

---

# Variational Low-Rank Adaptation Using IVON

---

Anonymous Author(s)

Affiliation

Address

email

## Abstract

1 We use variational learning to improve accuracy and calibration of low-rank adaptation (LoRA) finetuning in large language models. Specifically, we employ the  
2 Improved Variational Online Newton (IVON) optimizer, which is a drop-in replacement of AdamW but significantly improves the performance with negligible  
3 overhead. We test our method by finetuning a Llama 2 model with 7 billion parameters on a range of commonsense reasoning datasets. Compared to AdamW  
4 finetuning, IVON improves accuracy by 2.8% and ECE by 4.6% on average. Our  
5 work provides more evidence for the effectiveness of variational learning in large  
6 language models. A link to the code will be provided in the final paper.  
7  
8  
9

## 10 1 Introduction

11 Large Language Models (LLMs) exhibit impressive capabilities across a wide range of natural  
12 language generation and understanding tasks but adapting them to new data can be hard, because  
13 lots of compute and memory is required due to their large number of parameters. While techniques  
14 like Low-Rank Adaptation (LoRA) [9] can alleviate this and enable finetuning on a small compute  
15 budget, models trained with LoRA still tend to be badly calibrated [23]. Bayesian methods promise  
16 to alleviate this, and in fact several Bayesian adaptations of LoRA have been proposed, such as  
17 Laplace-LoRA [23], LoRA ensembles [21], or SWAG-LoRA [16]. However, they all require either  
18 additional postprocessing steps or multiple training runs.

19 We present a method that uses variational learning [7, 8] to improve both the accuracy and calibration  
20 of LoRA-trained models without training overhead or additional postprocessing steps. The proposed  
21 method is called IVON-LoRA and we use the IVON optimizer [18] to estimate a diagonal Gaussian  
22 distribution over the low-rank factors that are added in LoRA. Learning the diagonal Gaussian only  
23 induces a negligible overhead (around  $\approx 1\%$ , as discussed in Sec. 3) in terms of training speed when  
24 compared to AdamW but provides a posterior from which diverse models can be sampled. We find  
25 that IVON-LoRA improves both calibration and accuracy and averaging the predictions of multiple  
26 sampled models further significantly improves calibration.

27 Several recent works consider related approaches to improve language model finetuning. Following a  
28 PAC-Bayesian framework, Liu et al. [11] proposes to finetune the full model using perturbed gradient  
29 descent. Chen and Garner [2] uses variational learning to estimate parameter importance in adaptive  
30 LoRA [24]. However, neither of them has been shown to work for recent billion-scale LLMs. Similar  
31 to Liu et al. [11], Zhelnin et al. [25] shows that Gaussian noise injection can improve instruction  
32 tuning of LLMs. Different from our work, they finetune on a significantly larger instruction dataset,  
33 which is more resilient to bad calibration and overfitting.

34 We show the effectiveness of our method by finetuning a Llama 2 model with 7 billion parameters on  
35 a range of commonsense reasoning tasks. We compare our method to both standard LoRA training  
36 with AdamW, and various Bayesian variants of LoRA [16, 23]. Our results show that IVON-LoRA  
37 outperforms AdamW significantly in terms of both accuracy and calibration on all tasks. Our method

38 gives better accuracy than the other Bayesian adaptations of LoRA and provides competitive results  
 39 in calibration metrics but without requiring training overhead or additional postprocessing.  
 40 Altogether, IVON-LoRA is easy to implement and can be used as a plug-in replacement for LoRA  
 41 training with AdamW that enables better accuracy and calibration without overhead. Our work further  
 42 shows the effectiveness of variational learning for LLMs.

## 43 2 Variational low-Rank adaptation using IVON

44 The parameters of LLMs contain weight matrices  $\mathbf{W} = \{\mathbf{W}^{(1)}, \dots, \mathbf{W}^{(D)}\}$ , where  $\mathbf{W}^{(i)} \in \mathbb{R}^{d_i \times k_i}$ .  
 45 Instead of directly finetuning the weight matrices  $\mathbf{W}$ , low-rank adaptation (LoRA) [9] learns low-rank  
 46 increments  $\Delta\mathbf{W} = \{\Delta\mathbf{W}^{(1)}, \dots, \Delta\mathbf{W}^{(D)}\}$ . Each increment  $\Delta\mathbf{W}^{(i)} = \mathbf{B}^{(i)}\mathbf{A}^{(i)}$  is parametrized  
 47 by two smaller matrices  $\mathbf{B}^{(i)} \in \mathbb{R}^{d_i \times r}$ ,  $\mathbf{A}^{(i)} \in \mathbb{R}^{r \times k_i}$ .  $r$  is a hyperparameter that determines the  
 48 rank of the increments  $\Delta\mathbf{W}$ , and is smaller than  $d_i$  and  $k_i$ . LoRA initializes all  $\mathbf{A}^{(i)}$  as random  
 49 Gaussian matrices and  $\mathbf{B}^{(i)}$  to zero. Given a set of pretrained weight matrices  $\mathbf{W}_0$ , LoRA optimizes  
 50 the following loss function,

$$\mathcal{L}(\Delta\mathbf{W}) = \frac{1}{N} \sum_{n=1}^N \ell_n(\mathbf{W}_0 + (\alpha/r)\Delta\mathbf{W}), \quad (1)$$

