

# CALIBRATING THE RIGGED LOTTERY: MAKING ALL TICKETS RELIABLE

**Bowen Lei**

Texas A&M University  
bowenlei@stat.tamu.edu

**Ruqi Zhang**

Purdue University  
ruqiz@purdue.edu

**Dongkuan Xu**

North Carolina State University  
dxu27@ncsu.edu

**Bani Mallick**

Texas A&M University  
bmallick@stat.tamu.edu

## ABSTRACT

Although sparse training has been successfully used in various resource-limited deep learning tasks to save memory, accelerate training, and reduce inference time, the reliability of the produced sparse models remains unexplored. Previous research has shown that deep neural networks tend to be over-confident, and we find that sparse training exacerbates this problem. Therefore, calibrating the sparse models is crucial for reliable prediction and decision-making. In this paper, we propose a new sparse training method to produce sparse models with improved confidence calibration. In contrast to previous research that uses only one mask to control the sparse topology, our method utilizes two masks, including a deterministic mask and a random mask. The former efficiently searches and activates important weights by exploiting the magnitude of weights and gradients. While the latter brings better exploration and finds more appropriate weight values by random updates. Theoretically, we prove our method can be viewed as a hierarchical variational approximation of a probabilistic deep Gaussian process. Extensive experiments on multiple datasets, model architectures, and sparsities show that our method reduces ECE values by up to 47.8% and simultaneously maintains or even improves accuracy with only a slight increase in computation and storage burden.

## 1 INTRODUCTION

Sparse training is gaining increasing attention and has been used in various deep neural network (DNN) learning tasks (Evci et al., 2020; Dietrich et al., 2021; Bibikar et al., 2022). In sparse training, a certain percentage of connections are maintained being removed to save memory, accelerate training, and reduce inference time, enabling DNNs for resource-constrained situations. The sparse topology is usually controlled by a mask, and various sparse training methods have been proposed to find a suitable mask to achieve comparable or even higher accuracy compared to dense training (Evci et al., 2020; Liu et al., 2021; Schwarz et al., 2021). However, in order to deploy the sparse models in real-world applications, a key question remains to be answered: how reliable are these models?

There has been a line of work on studying the reliability of dense DNNs, which means that DNNs should know what it does not know (Guo et al., 2017; Nixon et al., 2019; Wang et al., 2021). In other words, a model’s confidence (the probability associated with the predicted class label) should reflect its ground truth correctness likelihood. A widely used reliability metric is Expected Calibration Error (ECE) (Guo et al., 2017), which measures the difference between confidence and accuracy, with a lower ECE indicating higher reliability. However, prior research has shown that DNNs tend to be over-confident (Guo et al., 2017; Rahaman et al., 2021; Patel et al., 2022), suggesting DNNs may be too confident to notice incorrect decisions, leading to safety issues in real-world applications, e.g., automated healthcare and self-driving cars (Jiang et al., 2012; Bojarski et al., 2016).

In this work, we for the *first* time identify and study the reliability problem of sparse training. We start with the question of how reliable the current sparse training is. We find that the over-confidence

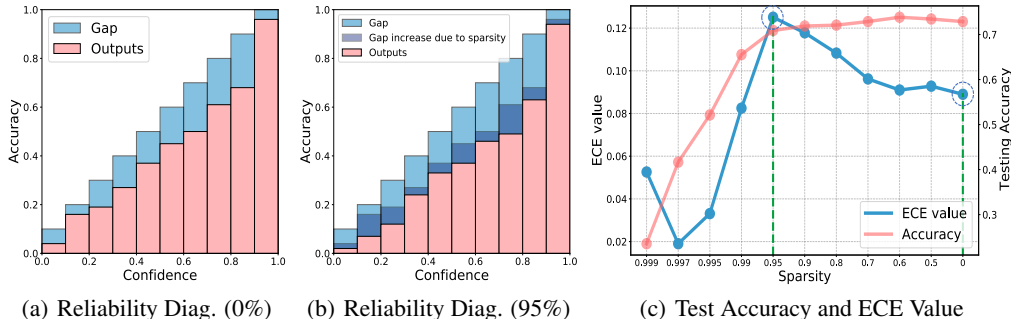


Figure 1: Reliability diagrams for (a) the dense model and (b) the sparse model. The sparse model is more over-confident than the dense model. (c) the scatter plot of test accuracy (%) and ECE value at different sparsities. From the high sparse model to the dense model, the ECE value first decreases, then increases, and then decreases again, showing a double descent pattern.

