of LoRA, in addition to making learning significantly harder. Using low-rank, Kronecker-factored, or diagonal-only covariances partially alleviates

^{*}Denotes primary contributors.

[†]Correspondence: t.kusmierczyk@uj.edu.pl



	Method	r	k	# Parameters
Þ	LoRA	2	-	0.2M
Standard	LoRA	8	-	0.8M
Star	LoRA-XS	8	-	6k
	LoRA-XS	25	-	60k
	SWAG-LoRA	2	10	2.4M
₽	SWAG-LoRA	8	10	9.4M
Bayesian	SWAG-LoRA	8	5	5.5M
Bay	B-LoRA-XS	8	10	74k
	B-LoRA-XS	25	10	0.7M
	B-LoRA-XS	25	5	0.4M

Figure 2: (**left**): Weight-adaptation approaches: LoRA vs. B-LoRA-XS. As indicated by the color coding, some parameters remain frozen (*blue*), others are trained (*orange*) or obtained via SVD (*green*). (**right**): Number of trainable parameters per method. XS variants remain computationally competitive even for ranks as large as r = 25.

the problem, but as we demonstrate in Sec. 3, this comes at the cost of results quality loss. Furthermore, even at rank = 2, the number of trainable parameters is quadrupled compared to vanilla LoRA. This creates a need for an alternative approach that retains covariance modeling capacity while reducing the number of required parameters.

We propose a method that learns Bayesian posteriors for weights projected onto a low-dimensional manifold, hence maintaining parameter efficiency. The thoughtfully selected projection allows for the effective representation of the covariances between weights through covariances between representations in the lower-dimensional space. In this design, we follow the work of Bałazy et al. (2024), who recently proposed a strategy for finding such projections with SVD. We prove that they are *effective* for learning Bayesian models as well.

Operating in such a reduced parameter space significantly improves the feasibility of Bayesian inference. We show that correlations between weights can be represented very efficiently – unlike in the original weight space, we can use covariance matrices with ranks as low as k=2. Thanks to the low number of parameters, training is also more stable. Finally, the method achieves superior calibration and accuracy at low budgets (e.g., see Fig. 1).

A key contribution of our work lies in introducing Bayesian learning within a low-rank projected subspace derived from pre-trained weights. While the individual components of our method, such as SVD projections and Bayesian inference, are established techniques, their synergistic application to learn Bayesian posteriors within a compressed

subspace constitutes a meaningful conceptual innovation. Our approach enables uncertainty-aware fine-tuning with strong parameter efficiency, yielding improved uncertainty quantification with minimal computational overhead.

In the Appendix, we supplement the results presented in the paper with a discussion of related work, a detailed overview of the experimental setup, and the exact numeric values for the figures in the main text.

The source code is available online¹.

2 Method: B-LoRA-XS

LoRA fine-tunes large pre-trained models by learning low-rank weight updates ΔW instead of training the weights W directly. For a pre-trained parameter matrix $W^0 \in \mathbb{R}^{m \times n}$ that is kept fixed, LoRA learns a rank-r update $\Delta W = AB$, where $A \in \mathbb{R}^{m \times r}$ and $B \in \mathbb{R}^{r \times n}$ have far fewer parameters. The effective weight is then: $W = W^0 + \Delta W = W^0 + AB$, where only A and B are trained. LoRA is then typically applied jointly for multiple layers l, yielding a set of updates $\{\Delta W_l\}$. Then, Bayesian treatment of LoRA can improve its calibration and uncertainty quantification.

Bayesian treatment of a neural network involves finding the posterior $p(\theta \mid \mathcal{D})$ given training data \mathcal{D} . By Bayes' theorem: $p(\theta \mid \mathcal{D}) = \frac{p(\mathcal{D}|\theta)\,p(\theta)}{p(\mathcal{D})}$, where θ represents the model's parameters (i.e., weights) considered random variables. Specifically, for the Bayesian LoRA setting, θ denotes a set of the *learned model updates*, while the remaining *frozen* weights are hid-

¹Source code: https://github.com/gmum/b-lora-xs

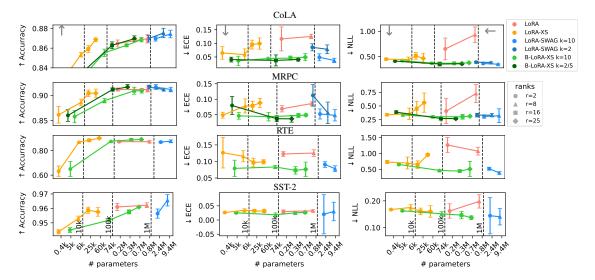


Figure 3: Median \pm s.d. accuracy (left), ECE (middle), and NLL (right) on 4 GLUE tasks (rows) vs. total parameter count for several methods and varying ranks r. B-LoRA-XS (our) achieves the accuracy and the calibration of SWAG-LoRA while using significantly fewer parameters than LoRA. See Fig. 1 for averaged results. The arrows in the figure indicate the direction of improvement ("towards better") in the standard manner. The exact numeric values underlying the plots are reported in Tables 1-4 in the Appendix.