51 where  $\alpha > 0$  is an additional tuning hyperparameter,  $\alpha/r$  makes hyperparameter choice consistent  
 52 when varying the rank  $r$  and  $\ell_n$  is the loss for data example  $n$ .

53 We propose a variational learning approach and search for mean-field Gaussian posterior distributions  
 54 over the low-rank factors  $\mathbf{A}$  and  $\mathbf{B}$ . Denoting  $\boldsymbol{\theta} = (\mathbf{A}, \mathbf{B})$ , we assume that  $q(\boldsymbol{\theta}) \sim \mathcal{N}(\mathbf{m}, \text{diag}(\mathbf{v}))$ ,  
 55 meaning that for every scalar entry of  $\mathbf{A}$  and  $\mathbf{B}$  we learn a scalar mean and variance. This amounts to  
 56 minimizing the variational objective

$$\mathcal{L}_{\text{variational}}(q) = \lambda \mathbb{E}_{q(\boldsymbol{\theta})} [\mathcal{L}(\Delta\mathbf{W})] + \mathbb{D}_{\text{KL}}(q(\boldsymbol{\theta}) \parallel p(\boldsymbol{\theta})), \quad (2)$$

57 over the posterior mean  $\mathbf{m}$  and variance  $\mathbf{v}$ . Here the prior  $p(\boldsymbol{\theta}) \sim \mathcal{N}(0, v_0\mathbf{I})$  is chosen as a Gaussian  
 58 with zero mean and a constant variance  $v_0$ , and  $\lambda > 0$  is weighting-parameter. Setting  $\lambda = N$   
 59 one approximates the Bayesian posterior, and larger values target a cold posterior. The variational  
 60 objective (2) can be optimized using the IVON [18] optimizer,<sup>1</sup> which we employ to learn the  
 61 variational posterior  $q(\boldsymbol{\theta})$  over the low-rank factors  $\mathbf{A}$ ,  $\mathbf{B}$ . Despite learning a Gaussian mean-field  
 62 posterior over  $\mathbf{A}$  and  $\mathbf{B}$ , the optimizer has negligible memory overhead when compared to Adam.  
 63 This is due to the variance playing an analogous role to the second-moment estimate in Adam, we  
 64 refer the interested reader to Shen et al. [18] for the details.

## 65 3 Experiments

66 To evaluate the effectiveness of the proposed method, we use IVON to finetune pretrained Llama 2  
 67 [19] model with 7 billion parameters on six datasets with commonsense reasoning multiple-choice  
 68 or true/false questions. These six datasets are WinoGrande-Small (WG-S), WinoGrande-Medium  
 69 (WG-M) [17], ARC-Challenge (ARC-C), ARC-Easy (ARC-E) [4], OpenBookQA (OBQA) [14], and  
 70 BoolQ [3]. We evaluate the performance of the trained LoRA adapters by calculating the accuracy,  
 71 expected calibration error (ECE) and negative log-likelihood (NLL) on the test set. We add NLL  
 72 as another metric for calibration, since ECE is found not always reliable for this purpose [1]. As  
 73 for baseline methods, we compare the performance of IVON-LoRA adapters with LoRA adapters  
 74 trained using AdamW. We also consider other methods for improving generalization and calibration,  
 75 including Monte Carlo Dropout (MC Dropout) [6], Laplace Approximation (LA) [23], Stochastic  
 76 Weight Averaging (SWA) [10, 16], and SWA-Gaussian (SWAG) [12, 16].

77 Compared to standard AdamW finetuning which learns a point estimation of the parameters, IVON  
 78 finetuning learns a *distribution* over the parameters, which allows sampling to obtain an ensemble of  
 79 models during inference. Considering this, we evaluate two variants of inference with IVON-LoRA  
 80 adapters: predicting at the mean of the posterior distribution, and predicting at an ensemble of 10

<sup>1</sup><https://github.com/team-approx-bayes/ivon>

Table 1: Comparison of techniques applied to finetuning/finetuned Llama-2 7B model across commonsense reasoning datasets. Results at the end of training are reported, with subscripts indicating standard error of the mean across 3 runs. We show the relative metric changes achieved by using IVON over AdamW in parentheses, with improvements in blue and degradation in red. The methods marked with \* do not require customized pipeline or additional computation during inference.