problem becomes even more pronounced when sparse training is applied to ResNet-50 on CIFAR-100. Figures 1 (a)-(b) show that the gap (blue area) between confidence and accuracy of the sparse model (95% sparsity) is larger than that of the dense model (0% sparsity), implying the sparse model is more over-confident than the dense model. Figure 1 (c) shows the test accuracy (pink curve) and ECE value (blue curve, a measure of reliability) (Guo et al., 2017) at different sparsities. When the accuracy is comparable to dense training (0%-95%), we observe that the ECE values increase with sparsity, implying that the problem of over-confidence becomes more severe at higher sparsity. And when the accuracy decreases sharply (>95%), the ECE value first decreases and then increases again. This leads to a *double descent* phenomenon (Nakkiran et al., 2021) when we view the ECE value curve from left to right (99.9%-0%) (see Section 6 for more discussion).

To improve the reliability, we propose a new sparse training method to produce well-calibrated predictions while maintaining a high accuracy performance. We call our method “The Calibrated Rigged Lottery” or CigL. Unlike previous sparse training methods with only one mask, our method employs two masks, including a deterministic mask and a random mask, to better explore the sparse topology and weight space. The deterministic one efficiently searches and activates important weights by exploiting the magnitude of weights/gradients. And the random one, inspired by dropout, adds more exploration and leads to better convergence. When near the end of training, we collect weights & masks at each epoch, and use the designed weight & mask averaging procedure to obtain one sparse model. From theoretical analysis, we show our method can be viewed as a hierarchical variational approximation (Ranganath et al., 2016) to a probabilistic deep Gaussian process (Gal & Ghahramani, 2016), which leads to a large family of variational distributions and better Bayesian posterior approximations. Our contributions are summarized as follows:

- We for the first time identify and study the reliability problem of sparse training and find that sparse training exacerbates the over-confidence problem of DNNs.
- We then propose CigL, a new sparse training method that improves confidence calibration with comparable and even higher accuracy.
- We prove that CigL can be viewed as a hierarchical variational approximation to a probabilistic deep Gaussian process which improves the calibration by better characterizing the posterior.
- We perform extensive experiments on multiple benchmark datasets, model architectures, and sparsities. CigL reduces ECE values by up to **47.8%** and simultaneously maintain or even improve accuracy with only a slight increase in computational and storage burden.

## 2 RELATED WORK

### 2.1 SPARSE TRAINING

As the scale of models continues to grow, there is an increasing attention to the sparse training which maintains sparse weights throughout the training process. Different sparse training methods have been investigated, and various pruning and growth criteria, such as weight/gradient magnitude,

are designed (Mocanu et al., 2018; Bellec et al., 2018; Frankle & Carbin, 2019; Mostafa & Wang, 2019; Dettmers & Zettlemoyer, 2019; Evci et al., 2020; Jayakumar et al., 2020; Liu et al., 2021; Özdenizci & Legenstein, 2021; Zhou et al., 2021; Schwarz et al., 2021; Yin et al., 2022). However, sparse training is more challenging in the weight space exploration because sparse constraints cut off update routes and produce spurious local minima (Evci et al., 2019; Sun & Li, 2021; He et al., 2022). There are some studies that have started to promote exploration, while they primarily pursue high accuracy and might add additional costs (Liu et al., 2021; Huang et al., 2022). Most sparse training methods use only one mask to determine the sparse topology, which is insufficient for achieving adequate exploration, and existing multi-mask methods are not designed for improved exploration (Xia et al., 2022; Bibikar et al., 2022) (more details in Section D.4).