den inside the model likelihood, given by $p(\mathcal{D} \mid \theta) = \prod_{i \in [\mathcal{D}]} p(y_i | x_i, \theta)$. The learned posterior allows Bayesian model averaging at inference as: $p(y_* \mid x_*, \mathcal{D}) = \int p(y_* \mid x_*, \theta) \, p(\theta \mid \mathcal{D}) \, d\theta \approx \frac{1}{S} \sum_{\theta \sim p(\theta \mid \mathcal{D})} p(y_* \mid x_*, \theta)$.

Bayesian LoRAs obtain the posterior for $\{\Delta W_l\}$ through the learned posterior for $\theta =$ $\bigcup_{l} \{A_{l} \cup B_{l}\}\$, where l indexes the weight updates (layers). The posterior itself is approximated either using a set of particles or a closed-form distribution. Due to its superior performance, we rely on the latter and assume $p(\theta|\mathcal{D}) \approx \mathcal{N}(\mu, \Sigma)$, where μ is the vector of means (of size equal to the number of learned parameters) and Σ is the covariance matrix, whose size grows quadratically with the total number of parameters. Notably, we aim to model cross-layer interdependencies, requiring covariance estimation also across weights in different layers $\{l\}$. This however results in an impractically large number of parameters. Consequently, we explore methods to reduce this cost by representing distributions $p(\{\Delta W_l\}|\mathcal{D})$ differently, e.g., using SVD-based projections.

In LoRA-XS (Bałazy et al., 2024), the adaptation matrices A and B are initialized using the truncated SVD of the corresponding pre-trained weight matrices W^0 . This initialization captures the most informative singular components of the original weights. Under the assumption that the fine-tuned task is similar to the original task, these projections retain the functional properties also for

downstream adaptations. LoRA-XS then freezes A and B and inserts a small trainable matrix $R \in \mathcal{R}^{r \times r}$ between them, reducing the number of trainable parameters to r^2 ($r^2 \ll (n+m) \cdot r$) per weight matrix. Then, the fine-tuning update is: $h = xW^0 + x\Delta W = xW^0 + xARB$, where $A \in \mathbb{R}^{m \times r}$ and $B \in \mathbb{R}^{r \times n}$ are low-rank matrices obtained from the truncated SVD of W^0 , specifically $A = U_r S_r$ and $B = V_r^T$.

B-LoRA-XS, proposed in this paper, leverages the frozen projections A and B for effective and efficient Bayesian learning. Its core idea is to apply the Bayesian treatment in the extremely compressed parameter space given by SVD-based projections, making Bayesian inference tractable and highly efficient. Although this specific idea is novel, it can be related to the framework of subspace Bayesian Inference.

Subspace Inference (SI) (Izmailov et al., 2020) was proposed as a remedy for the intractability of full-dimensional posteriors in modern networks. Low-dimensional affine manifold, carefully centred on a well-trained solution, already contains a rich family of high-performing weight vectors. Performing Bayesian integration restricted to that manifold restores calibrated uncertainty without revisiting the entire parameter space. In particular, given a well-trained reference point $\bar{w} \in \mathbb{R}^d$ and $K \ll d$ orthonormal basis vectors-e.g. the leading PCA directions of an SGD trajectory - the learning subspace is $\mathcal{S} := \{w = \bar{w} + Pz \mid z \in \mathbb{R}^K\}$,

 $P \in \mathbb{R}^{d \times K}$, where Bayesian parameters are the low-dimensional representations $\theta = \{z\}$.

B-LoRA-XS defines a *layer-local* manifold that leaves the pretrained backbone untouched. For each frozen weight matrix W_l^0 , an SVD $W_l^0 = U_l S_l V_l^T$ is computed once, and the top r singular directions yield fixed projectors $A_l = U_{l,r} S_{l,r}$ and $B_l = V_{l,r}^T$. Note that, whereas SI *learns* P from the training dynamics, B-LoRA-XS *projects* onto directions already favored by the pretrained backbone. Then, all task-specific variability is captured by these square adapters $R_l \in \mathbb{R}^{r \times r}$.

Vectorising and stacking these layer-projections defines $w = \bar{w} + P_B z_B$, $P_B = \operatorname{blockdiag}(B_l^T \otimes A_l)$, where $w = [\operatorname{vec}(W_1)^T, \ldots, \operatorname{vec}(W_l)^T, \ldots]^T$, $\bar{w} = [\operatorname{vec}(W_1^0)^T, \ldots, \operatorname{vec}(W_l^0)^T, \ldots]^T$, and P_B is a block-diagonal matrix with blocks $B_l^T \otimes A_l$, one for each layer l.