Metrics	Methods	WG-S	ARC-C	ARC-E	WG-M	OBQA	BoolQ	Average	
ACC $\uparrow$	AdamW*	66.5 <sub>0.4</sub>	66.7 <sub>0.5</sub>	84.9 <sub>0.2</sub>	73.5 <sub>0.4</sub>	78.9 <sub>0.7</sub>	85.8 <sub>0.1</sub>	76.1	
	+ MC Drop	66.7 <sub>0.4</sub>	67.3 <sub>0.5</sub>	84.8 <sub>0.4</sub>	73.7 <sub>0.2</sub>	79.3 <sub>0.5</sub>	85.9 <sub>0.2</sub>	76.3	
	+ LA (KFAC)	66.6 <sub>0.3</sub>	66.0 <sub>1.4</sub>	84.3 <sub>0.4</sub>	73.2 <sub>0.3</sub>	78.6 <sub>0.9</sub>	85.7 <sub>0.2</sub>	75.7	
	+ LA (diag)	66.2 <sub>0.3</sub>	61.2 <sub>1.9</sub>	81.8 <sub>0.5</sub>	73.3 <sub>0.3</sub>	79.7 <sub>0.8</sub>	85.7 <sub>0.2</sub>	74.7	
	+ SWA*	69.7 <sub>0.6</sub>	67.2 <sub>1.3</sub>	85.2 <sub>0.1</sub>	75.6 <sub>0.2</sub>	79.8 <sub>0.5</sub>	85.5 <sub>0.1</sub>	77.2	
	+ SWAG	69.4 <sub>0.6</sub>	68.4 <sub>1.3</sub>	85.1 <sub>0.2</sub>	75.2 <sub>0.4</sub>	80.1 <sub>0.1</sub>	85.2 <sub>0.2</sub>	77.2	
	IVON@mean*	(+5.6) 72.1 <sub>0.5</sub>	(+3.2) 69.9 <sub>0.7</sub>	(+2.6) 87.5 <sub>0.6</sub>	(+3.1) 76.6 <sub>0.5</sub>	(+2.0) 80.9 <sub>0.6</sub>	(+0.3) 86.1 <sub>0.2</sub>	(+2.8) 78.9	
	IVON	(+5.7) 72.2 <sub>0.5</sub>	(-0.4) 66.3 <sub>0.6</sub>	(+0.8) 85.7 <sub>0.3</sub>	(+2.9) 76.4 <sub>0.6</sub>	(+1.5) 80.4 <sub>0.4</sub>	(-0.1) 85.7 <sub>0.2</sub>	(+1.7) 77.8	
	ECE ( $\times 100$ ) $\downarrow$	AdamW*	32.8 <sub>0.5</sub>	31.4 <sub>0.6</sub>	14.5 <sub>0.3</sub>	25.3 <sub>0.4</sub>	19.1 <sub>0.8</sub>	7.6 <sub>0.2</sub>	21.8
		+ MC Drop	30.7 <sub>0.3</sub>	28.8 <sub>0.7</sub>	13.4 <sub>0.4</sub>	23.6 <sub>0.1</sub>	17.5 <sub>0.6</sub>	7.6 <sub>0.3</sub>	20.2
+ LA (KFAC)		5.2 <sub>2.0</sub>	12.4 <sub>2.5</sub>	5.4 <sub>1.6</sub>	11.1 <sub>1.2</sub>	5.5 <sub>0.2</sub>	3.9 <sub>0.1</sub>	7.3	
+ LA (diag)		12.4 <sub>1.2</sub>	16.3 <sub>1.7</sub>	24.2 <sub>3.6</sub>	5.8 <sub>0.6</sub>	13.1 <sub>1.2</sub>	19.7 <sub>0.2</sub>	15.3	
+ SWA*		19.7 <sub>0.5</sub>	24.6 <sub>0.8</sub>	9.8 <sub>0.2</sub>	12.3 <sub>1.4</sub>	9.7 <sub>0.4</sub>	2.4 <sub>0.2</sub>	13.1	
+ SWAG		12.9 <sub>1.1</sub>	15.6 <sub>1.1</sub>	5.7 <sub>0.3</sub>	8.0 <sub>1.2</sub>	5.1 <sub>0.4</sub>	1.1 <sub>0.4</sub>	8.1	
IVON@mean*		(-5.3) 27.5 <sub>0.4</sub>	(-5.6) 25.8 <sub>0.4</sub>	(-4.4) 10.1 <sub>0.4</sub>	(-2.3) 23.0 <sub>0.5</sub>	(-7.9) 11.2 <sub>0.5</sub>	(-2.0) 5.6 <sub>0.1</sub>	(-4.6) 17.2	
IVON		(-11.0) 21.8 <sub>0.8</sub>	(-20.7) 10.7 <sub>0.4</sub>	(-10.9) 3.6 <sub>0.6</sub>	(-3.9) 21.4 <sub>0.5</sub>	(-15.8) 3.3 <sub>0.9</sub>	(-5.0) 2.6 <sub>0.2</sub>	(-11.2) 10.6	
NLL $\downarrow$		AdamW*	4.19 <sub>0.43</sub>	3.71 <sub>0.49</sub>	1.52 <sub>0.05</sub>	2.03 <sub>0.06</sub>	1.54 <sub>0.05</sub>	0.44 <sub>0.01</sub>	2.24
		+ MC Drop	3.75 <sub>0.33</sub>	3.25 <sub>0.38</sub>	1.36 <sub>0.07</sub>	1.85 <sub>0.05</sub>	1.40 <sub>0.04</sub>	0.43 <sub>0.01</sub>	2.01
	+ LA (KFAC)	0.63 <sub>0.01</sub>	0.96 <sub>0.03</sub>	0.49 <sub>0.02</sub>	0.76 <sub>0.01</sub>	0.68 <sub>0.00</sub>	0.37 <sub>0.00</sub>	0.65	
	+ LA (diag)	0.66 <sub>0.01</sub>	1.05 <sub>0.05</sub>	0.70 <sub>0.05</sub>	0.57 <sub>0.01</sub>	0.65 <sub>0.01</sub>	0.47 <sub>0.00</sub>	0.68	
	+ SWA*	0.87 <sub>0.02</sub>	1.34 <sub>0.06</sub>	0.55 <sub>0.00</sub>	0.63 <sub>0.04</sub>	0.65 <sub>0.02</sub>	0.34 <sub>0.00</sub>	0.73	
	+ SWAG	0.68 <sub>0.02</sub>	1.00 <sub>0.04</sub>	0.46 <sub>0.00</sub>	0.56 <sub>0.02</sub>	0.55 <sub>0.02</sub>	0.34 <sub>0.00</sub>	0.60	
	IVON@mean*	(-0.69) 3.50 <sub>0.07</sub>	(-1.74) 1.97 <sub>0.03</sub>	(-0.83) 0.69 <sub>0.00</sub>	(+0.40) 2.43 <sub>0.04</sub>	(-0.88) 0.66 <sub>0.02</sub>	(-0.08) 0.36 <sub>0.00</sub>	(-0.64) 1.60	
	IVON	(-1.94) 2.25 <sub>0.09</sub>	(-2.71) 1.00 <sub>0.03</sub>	(-1.12) 0.40 <sub>0.00</sub>	(+0.13) 2.16 <sub>0.01</sub>	(-1.00) 0.54 <sub>0.01</sub>	(-0.11) 0.33 <sub>0.00</sub>	(-1.13) 1.11	