## 2.2 CONFIDENCE CALIBRATION IN DNNs

Many studies have investigated whether the confidences of DNNs are well-calibrated (Guo et al., 2017; Nixon et al., 2019; Zhang et al., 2020), and existing research has found DNNs tend to be over-confident (Guo et al., 2017; Rahaman et al., 2021; Patel et al., 2022), which may mislead our choices and cause unreliable decisions in real-world applications. To improve confidence calibration, a widely-used method is temperature scaling (Guo et al., 2017), which adds a scaling parameter to the softmax formulation and adjusts it on a validation set. Some other works incorporate regularization in the training, such as Mixup (Zhang et al., 2017) and label smoothing (Szegedy et al., 2016). In addition, Bayesian methods also show the ability to improve calibration, such as Monte Carlo Dropout (Gal & Ghahramani, 2016) and Bayesian deep ensembles (Ashukha et al., 2020). However, they mainly focus on dense training. Studies have been conducted on reliability of sparse DNNs (more details in Section D.5), but they target on pruning, starting with a dense model to gradually increase sparsity, which reduces the exploration challenge (Venkatesh et al., 2020; Chen et al., 2022). They still find that uncertainty measures are more sensitive to pruning than generalization metrics, indicating the sensitivity of reliability to sparsity. Yin et al. (2022) studies sparse training, but it aims to boost the performance and brings limited improvement in reliability. Therefore, how to obtain a well-calibrated DNN in sparse training is more challenging and remains unknown.

## 3 METHOD

We propose a new sparse training method, CigL, to improve the confidence calibration of the produced sparse models, which simultaneously maintains comparable or even higher accuracy. Specifically, CigL starts with a random sparse network and uses two masks to control the sparse topology and explore the weight space, including a deterministic mask and a random mask. The former is updated periodically to determine the non-zero weights, while the latter is sampled randomly in each iteration to bring better exploration in the model update. Then, with the designed weight & mask averaging, we combine information about different aspects of the weight space to obtain a single output sparse model. Our CigL method is outlined in Algorithm 1.

### 3.1 DETERMINISTIC MASK & RANDOM MASK

In our CigL, we propose to utilize two masks, a deterministic mask  $M$  and a random mask  $Z$ , to search for a sparse model with improved confidence calibration and SOTA accuracy. We will first describe the two masks in detail and discuss how to set their sparsity.

**The deterministic mask** controls the entire sparse topology with the aim of finding a well-performing sparse model. That is, the mask determines which weights should be activated and which should not. Inspired by the widely-used sparse training method RigL (Evci et al., 2020), we believe a larger weight/gradient magnitude implies that the weight is more helpful for loss reduction and needs to be activated. Thus, CigL removes a portion of the weights with small magnitudes, and activates new weights with large gradient magnitudes at fixed time intervals  $\Delta T$ .

**The random mask** allows the model to better explore the weight space under sparsity constraints. In each iteration prior to backpropagation, the mask is randomly drawn from Bernoulli distribution. In this way, the mask randomly selects a portion of the non-zero weights to be temporarily deactivated and forces the model to explore more in other directions of the weight space, which adds more randomness in the weight update step and leads to a better exploration of the weight space compared

to one mask strategy. As a result, the model is more likely to jump out of spurious local minima while avoiding deviations from the sparse topology found by the deterministic mask.

**The sparsity setting** of the two masks is illustrated as below. On the one hand, the deterministic mask is responsible for the overall sparsity of the output sparse model. Suppose we want to train a network with 95% sparsity, the deterministic mask will also have the same sparsity, with 5% of the elements being 1. On the other hand, the random mask deactivates some non-zero weights during the training process, producing temporary models with increasing sparsity. Since highly sparse models (like 95% sparsity) are sensitive to further increases in sparsity, we set a low sparsity, such as 10%, for the random mask so that no significant increases in sparsity and no dramatic degradation in performance occurs in these temporary models.