Thus B-LoRA-XS explores an affine subspace $\mathcal{S}_B:=\{w=w^0+P_Bz_B|z_B\in\mathbb{R}^{\sum_l r^2}\}$, whose dimension scales with r^2 per adapted layer.

In Sec. 3, we empirically demonstrate that A_l and B_l , obtained from the SVD of the pre-trained weights, are not only effective for point-wise finetuning but also enable effective uncertainty quantification for $\{\Delta W_l\}$ through modeling covariances for $\{R_l\}$. Although we never compute it explicitly, the covariance matrix for individual ΔW is expressed as $\Sigma_{\Delta W_l} = (B_l^T \otimes A_l) \Sigma_{R_l} (B_l^T \otimes A_l)^T$,, where Σ_R is the (intra-layer) covariance matrix for R and \otimes denotes the Kronecker product.

In practice, we simply learn the *joint posterior* $p(\theta = \cup_l R_l | \mathcal{D}) \approx \mathcal{N}(\mu, \Sigma)$ (only) for the inner matrices R. The covariance matrix Σ captures both inter-layer and intra-layer dependencies, allowing the model to learn complex relationships. At inference, similar to LoRA, we use samples of R along with the respective projections R and R to obtain R, as realized through samples of R0, however without ever computing it explicitly.

The parameters μ and Σ are learned efficiently using **SWAG** (Maddox et al., 2019) (though Variational Inference or the Laplace approximation could also be used). After a burn-in phase (a fixed 10 or 25 epochs) of the gradient-based optimization, the algorithm maintains $\hat{\mu}$ – a running average of θ – along with k vectors of differences $\hat{D}_{last} = \theta_{last} - \hat{\mu}$ for the last k values of θ , and a running average of θ^2 . Based on these averages, we estimate the variances $\hat{\sigma}^2$ for individual parameters and approximate the covariance as $\hat{\Sigma} \approx \frac{1}{2}(\hat{D}\cdot\hat{D}^T+diag(\hat{\sigma}^2))$, which constitutes a

rank-k approximation to the covariance matrix.

We illustrate B-LoRA-XS method in Fig. 2. Our method uses the total of $|\theta| \cdot (k+2)$ parameters, where $|\theta| = \sum_l r^2$.

3 Experiments

Setup: We performed experimental evaluation on four GLUE tasks (Wang et al., 2018) (RTE, MRPC, CoLA, and SST-2) using RoBERTa-large (Liu et al., 2019). We compare our method (B-LoRA-XS) against the standard LoRA, LoRA-XS – a parameter efficient variant, and against SWAG-LoRA (Onal et al., 2024) – a Bayesian variant. For LoRA-XS and B-LoRA-XS we considered ranks $r \in \{2, 8, 16, 25\}$ and for LoRA and SWAG-LoRA due to limited budget we were able to test $r \in \{2, 8\}$. The number of parameters (a *proxy for storage and computation costs*) as a function of ranks r and k is summarized in Fig. 2. We report accuracy (higher is better), ECE and NLL (lower is better) of median \pm s.d. across 5 runs.

Performance analysis: Fig. 3 compares accuracy, Expected Calibration Error (ECE), and Negative Log-Likelihood (NLL) against total parameter count across 4 datasets. Our main claim is that B-LoRA-XS improves overall model performance, with a particular focus on calibration metrics. Indeed, Figure 3 (middle and right) demonstrates that B-LoRA-XS consistently yields lower ECE and NLL compared to standard LoRA across all parameter scales. Regarding accuracy (Figure 3: left), while standard LoRA shows marginally better results for a few configurations at moderate parameter scales, the majority of configurations show B-LoRA-XS matching or exceeding the accuracy of standard LoRA. More importantly, in no setting does standard LoRA significantly outperform B-LoRA-XS in terms of calibration, which is a primary focus of our work. Bayesian variants, including B-LoRA-XS and SWAG-LoRA, generally outperform their non-Bayesian counterparts in ECE and NLL. However, our model achieves these strong calibration results with 5-15 times fewer parameters than SWAG-LoRA. Moreover, while SWAG-LoRA sometimes performs well, its results vary significantly between runs. In contrast, B-LoRA-XS exhibits stable and consistent convergence. Finally, as results for MRPC and CoLA suggest, its performance remains robust across different values of k, whereas SWAG-LoRA's ECE deteriorates significantly at k=2.

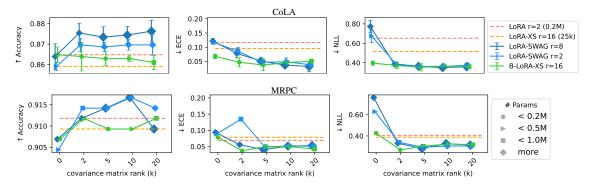


Figure 4: Impact of the posterior covariance matrix rank (k = 0 indicates the case with no off-diagonal elements) for CoLA (top) and MRPC (bottom). For brevity, confidence bars (\pm standard deviation) are omitted for MRPC. The colored lines represent non-Bayesian baselines (e.g., standard LoRA or LoRA-XS at a given rank r). The exact numeric values underlying the plots are reported in Tables 5 and 6 in the Appendix.