81 samples from the posterior distribution (referred to as IVON@mean and IVON, respectively). For a  
82 fair comparison, we use the same number of samples for MC Dropout, SWA and SWAG.

83 We present the results in Table 1. First, we observe that IVON, as an alternative to AdamW,  
84 significantly improves the generalization of LoRA finetuning. When evaluated at the mean, IVON  
85 outperforms standard AdamW finetuning and other Bayesian adaptations of LoRA on all datasets in  
86 terms of accuracy, often by a large margin. We also observe that IVON exhibits improved calibration  
87 compared to AdamW and MC Dropout baselines, as indicated by the lower ECE and NLL values.  
88 We want to highlight that these improvements can be achieved with minimal changes to the codebase,  
89 that is, replacing the AdamW optimizer with IVON, and without any additional postprocessing steps  
90 or noticeable overhead.

91 Next, we observe that ensembling with samples from IVON’s posterior distribution further improves  
92 calibration. When evaluate at an ensemble of 10 samples, IVON outperforms all other methods and is  
93 comparable to the best-performing LA (with a Kronecker-factored Hessian) and SWAG on ECE and  
94 NLL. Notably, IVON achieves this despite using a diagonal Hessian and without an additional pass  
95 through the data for computing Hessians at the converged point as in Laplace methods. With this  
96 improvement in calibration, IVON still maintains comparable or better accuracy over other methods.

97 We also notice that by scaling the learned variance during inference (sampling from  $\mathcal{N}(\mathbf{m}, \text{diag}(\tau\mathbf{v}))$   
98 with  $\tau$  being a scaling factor), the ensemble of IVON samples can have different behaviors in terms of  
99 accuracy and calibration. In Table 2, we present the results of IVON with different choices of  $\tau$  during  
100 inference. When a smaller  $\tau$  is used, the ensemble of IVON samples achieves slight improvements  
101 both in accuracy and calibration. On the other hand, the ensemble is further improved in calibration  
102 at the cost of accuracy when a larger  $\tau$  is used. This can be useful in tweaking the trade-off between  
103 accuracy and calibration to suit the needs of different applications.