### 3.2 WEIGHT & MASK AVERAGING

With the two masks designed above, we propose a weight & mask averaging procedure to obtain one single sparse model with improved confidence calibration and comparable or even higher accuracy. We formalize this procedure as follows. We first iteratively update the two masks and model weights. Consistent with widely used sparse training methods (Evcı et al., 2020; Liu et al., 2021), the deterministic mask stops updating near the end of the training process. While we still continuously draw different random masks from the Bernoulli distribution and collect a pair of sparse weights and random masks  $\{\mathbf{Z}^{(t)}, \mathbf{W}^{(t)}\}$  at each epoch after  $T^*$ -th epoch. Then, we can produce multiple temporary sparse models  $\mathbf{Z}^{(t)} \odot \mathbf{W}^{(t)}$  with different weight values and different sparse topologies, which contain more knowledge about the weight space than the single-mask training methods. Finally, inspired by a popular way of combining models (Izmailov et al., 2018; Wortsman et al., 2022), we obtain the single output sparse model by averaging the weights of these temporary sparse models, which can be viewed as a mask-based weighted averaging.

## 4 MAKING ALL TICKETS RELIABLE

### 4.1 CIGL WITH BETTER CONFIDENCE CALIBRATION

Obtaining reliable DNNs is more challenging in sparse training, and we will show why our CigL provides a solution to this problem. Bayesian methods have shown the ability to improve confidence calibration (Gal & Ghahramani, 2016; Ashukha et al., 2020), but they become more difficult to fit the posterior well under sparse constraints, limiting their ability to solve unreliable problems. We find that CigL can be viewed as a hierarchical Bayesian method (shown in Section 4.2), which improves confidence calibration by performing better posterior approximations in two ways discussed below.

On the one hand, the model is more challenging to fully explore the weight space due to sparsity constraint, and inappropriate weight values can also negatively affect the mask search. During sparse training, a large percentage of connections are removed, cutting off the update routes and thus narrowing the family of Bayesian proposal distributions. This leads to more difficult optimization and sampling. To overcome this issue, our CigL adds a hierarchical structure to the variational distributions so that we have a larger family of distributions, allowing it to capture more complex marginal distributions and reducing the difficulty of fitting the posterior.

On the other hand, when sparsity constraints are added, the posterior landscape changes, leading to more complex posterior distributions. One example is the stronger correlation between hidden

---

### Algorithm 1 CigL

---

**Input:** initialize  $\mathbf{W}^{(0)}$ ,  $\mathbf{M}$ , and  $\mathbf{W}_{\text{CigL}} = \text{None}$ , set epoch length  $m$ , update interval  $\Delta T$ , number of iterations  $T$ , start iteration for weight & mask averaging  $T^*$ , random mask rate  $p$ , and learning rate  $\alpha_t$

**for**  $t = 1$  **to**  $T$  **do**

Sample a mini-batch data  $\mathbf{B}_t$  with size  $n$

**if**  $t \bmod \Delta T = 0$  **then**

Update mask  $\mathbf{M}$  using weights and gradients

Prune and regrow weights  $\mathbf{W}^{(t)}$  based on  $\mathbf{M}$

**end if**

Sample mask  $\mathbf{Z}^{(t)}$  and  $\mathbf{Z}_{ij}^{(t)} = \text{Bernoulli}(p)$

Update sparse weights:  $\mathbf{W}^{(t)} = \mathbf{W}^{(t-1)} - \alpha_t \mathbf{M} \odot \mathbf{Z}^{(t)} \odot \nabla L(\mathbf{M} \odot \mathbf{Z}^{(t)} \odot \mathbf{W}^{(t-1)}; \mathbf{B}_t)$

**if**  $t \bmod m = 0$  and  $t > T^*$  **then**

**if**  $\mathbf{W}_{\text{CigL}} = \text{None}$  **then**

$\mathbf{W}_{\text{CigL}} = \mathbf{M} \odot \mathbf{Z}^{(t)} \odot \mathbf{W}^{(t)}$