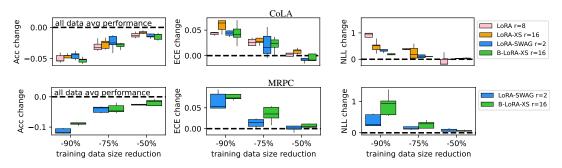


Figure 5: Accuracy, ECE and NLL change as the training set is progressively reduced (e.g. -90% means using only 10% of the data for training). The dashed line marks the model's performance when trained on the full dataset.

Covariance matrix rank analysis: Figure 4 compares the sensitivity of the Bayesian LoRA variants to changes in covariance matrix rank k. Markers indicate model sizes (e.g., SWAG-LoRA \gg B-LoRA-XS). As expected, SWAG-LoRA deteriorates proportionally as rank decreases. On the other hand, B-LoRA-XS maintains its performance across a wide range of k. Significant degradation occurs only when off-diagonal covariance values are entirely ignored (i.e., at k=0). Notably, B-LoRA-XS achieves the best calibration at low ranks, particularly at k=2 or k=5. This demonstrates that the SVD-based projection effectively disentangles parameters, enabling low-dimensional uncertainty modeling.

Data size reduction analysis: Figure 5 compares how accuracy, ECE, and NLL degrade when training data is subsampled. All methods predictably lose accuracy as data size decreases, with little difference between the various LoRA-based approaches. We conclude that Bayesian learning does not improve robustness in this case. However, we note variations across datasets in terms of accuracy. For example, in MRPC, the decline is more pronounced, likely due to the dataset smaller size.

4 Conclusion

B-LoRA-XS addresses the lack of uncertainty quantification in LoRA fine-tuning while maintaining parameter efficiency. It utilizes truncated SVD to project model updates into a lower-dimensional space and leverages the Bayesian framework to enhance uncertainty estimation.

Our method's primary strength lies in its calibration capabilities; it consistently achieves lower expected calibration error and negative log-likelihood compared to standard LoRA and LoRA-XS across various parameter scales. While standard LoRA may exhibit marginally better accuracy in a few specific configurations, B-LoRA-XS generally matches or exceeds its accuracy in most settings, and critically, always provides superior or equal calibration. Compared to the Bayesian LoRA baseline, B-LoRA-XS matches or surpasses its accuracy and calibration performance while using significantly fewer parameters, exhibiting greater training stability, and relying on simpler, lower-rank covariance representations.

Limitations

While B-LoRA-XS demonstrates promising results in parameter-efficient uncertainty quantification, several limitations should be acknowledged. First, the effectiveness of B-LoRA-XS inherently depends on the quality of the initial SVD projection derived from pre-trained weights (as in LoRA-XS). If the principal components of the pre-trained model are not well-aligned with the requirements of a significantly different downstream task, the performance might be suboptimal. Second, our method employs SWAG with a low-rank approximation for the covariance matrix. While efficient, this is one specific approach to approximate Bayesian inference. Other techniques (e.g., more sophisticated variational inference methods or different posterior approximations) might yield different trade-offs between performance, calibration, and computational cost, and were not explored in this work. Third, although B-LoRA-XS significantly reduces the number of trainable parameters for Bayesian adaptation, the inference process still requires multiple forward passes for sampling, which increases computational cost compared to non-Bayesian LoRA or LoRA-XS. Fourth, our empirical validation is conducted on GLUE classification tasks using RoBERTa-Large. The generalizability of B-LoRA-XS's benefits to other model architectures, much larger model scales, or different task types (such as text generation or more complex reasoning tasks) warrants further investigation. Finally, the optimal choice of LoRA rank r and SWAG covariance rank k might require careful tuning for different datasets and models, potentially adding to the practical overhead of applying the method effectively.

Acknowledgments

This research is part of the project No. 2022/45/P/ST6/02969 co-funded by the National Science Centre and the European Union Framework Programme for Research and Innovation Horizon 2020 under the Marie Skłodowska-Curie grant agreement No. 945339. For the purpose of Open Access, the author has applied a CC-BY public copyright licence to any Author Accepted Manuscript (AAM) version arising from this submission.



We gratefully acknowledge Polish high-performance computing infrastructure PLGrid (HPC Center: ACK Cyfronet AGH) for providing computer facilities and support within computational grant no. PLG/2024/017893

The work of Klaudia Bałazy was supported by the National Centre of Science (Poland) Grant No. 2020/39/D/ST6/ 01332. Klaudia Bałazy is affiliated with Doctoral School of Exact and Natural Sciences at the Jagiellonian University.