104 Finally, we observe that the overhead induced by IVON is negligible. To investigate this, we profile  
105 our training code on an NVIDIA RTX 6000 Ada GPU. In our test run, the forward pass, loss  
106 computation, and backward pass of a training step take in total 316.3ms on average. As for the  
107 overhead of IVON, the sampling procedure and the optimization step of each training step take 1.8ms  
108 and 1.0ms on average, respectively, which is less than 1% of the running time of a training step.

Table 2: Comparison of different variance scaling factor  $\tau$  used during IVON inference. For  $\tau = 0$ , it is equivalent to inference at the posterior mean. For  $\tau = 1$ , it recovers the standard sampling. Results at the end of training are reported. The subscripts indicate standard error of the mean across 3 runs.

Metrics	$\tau$	WG-S	ARC-C	ARC-E	WG-M	OBQA	BoolQ
ACC $\uparrow$	0 (IVON@mean)	72.1 <sub>0.5</sub>	69.9 <sub>0.7</sub>	<b>87.5</b> <sub>0.6</sub>	76.6 <sub>0.5</sub>	<b>80.9</b> <sub>0.6</sub>	86.1 <sub>0.2</sub>
	0.25	72.3 <sub>0.6</sub>	<b>71.1</b> <sub>1.1</sub>	87.4 <sub>0.6</sub>	76.6 <sub>0.5</sub>	<b>80.9</b> <sub>0.7</sub>	<b>86.2</b> <sub>0.2</sub>
	0.5	<b>72.4</b> <sub>0.6</sub>	70.1 <sub>1.1</sub>	87.2 <sub>0.6</sub>	<b>76.7</b> <sub>0.5</sub>	80.8 <sub>0.8</sub>	86.1 <sub>0.2</sub>
	0.75	72.2 <sub>0.7</sub>	68.6 <sub>1.0</sub>	86.7 <sub>0.4</sub>	76.5 <sub>0.4</sub>	80.7 <sub>0.4</sub>	86.1 <sub>0.1</sub>
	1 (IVON)	72.2 <sub>0.5</sub>	66.3 <sub>0.6</sub>	85.7 <sub>0.3</sub>	76.4 <sub>0.6</sub>	80.4 <sub>0.4</sub>	85.7 <sub>0.2</sub>
	1.25	71.7 <sub>0.5</sub>	56.0 <sub>0.6</sub>	81.5 <sub>0.6</sub>	76.3 <sub>0.5</sub>	75.6 <sub>0.4</sub>	84.8 <sub>0.1</sub>
	1.5	70.5 <sub>0.4</sub>	30.2 <sub>0.4</sub>	48.6 <sub>2.8</sub>	76.6 <sub>0.4</sub>	53.5 <sub>1.9</sub>	82.7 <sub>0.4</sub>
ECE ( $\times 100$ ) $\downarrow$	0 (IVON@mean)	27.5 <sub>0.4</sub>	25.8 <sub>0.4</sub>	10.1 <sub>0.4</sub>	23.0 <sub>0.5</sub>	11.2 <sub>0.5</sub>	5.6 <sub>0.1</sub>
	0.25	25.6 <sub>0.5</sub>	21.9 <sub>0.7</sub>	9.2 <sub>0.4</sub>	21.7 <sub>0.5</sub>	10.4 <sub>0.8</sub>	5.1 <sub>0.2</sub>
	0.5	24.3 <sub>0.5</sub>	19.3 <sub>0.9</sub>	8.1 <sub>0.5</sub>	21.6 <sub>0.3</sub>	8.8 <sub>0.6</sub>	4.7 <sub>0.1</sub>
	0.75	23.1 <sub>0.7</sub>	15.7 <sub>1.0</sub>	6.2 <sub>0.3</sub>	21.6 <sub>0.3</sub>	6.3 <sub>0.8</sub>	3.9 <sub>0.3</sub>
	1 (IVON)	21.8 <sub>0.8</sub>	10.7 <sub>0.4</sub>	<b>3.6</b> <sub>0.6</sub>	21.4 <sub>0.5</sub>	<b>3.3</b> <sub>0.9</sub>	2.6 <sub>0.2</sub>
	1.25	20.4 <sub>0.9</sub>	<b>7.5</b> <sub>0.6</sub>	4.9 <sub>0.7</sub>	21.3 <sub>0.4</sub>	9.2 <sub>0.8</sub>	<b>1.1</b> <sub>0.3</sub>
	1.5	<b>19.2</b> <sub>0.6</sub>	14.2 <sub>2.0</sub>	5.6 <sub>1.9</sub>	<b>20.9</b> <sub>0.2</sub>	15.0 <sub>1.6</sub>	2.1 <sub>0.4</sub>
NLL $\downarrow$	0 (IVON@mean)	3.50 <sub>0.07</sub>	1.97 <sub>0.03</sub>	0.69 <sub>0.00</sub>	2.43 <sub>0.04</sub>	0.66 <sub>0.02</sub>	0.36 <sub>0.00</sub>
	0.25	3.11 <sub>0.05</sub>	1.70 <sub>0.03</sub>	0.62 <sub>0.01</sub>	2.22 <sub>0.01</sub>	0.63 <sub>0.02</sub>	0.35 <sub>0.00</sub>
	0.5	2.87 <sub>0.08</sub>	1.41 <sub>0.02</sub>	0.53 <sub>0.01</sub>	2.22 <sub>0.02</sub>	0.59 <sub>0.02</sub>	0.35 <sub>0.00</sub>
	0.75	2.52 <sub>0.07</sub>	1.17 <sub>0.02</sub>	0.46 <sub>0.01</sub>	2.20 <sub>0.02</sub>	0.54 <sub>0.01</sub>	<b>0.34</b> <sub>0.00</sub>
	1 (IVON)	2.20 <sub>0.06</sub>	<b>0.99</b> <sub>0.01</sub>	<b>0.41</b> <sub>0.02</sub>	2.17 <sub>0.02</sub>	<b>0.53</b> <sub>0.01</sub>	<b>0.34</b> <sub>0.00</sub>
	1.25	1.88 <sub>0.08</sub>	1.11 <sub>0.01</sub>	0.51 <sub>0.02</sub>	2.14 <sub>0.03</sub>	0.65 <sub>0.01</sub>	0.35 <sub>0.00</sub>
	1.5	<b>1.60</b> <sub>0.10</sub>	1.47 <sub>0.01</sub>	1.19 <sub>0.05</sub>	<b>2.11</b> <sub>0.03</sub>	1.17 <sub>0.01</sub>	0.39 <sub>0.01</sub>