$n_{\text{models}} = 1$

**else**

$\mathbf{W}_{\text{CigL}} = \frac{\mathbf{W}_{\text{CigL}} \cdot n_{\text{models}} + \mathbf{M} \odot \mathbf{Z}^{(t)} \odot \mathbf{W}^{(t)}}{n_{\text{models}} + 1}$

$n_{\text{models}} = n_{\text{models}} + 1$

**end if**

**end if**

**end for**

**Output:** Sparse Model Weights  $\mathbf{W}_{\text{CigL}}$

---

variables, such as the random mask  $\mathbf{Z}$  and weight  $\mathbf{W}$  (shown in the Appendix C.1). In a dense model, the accuracy does not change much if we randomly draw  $\mathbf{Z}$  and use  $\mathbf{Z} \odot \mathbf{W}$  compared to using  $\mathbf{W}$ . However, in high sparsity like 95%, we see a significant accuracy drop when  $\mathbf{Z} \odot \mathbf{W}$  is used compared to using  $\mathbf{W}$ . Thus, in CigL, the pairings of  $\mathbf{Z}$  and  $\mathbf{W}$  are collected to capture the correlation, leading to a better posterior approximation.

#### 4.2 CIGL AS A HIERARCHICAL BAYESIAN APPROXIMATION

We prove that training sparse neural networks with our CigL are mathematically equivalent to approximating the probabilistic deep GP (Damianou & Lawrence, 2013; Gal & Ghahramani, 2016) with hierarchical variational inference. We show that the objective of CigL is actually to minimize the Kullback–Leibler (KL) divergence between a hierarchical variational distribution and the posterior of a deep GP. During our study, we do not restrict the architecture, allowing the results applicable to a wide range of applications. Detailed derivation is shown in Appendix B.

We first present the minimisation objective function of CigL for a sparse neural network (NN) model with  $L$  layers and loss function  $E$ . The sparse weights and bias of the  $l$ -th layer are denoted by  $\mathbf{W}_l \in \mathbb{R}^{K_l \times K_{l-1}}$  and  $\mathbf{b}_l \in \mathbb{R}^{K_l}$  ( $l = 1, \dots, L$ ), and the output prediction is denoted by  $\hat{\mathbf{y}}_i$ . Given data  $\{\mathbf{x}_i, y_i\}$ , we train the NN model by iteratively update the deterministic mask and the sparse weights. Since the random mask is drawn from a Bernoulli distribution, it has no parameters that need to be updated. For deterministic mask updates, we prune weights with the smaller weight magnitude and regrow weights with larger gradient magnitude. For the weight update, we minimise Eq. (1) which is composed of the difference between  $\hat{\mathbf{y}}_i$  and the true label  $\mathbf{y}_i$  and a  $L_2$  regularisation.

$$\mathcal{L}_{\text{CigL}} := \frac{1}{N} \sum_{i=1}^N E(y_i, \hat{\mathbf{y}}_i) + \lambda \sum_{l=1}^L (\|\mathbf{W}_l\|_2^2 + \|\mathbf{b}_l\|_2^2). \quad (1)$$

Then, we derive the minimization objective function of Deep GP which is a flexible probabilistic NN model that can model the distribution of functions (Gal & Ghahramani, 2016). Taking regression as an example, we assume that  $\mathbf{W}_l$  is a random matrix and  $\mathbf{w} = \{\mathbf{W}_l\}_{l=1}^L$ , and denote the prior by  $p(\mathbf{w})$ . Then, the predictive distribution of the deep GP can be expressed as Eq. (2) where  $\tau > 0$

$$p(\mathbf{y}|\mathbf{x}, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}|\mathbf{x}, \mathbf{w})p(\mathbf{w}|\mathbf{X}, \mathbf{Y})d\mathbf{w}, \quad (2)$$

$$p(\mathbf{y}|\mathbf{x}, \mathbf{w}) = \mathcal{N}(\mathbf{y}; \hat{\mathbf{y}}, \tau^{-1}\mathbf{I}), \quad \hat{\mathbf{y}} = \sqrt{\frac{1}{K_L}}\mathbf{W}_L\sigma\left(\dots\sqrt{\frac{1}{K_1}}\mathbf{W}_2\sigma(\mathbf{W}_1\mathbf{x} + \mathbf{u}_1)\right).$$