References

- Laurence Aitchison, Adam Yang, and Sebastian W Ober. 2021. Deep kernel processes. In *International Conference on Machine Learning*, pages 130–140. PMLR.
- Klaudia Bałazy, Mohammadreza Banaei, Karl Aberer, and Jacek Tabor. 2024. LoRA-XS: Low-rank adaptation with extremely small number of parameters. *arXiv preprint arXiv:2405.17604*.
- Charles Blundell, Julien Cornebise, Koray Kavukcuoglu, and Daan Wierstra. 2015. Weight uncertainty in neural network. In *International Conference on Machine Learning*, pages 1613–1622. PMLR.
- Tim Dettmers, Artidoro Pagnoni, Ari Holtzman, and Luke Zettlemoyer. 2024. Qlora: Efficient finetuning of quantized llms. *Advances in Neural Information Processing Systems*, 36.
- Bao Gia Doan, Afshar Shamsi, Xiao-Yu Guo, Arash Mohammadi, Hamid Alinejad-Rokny, Dino Sejdinovic, Damien Teney, Damith C. Ranasinghe, and Ehsan Abbasnejad. 2025. Bayesian low-rank learning (bella): A practical approach to bayesian neural networks. *Proceedings of the AAAI Conference on Artificial Intelligence*, 39(15):16298–16307.
- Demi Guo, Alexander Rush, and Yoon Kim. 2021. Parameter-efficient transfer learning with diff pruning. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4884–4896, Online. Association for Computational Linguistics.
- Guande He, Jianfei Chen, and Jun Zhu. 2023. Preserving pre-trained features helps calibrate fine-tuned language models. In *International Conference on Learning Representations (ICLR)*.
- Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. 2019. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR.
- Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. 2022. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3.
- Pavel Izmailov, Wesley J Maddox, Polina Kirichenko, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. 2020. Subspace inference for bayesian deep learning. In *Uncertainty in Artificial Intelligence*, pages 1169–1179. PMLR.
- Pavel Izmailov, Sharad Vikram, Matthew D Hoffman, and Andrew Gordon Gordon Wilson. 2021. What are bayesian neural network posteriors really like?

- In *International Conference on Machine Learning*, pages 4629–4640. PMLR.
- Zhengbao Jiang, Jun Araki, Haibo Ding, and Graham Neubig. 2021. How can we know when language models know? on the calibration of language models for question answering. *Transactions of the Association for Computational Linguistics*, 9:962–977.
- Dawid J. Kopiczko, Tijmen Blankevoort, and Yuki M. Asano. 2024. Vera: Vector-based random matrix adaptation. In *International Conference on Learning Representations (ICLR)*.
- Agustinus Kristiadi, Matthias Hein, and Philipp Hennig. 2020. Being bayesian, even just a bit, fixes overconfidence in relu networks. In *Proceedings of the 37th International Conference on Machine Learning*.
- Brian Lester, Rami Al-Rfou, and Noah Constant. 2021. The power of scale for parameter-efficient prompt tuning. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 3045–3059, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Xiang Lisa Li and Percy Liang. 2021. Prefix-tuning: Optimizing continuous prompts for generation. In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4582–4597, Online. Association for Computational Linguistics.
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. Roberta: A robustly optimized bert pretraining approach. *Preprint*, arXiv:1907.11692.
- Ilya Loshchilov and Frank Hutter. 2019. Decoupled weight decay regularization. *Preprint* arXiv:1711.05101.
- Wesley J Maddox, Pavel Izmailov, Timur Garipov, Dmitry P Vetrov, and Andrew Gordon Wilson. 2019. A simple baseline for bayesian uncertainty in deep learning. In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Cristian Meo, Ksenia Sycheva, Anirudh Goyal, and Justin Dauwels. 2024. Bayesian-loRA: LoRA based parameter efficient fine-tuning using optimal quantization levels and rank values trough differentiable bayesian gates. In 2nd Workshop on Advancing Neural Network Training: Computational Efficiency, Scalability, and Resource Optimization (WANT@ICML 2024).
- Emre Onal, Klemens Flöge, Emma Caldwell, Arsen Sheverdin, and Vincent Fortuin. 2024. Gaussian stochastic weight averaging for bayesian low-rank adaptation of large language models. In Sixth Symposium on Advances in Approximate Bayesian Inference Non Archival Track.