## 109 4 Discussion

110 Our direct variational learning approach using IVON is surprisingly effective for improving calibration  
 111 and accuracy in LoRA finetuning. Given the strong results, we hope that this work invigorates research  
 112 in variational methods for LLMs. Reasons for IVON’s success are not fully understood, but one  
 113 hypothesis is the prevention of overfitting as the finetuning datasets are often comparably small.  
 114 This may be attributed to the preference for simpler solutions (flatter minima) which is inherent in  
 115 variational learning [8, 7].

116 On most of the datasets, ensemble of IVON samples outperforms IVON evaluated at the posterior  
 117 mean on ECE and NLL, but at the cost of a slight decrease in accuracy. We want to point out this  
 118 is possibly due to the limited number of samples used in the ensemble. We draw 10 samples for  
 119 all the ensemble-based methods in our experiments, both to follow the setting in Yang et al. [23]  
 120 and to keep the computational cost manageable. It is possible that using more IVON samples could  
 121 further improve the performance of the ensemble, which is reported in Shen et al. [18] on image  
 122 classification tasks. Nevertheless, the parameter uncertainty obtained by IVON is expected to be  
 123 useful for several downstream tasks such as sensitivity analysis [15] and model merging [5], which  
 124 will be explored in future work.

125 A limitation shared with other Bayesian LoRA methods [23, 16] is that the learned posterior over the  
 126 increment  $\Delta\mathbf{W}$  is non-Gaussian, as it is the product of two Gaussian random variables. Therefore,  
 127 it cannot be easily merged down into or easily combined with Gaussian uncertainty on the full  
 128 weights  $\mathbf{W}$ . A different approach would be to use a variational low-rank correction to correct the  
 129 mean and variance of a Laplace approximation of the original model. van Niekerk and Rue [20]  
 130 propose such a low-rank approach in the context of latent Gaussian models, and adapting these ideas  
 131 to large language models may be an interesting direction for future work.

132 IVON also has some practical limitations. The method introduces two new hyperparameters over  
 133 AdamW, which are  $\lambda$  in (2) and the initialization of the posterior variance. This makes tuning IVON  
 134 a bit more involved than tuning AdamW and the results depend on setting these parameters well.  
 135 While a good heuristic is to set  $\lambda$  as small as possible while still retaining stable training and setting  
 136 the posterior initialization in the order of magnitude of the final posterior variance, more principled  
 137 or automatic ways to set them reliably would be desirable.