- OpenAI. 2023. GPT-4 technical report. *Preprint*, arXiv:2303.08774.
- Katherine Tian, Eric Mitchell, Allan Zhou, Archit Sharma, Rafael Rafailov, Huaxiu Yao, Chelsea Finn, and Christopher Manning. 2023. Just ask for calibration: Strategies for eliciting calibrated confidence scores from language models fine-tuned with human feedback. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing*, pages 5433–5442, Singapore. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. In *Advances in Neural Information Processing Systems*, volume 30. Curran Associates, Inc.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. GLUE: A multi-task benchmark and analysis platform for natural language understanding. In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- Yibin Wang, Haizhou Shi, Ligong Han, Dimitris N. Metaxas, and Hao Wang. 2024. BLob: Bayesian low-rank adaptation by backpropagation for large language models. In *The Thirty-eighth Annual Conference on Neural Information Processing Systems*.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. Transformers: State-of-the-art natural language processing. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations, pages 38–45, Online. Association for Computational Linguistics.
- Yuxin Xiao, Paul Pu Liang, Umang Bhatt, Willie Neiswanger, Ruslan Salakhutdinov, and Louis-Philippe Morency. 2022. Uncertainty quantification with pre-trained language models: A large-scale empirical analysis. In *EMNLP*.
- Adam Yang, Maxime Robeyns, Xi Wang, and Laurence Aitchison. 2024. Bayesian low-rank adaptation for large language models. In *International Conference on Representation Learning*, volume 2024, pages 1812–1842.
- Qingru Zhang, Minshuo Chen, Alexander Bukharin, Pengcheng He, Yu Cheng, Weizhu Chen, and Tuo Zhao. 2023. Adaptive budget allocation for parameter-efficient fine-tuning. In *The Eleventh International Conference on Learning Representations*.

A Related Work

PEFT: As large language models continue to grow, parameter-efficient fine-tuning (PEFT) has become a popular approach to reducing computational and storage costs. Among various methods (Houlsby et al., 2019; Guo et al., 2021; Li and Liang, 2021; Lester et al., 2021), LoRA (Hu et al., 2022) has emerged as one of the most widely used. Building on its success, several approaches have been proposed to enhance different aspects of PEFT (Kopiczko et al., 2024; Zhang et al., 2023; Dettmers et al., 2024). One such method, LoRA-XS (Bałazy et al., 2024), further optimizes parameter efficiency by enabling flexible control over the number of trainable parameters per adaptation module. B-LoRA-XS reuses the idea of SVD-based projections to reduce the parameter space dimensionality.

Bayesian LoRAs: Standard LoRA (Hu et al., 2022) does not account for uncertainty, making fine-tuned models susceptible to miscalibration. Then, Bayesian LoRA approaches integrate Bayesian inference techniques into LoRA to improve uncertainty estimation and generalization.

Several Bayesian LoRA methods have been proposed, each employing different Bayesian techniques to address these challenges. SWAG-LoRA (Onal et al., 2024) combines Stochastic Weight Averaging-Gaussian (SWAG) with LoRA to enable approximate Bayesian inference, significantly improving model calibration and reducing overconfidence. Laplace-LoRA (Yang et al., 2024) applies a Laplace approximation to the posterior over LoRA parameters. Bella (Doan et al., 2025) introduces an approach that reduces the cost of Bayesian deep ensembles by applying multiple low-rank perturbations to a pre-trained model. BLoB (Bayesian Low-Rank Adaptation by Backpropagation) (Wang et al., 2024) jointly learns both the mean and covariance of model parameters throughout the fine-tuning process using Variational Inference. B-LoRA (Meo et al., 2024) introduces a Bayesian perspective to both quantization and rank selection by using a prior distribution over these hyperparameters, optimizing model efficiency and reducing bit operations.

The key challenge lies in balancing uncertainty modeling with parameter efficiency, as Bayesian inference typically increases both the number of trainable parameters and computational cost. Despite their advantages, Bayesian LoRA methods face challenges related to increased parameter count and computational cost. One major issue is the higher storage and memory requirements, as Bayesian methods often require additional parameters to model uncertainty, particularly those involving covariance estimation, such as SWAG-LoRA. Scalability remains a concern for methods that explicitly model uncertainty across a large number of parameters.

B Scientific Artifacts Licenses

Listed below are the licenses for the scientific artifacts used in this research. For complete information, please use the links below and refer to the original sources.

Scientific Artifacts: RoBERTa-large (MIT), MRPC (Unknown), RTE (Unknown), CoLA (Unknown), SST-2 (Unknown), HuggingFace Transformers Library (Apache-2.0), SWAG-LoRA repository² (MIT), LoRA-XS repository³ (Unknown).

C Model Size And Budget

• RoBERTa-large: 355M parameters

• GPUs: RTX4090 and V100-SXM2-32GB, each run was performed on a single GPU

• GPU total time: \approx 63 days

D Statistics For Data

We followed the original GLUE train-validation split

• MRPC - train: 3'668, val: 408

• RTE - train: 2'490, val: 277

• CoLA - train: 8'551, val: 1043

• SST2 - train: 67'349, val: 872

E Experimental Setup Details

The study was conducted on a subset of the GLUE benchmark (Wang et al., 2018), specifically on RTE, MRPC, CoLA, and SST-2 tasks (with the original train-validation split), using RoBERTalarge (Liu et al., 2019) checkpoints from the HuggingFace Transformers library (Wolf et al., 2020). For the RTE and MRPC tasks, we followed LoRA-XS and initialized LoRA-XS modules with

²https://github.com/fortuinlab/swag-lora
3https://github.com/MohammadrezaBanaei/
LoRA-XS

weights fine-tuned on the MNLI task. We integrated B-LoRA-XS/LoRA-XS modules into the Query, Value, Attention Output, and Output Fully Connected weight matrices in all transformer layers (Vaswani et al., 2017), whereas due to budget limits, standard LoRA and SWAG-LoRA modules were added only to the Query and Value matrices. Note this is sufficient for SWAG-LoRA to achieve its best performance.

For each dataset, for the burn-in stage of training, we adopted hyperparameters from the LoRA-XS paper. These include: learning rate, batch size, AdamW optimizer (Loshchilov and Hutter, 2019), linear scheduler with warm-up, dropout, and the LoRA scaling factor α . For standard LoRA we followed the same setup, except for the learning rate, which was determined through grid search. Similarly, the SWAG starting epoch (e.g. 10 or 25) was selected through grid search. Based on the findings from SWAG-LoRA, we used a constant learning rate scheduler (SWALR) with warm-up. The SWAG learning rate was set to the maximum learning rate from the first (burn-in) phase of training. Unless stated otherwise, we used a low-rank covariance matrix approximation with the rank k = 10. In all our experiments, SWAG estimation was applied exclusively to the LoRA modules, and SWAG predictions were consistently obtained with S=15model samples.

F Numeric Results

Tables 1-6 present exact numeric values for the plots presented in Figures 3 and 4.

G Acknowledgments

We acknowledge the use of ChatGPT for grammar checking and generation of the initial version of the plotting code.

			Ac	curacy		NLL		
			median	s.d.	median	s.d.	median	s.d.
method	r	k						
LoRA	2	10	0.865	0.004	0.116	0.044	0.652	0.215
	8	10	0.870	0.004	0.125	0.008	0.933	0.154
LoRA-SWAG	2	2	0.870	0.003	0.086	0.013	0.390	0.018
		10	0.870	0.002	0.049	0.013	0.365	0.015
	8	2	0.875	0.005	0.078	0.016	0.384	0.012
		10	0.874	0.004	0.037	0.008	0.351	0.009
LoRA-XS	2	10	0.822	0.009	0.065	0.023	0.451	0.022
	8	10	0.853	0.002	0.059	0.021	0.440	0.052
	16	10	0.859	0.003	0.096	0.016	0.516	0.067
	25	10	0.869	0.002	0.099	0.021	0.465	0.102
B-LoRA-XS	2	5	0.822	0.002	0.040	0.009	0.412	0.003
		10	0.822	0.005	0.036	0.017	0.422	0.016
	8	10	0.855	0.004	0.044	0.005	0.372	0.018
	16	5	0.863	0.003	0.038	0.020	0.354	0.037
		10	0.863	0.001	0.046	0.007	0.367	0.006
	25	5	0.870	0.003	0.041	0.006	0.360	0.021
		10	0.869	0.002	0.049	0.013	0.378	0.016

Table 1: Numeric values for CoLA dataset.

			Ac	curacy		NLL		
			median	s.d.	median	s.d.	median	s.d.
method	r	k						
LoRA	2	10	0.912	0.003	0.069	0.010	0.406	0.230
	8	10	0.912	0.005	0.086	0.006	0.727	0.165
LoRA-SWAG	2	2	0.917	0.004	0.112	0.035	0.332	0.034
		10	0.917	0.005	0.052	0.016	0.306	0.031
	8	2	0.912	0.003	0.056	0.035	0.331	0.031
		10	0.912	0.004	0.048	0.018	0.321	0.127
LoRA-XS	2	10	0.861	0.017	0.048	0.010	0.338	0.022
	8	10	0.886	0.007	0.078	0.023	0.355	0.105
	16	10	0.904	0.006	0.079	0.015	0.450	0.135
	25	10	0.904	0.008	0.088	0.015	0.560	0.176
B-LoRA-XS	2	2	0.860	0.012	0.080	0.029	0.386	0.027
		10	0.858	0.012	0.046	0.011	0.336	0.025
	8	10	0.890	0.004	0.043	0.014	0.304	0.023
	16	2	0.912	0.003	0.037	0.010	0.270	0.030
		10	0.909	0.007	0.047	0.007	0.301	0.044
	25	2	0.917	0.005	0.036	0.011	0.268	0.020
		10	0.909	0.005	0.049	0.004	0.312	0.013

Table 2: Numeric values for MPRC dataset.