## References

- 138
- 139 [1] Joris Baan, Wilker Aziz, Barbara Plank, and Raquel Fernández. Stop measuring calibration when humans  
140 disagree. In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*,  
141 pages 1892–1915, 2022.
- 142 [2] Haolin Chen and Philip N Garner. A Bayesian interpretation of adaptive low-rank adaptation.  
143 *arXiv:2409.10673*, 2024.
- 144 [3] Christopher Clark, Kenton Lee, Ming-Wei Chang, Tom Kwiatkowski, Michael Collins, and Kristina  
145 Toutanova. Boolq: Exploring the surprising difficulty of natural yes/no questions. In *Proceedings of the*  
146 *2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human*  
147 *Language Technologies, Volume 1 (Long and Short Papers)*, pages 2924–2936, 2019.
- 148 [4] Peter Clark, Isaac Cowhey, Oren Etzioni, Tushar Khot, Ashish Sabharwal, Carissa Schoenick, and Oyvind  
149 Tafjord. Think you have solved question answering? try arc, the ai2 reasoning challenge. *arXiv preprint*  
150 *arXiv:1803.05457*, 2018.
- 151 [5] Nico Daheim, Thomas Möllenhoff, Edoardo Maria Ponti, Iryna Gurevych, and Mohammad Emtiyaz  
152 Khan. Model merging by uncertainty-based gradient matching. In *International Conference on Learning*  
153 *Representations (ICLR)*, 2024.
- 154 [6] Yarín Gal and Zoubin Ghahramani. Dropout as a Bayesian approximation: Representing model uncertainty  
155 in deep learning. In *International Conference on Machine Learning (ICML)*, 2016.
- 156 [7] Alex Graves. Practical variational inference for neural networks. In *Advances in Neural Information*  
157 *Processing Systems (NeurIPS)*, 2011.
- 158 [8] Geoffrey E Hinton and Drew Van Camp. Keeping the neural networks simple by minimizing the description  
159 length of the weights. In *Conference on Learning Theory (COLT)*, 1993.
- 160 [9] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and  
161 Weizhu Chen. LoRA: Low-rank adaptation of large language models. *arXiv preprint arXiv:2106.09685*,  
162 2021.
- 163 [10] Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averag-  
164 ing weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018.
- 165 [11] Guangliang Liu, Zhiyu Xue, Xitong Zhang, Kristen Marie Johnson, and Rongrong Wang. PAC-tuning:  
166 Fine-tuning pretrained language models with PAC-driven perturbed gradient descent. *arXiv:2310.17588*,  
167 2023.
- 168 [12] Wesley J Maddox, Pavel Izmailov, Timur Garipov, Dmitry P Vetrov, and Andrew Gordon Wilson. A simple  
169 baseline for bayesian uncertainty in deep learning. *Advances in neural information processing systems*, 32,  
170 2019.
- 171 [13] Sourab Mangrulkar, Sylvain Gugger, Lysandre Debut, Younes Belkada, Sayak Paul, and Benjamin Bossan.  
172 Peft: State-of-the-art parameter-efficient fine-tuning methods. <https://github.com/huggingface/peft>, 2022.  
173
- 174 [14] Todor Mihaylov, Peter Clark, Tushar Khot, and Ashish Sabharwal. Can a suit of armor conduct electricity?  
175 a new dataset for open book question answering. In *Proceedings of the 2018 Conference on Empirical*  
176 *Methods in Natural Language Processing*, pages 2381–2391, 2018.
- 177 [15] Peter Nickl, Lu Xu, Dharmesh Tailor, Thomas Möllenhoff, and Mohammad Emtiyaz Khan. The memory  
178 perturbation equation: Understanding model’s sensitivity to data. In *Advances in Neural Information*  
179 *Processing Systems (NeurIPS)*, 2023.
- 180 [16] Emre Onal, Klemens Flöge, Emma Caldwell, Arsen Sheverdin, and Vincent Fortuin. Gaussian stochastic  
181 weight averaging for bayesian low-rank adaptation of large language models. In *Sixth Symposium on*  
182 *Advances in Approximate Bayesian Inference-Non Archival Track*, 2024.
- 183 [17] Keisuke Sakaguchi, Ronan Le Bras, Chandra Bhagavatula, and Yejin Choi. Winogrande: An adversarial  
184 winograd schema challenge at scale. *Communications of the ACM*, 64(9):99–106, 2021.
- 185 [18] Yuesong Shen, Nico Daheim, Bai Cong, Peter Nickl, Gian Maria Marconi, Bazan Clement Emile Marcel  
186 Raoul, Rio Yokota, Iryna Gurevych, Daniel Cremers, Mohammad Emtiyaz Khan, and Thomas Möllenhoff.  
187 Variational learning is effective for large deep networks. In *International Conference on Machine Learning*,  
188 2024.

- 189 [19] Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay  
190 Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, et al. Llama 2: Open foundation and  
191 fine-tuned chat models. *arXiv preprint arXiv:2307.09288*, 2023.
- 192 [20] Janet van Niekerk and Haavard Rue. Low-rank variational bayes correction to the laplace method. *J. Mach.*  
193 *Learn. Res. (JMLR)*, 25(62):1–25, 2024.
- 194 [21] Xi Wang, Laurence Aitchison, and Maja Rudolph. LoRA ensembles for large language model fine-tuning.  
195 *arXiv preprint arXiv:2310.00035*, 2023.
- 196 [22] Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric  
197 Cistac, Tim Rault, Rémi Louf, Morgan Funtowicz, et al. Transformers: State-of-the-art natural language  
198 processing. In *Proceedings of the 2020 conference on empirical methods in natural language processing:*  
199 *system demonstrations*, pages 38–45, 2020.
- 200 [23] Adam X Yang, Maxime Robeyns, Xi Wang, and Laurence Aitchison. Bayesian low-rank adaptation for  
201 large language models. In *The Twelfth International Conference on Learning Representations*, 2024.
- 202 [24] Q Zhang, M Chen, A Bukharin, P He, Y Cheng, W Chen, and T Zhao. Adaptive budget allocation for  
203 parameter-efficient fine-tuning. In *International Conference on Learning Representations (ICLR)*, 2023.
- 204 [25] Maxim Zhelmin, Viktor Moskvoretskii, Egor Shvetsov, Egor Venediktov, Mariya Krylova, Aleksandr  
205 Zuev, and Evgeny Burnaev. GIFT-SW: Gaussian noise injected fine-tuning of salient weights for LLMs.  
206 *arXiv:2408.15300*, 2024.

Table 3: IVON hyperparameters used in experiments.

Hyperparameter	WG-S	ARC-C	ARC-E	WG-M	OBQA	BoolQ
Effective sample size	$1 \times 10^7$	$1 \times 10^6$	$1 \times 10^6$	$1 \times 10^8$	$1 \times 10^6$	$1 \times 10^7$
Hessian initialization	$3 \times 10^{-4}$	$1 \times 10^{-3}$	$1 \times 10^{-3}$	$3 \times 10^{-4}$	$1 \times 10^{-3}$	$3 \times 10^{-4}$
Learning rate				0.03		
Gradient momentum				0.9		
Hessian momentum				$1 - 10^{-5}$		
Clip radius				$10^{-3}$		

207 **Supplementary Material**

208 **Details on experimental setup**

209 Our experimental design is based on Yang et al. [23]. We utilize the PEFT [13] library for LoRA  
 210 adaptation, and apply LoRA to the query and value weights of the attention layers. Unlike in Yang  
 211 et al. [23], we do not apply LoRA to the output layer due to numerical instability encountered in  
 212 some preliminary experiments. The base model is quantized to 8-bit precision, with LoRA weights  
 213 maintained in 16-bit precision. Finetuning is performed on a single NVIDIA RTX 6000 Ada GPU  
 214 with a batch size of 4 for 10,000 steps, without gradient accumulation.

215 To finetune a pretrained language model which predicts the next token in a sequence for solving  
 216 multiple-choice or true/false questions, we need to wrap the text and the choice of each question  
 217 with predefined prompt templates to an instruction. We then use the pretrained model to predict  
 218 the next token of the wrapped instruction, and extract the output logits for the tokens standing for  
 219 "True"/"False" or "A"/"B"/"C"/"D" choices. For the prompt templates, we use the same ones as in  
 220 Yang et al. [23]. An example of such a prompt (used for WG-S and WG-M datasets) is as follows:

221           Select one of the choices that answers the following question: {question}  
 222           Choices: A. {option1}. B {option2}. Answer:

223 **Hyperparameters**

224 As for the hyperparameters of LoRA and AdamW finetuning, we use the same settings as in Yang  
 225 et al. [23], which are also the default settings in Huggingface’s Transformers [22] and PEFT [13]  
 226 library. For LoRA, we set the rank  $r$  to 8,  $\alpha$  to 16, and the dropout rate to 0.1. For AdamW optimizer,  
 227 we set the initial learning rate to  $5 \times 10^{-5}$ , weight decay to 0, and use a linear learning rate scheduler  
 228 which decays the learning rate to 0 at the end of the training.

229 Working IVON hyperparameters and guidelines for choosing them are discussed in Shen et al. [18].  
 230 Still, it is not well understood how to choose them in the context of LoRA finetuning. We empirically  
 231 find that setting  $\lambda$  as small as possible while still retaining stable training is a good heuristic. To  
 232 choose the initialization value  $v_0$  of the posterior variance, we track the mean value of the running  
 233 average of the posterior variance for the first few training steps. We notice that if the mean value  
 234 changes significantly during the first few steps, then the initialization value is likely too far from a  
 235 reasonable one. We follow the guideline in Shen et al. [18] and set the learning rate of IVON to 0.03,  
 236 Hessian momentum to  $1 - 10^{-5}$ , and clip radius to  $10^{-3}$ . Finally, We summarize the hyperparameters  
 237 of IVON used in our experiments in Table 3.