			Ac	curacy	ECE			NLL
			median	s.d.	median	s.d.	median	s.d.
method	r	k						
LoRA	2	10	0.874	0.008	0.123	0.007	1.264	0.239
	8	10	0.874	0.010	0.125	0.010	1.072	0.123
LoRA-SWAG	2	10	0.870	0.007	0.091	0.009	0.518	0.046
	8	10	0.877	0.011	0.078	0.009	0.388	0.039
LoRA-XS	2	10	0.632	0.043	0.126	0.046	0.730	0.052
	8	10	0.870	0.005	0.116	0.017	0.698	0.188
	16	10	0.884	0.007	0.097	0.013	0.644	0.123
	25	10	0.902	0.005	0.099	0.006	0.957	0.045
B-LoRA-XS	2	10	0.650	0.062	0.079	0.025	0.652	0.024
	8	10	0.877	0.003	0.083	0.004	0.465	0.029
	16	10	0.888	0.007	0.073	0.011	0.446	0.030
	25	10	0.892	0.005	0.076	0.022	0.510	0.239

Table 3: Numeric values for RTE dataset.

			Ac	curacy	ECE			NLL
			median	s.d.	median	s.d.	median	s.d.
method	r	k						
LoRA	2	10	0.961	0.003	0.030	0.006	0.162	0.027
	8	10	0.962	0.002	0.032	0.004	0.198	0.025
LoRA-SWAG	2	10	0.956	0.003	0.020	0.069	0.145	0.068
	8	10	0.966	0.004	0.030	0.033	0.141	0.031
LoRA-XS	2	10	0.944	0.001	0.026	0.002	0.168	0.003
	8	10	0.953	0.002	0.034	0.005	0.175	0.011
	16	10	0.959	0.002	0.032	0.003	0.161	0.012
	25	10	0.958	0.003	0.032	0.005	0.160	0.021
B-LoRA-XS	2	10	0.945	0.002	0.025	0.003	0.163	0.003
	8	10	0.952	0.001	0.019	0.006	0.152	0.008
	16	10	0.958	0.001	0.025	0.006	0.147	0.014
	25	10	0.961	0.000	0.027	0.005	0.137	0.007

Table 4: Numeric values for SST-2 dataset.

			Ac	curacy		NLL		
			median	s.d.	median	s.d.	median	s.d.
method	r	k						
LoRA	2	-	0.865	0.004	0.116	0.047	0.652	0.228
LoRA-SWAG	2	0	0.859	0.002	0.117	0.007	0.675	0.067
		2	0.870	0.003	0.086	0.013	0.390	0.018
		5	0.869	0.002	0.047	0.012	0.362	0.010
		10	0.870	0.003	0.049	0.014	0.365	0.016
		20	0.870	0.005	0.039	0.003	0.373	0.018
	8	0	0.864	0.006	0.122	0.008	0.773	0.064
		2	0.875	0.005	0.078	0.016	0.384	0.012
		5	0.873	0.005	0.051	0.008	0.364	0.013
		10	0.874	0.004	0.037	0.009	0.351	0.010
		20	0.876	0.005	0.032	0.019	0.360	0.034
LoRA-XS	16	-	0.859	0.004	0.096	0.017	0.516	0.071
B-LoRA-XS	16	0	0.865	0.004	0.068	0.008	0.396	0.021
		2	0.864	0.007	0.047	0.016	0.372	0.031
		5	0.863	0.003	0.038	0.020	0.354	0.037
		10	0.863	0.001	0.046	0.008	0.367	0.006
		20	0.861	0.003	0.051	0.005	0.360	0.015

Table 5: Covariance matrix rank k analysis for CoLA dataset.

			Ac	ccuracy		NLL		
			median	s.d.	median	s.d.	median	s.d.
method	r	k						
LoRA	2	-	0.912	0.003	0.069	0.011	0.406	0.244
LoRA-SWAG	2	0	0.904	0.003	0.089	0.003	0.628	0.077
		2	0.914	0.004	0.135	0.045	0.340	0.040
		5	0.914	0.007	0.048	0.005	0.294	0.028
		10	0.917	0.005	0.052	0.017	0.306	0.033
		20	0.914	0.005	0.051	0.015	0.306	0.086
	8	0	0.907	0.005	0.094	0.005	0.759	0.123
		2	0.912	0.003	0.056	0.035	0.331	0.031
		5	0.914	0.003	0.040	0.019	0.283	0.103
		10	0.917	0.004	0.051	0.017	0.328	0.128
		20	0.909	0.006	0.053	0.018	0.314	0.177
LoRA-XS	16	-	0.909	0.007	0.079	0.011	0.388	0.095
B-LoRA-XS	16	0	0.907	0.006	0.078	0.009	0.426	0.037
		2	0.912	0.003	0.037	0.010	0.270	0.030
		5	0.909	0.009	0.051	0.010	0.304	0.032
		10	0.909	0.004	0.050	0.008	0.325	0.035
		20	0.912	0.004	0.042	0.010	0.318	0.014

Table 6: Covariance matrix rank \boldsymbol{k} analysis for MRPC dataset.