The posterior distribution  $p(\mathbf{w}|\mathbf{X}, \mathbf{Y})$  is intractable, and one way of training the deep GP is variational inference where a family of tractable distributions  $q(\mathbf{w})$  is chosen to approximate the posterior. Specifically, we define the hierarchy of  $q(\mathbf{w})$  as Eq. (3):

$$q(\mathbf{W}_{lij}|\mathbf{Z}_{lij}, \mathbf{U}_{lij}, \mathbf{M}_l) \sim \mathbf{Z}_{lij} \cdot \mathcal{N}(\mathbf{M}_{lij}\mathbf{U}_{lij}, \sigma^2) + (1 - \mathbf{Z}_{lij}) \cdot \mathcal{N}(0, \sigma^2),$$

$$q(\mathbf{M}_l|\mathbf{U}_l) \propto \exp(\mathbf{M}_l \odot (|\mathbf{U}_l| + |\nabla\mathbf{U}_l|)), \quad \mathbf{U}_{lij} \sim \mathcal{N}(\mathbf{V}_{lij}, \sigma^2), \quad \mathbf{Z}_{lij} \sim \text{Bernoulli}(p_l), \quad (3)$$

where  $l$ ,  $i$  and  $j$  denote the layer, row, and column index,  $\mathbf{M}_l$  is a matrix with 0's and constrained 1's,  $\mathbf{W}_l$  is the sparse weights,  $\mathbf{U}_l$  is the variational parameters, and  $\mathbf{V}_l$  is the variational hyper parameters.

Then, we iteratively update  $\mathbf{M}_l$  and  $\mathbf{W}_l$  to approximate the posterior. For the update of  $\mathbf{M}_l$ , we obtain a point estimate by maximising  $q(\mathbf{M}_l|\mathbf{U}_l)$  under the sparsity constraint. In pruning step, since the gradient magnitudes  $|\nabla\mathbf{U}_l|$  can be relatively small compared to the weight magnitudes  $|\mathbf{U}_l|$  after training, we can use  $\exp(\mathbf{M}_l \odot |\mathbf{U}_l|)$  to approximate the distribution. In regrowth step, since the inactive weights are zero, we directly compare the gradient magnitudes  $\exp(\mathbf{M}_l \odot |\nabla\mathbf{U}_l|)$ . Thus, the update of  $\mathbf{M}_l$  is aligned with the update in CigL.

For  $\mathbf{W}_l$ , we minimise the KL divergence between  $q(\mathbf{w})$  and the posterior of deep GP as Eq. (4)

$$- \int q(\mathbf{w}) \log p(\mathbf{Y}|\mathbf{X}, \mathbf{w})d\mathbf{w} + D_{\text{KL}}(q(\mathbf{w})\|p(\mathbf{w})). \quad (4)$$

For the first term in Eq. (4), we can first rewrite it as  $-\sum_{n=1}^N \int q(\mathbf{w}) \log p(y_n|\mathbf{x}_n, \mathbf{w})$ . Then, we can approximate each integration in the sum with a single estimate  $\hat{\mathbf{w}}$ . For the second term in Eq. (4),

we can approximate it as  $\sum_{l=1}^L (\frac{p_l}{2} \|U_l\|_2^2 + \frac{1}{2} \|\mathbf{u}_l\|_2^2)$ . As a result, we can derive the objective as

$$\mathcal{L}_{\text{GP}} := \frac{1}{N} \sum_{i=1}^N \frac{-\log p(\mathbf{y}_n | \mathbf{x}_n, \hat{\mathbf{w}})}{\tau} + \sum_{l=1}^L (\frac{p_l}{2} \|U_l\|_2^2 + \frac{1}{2} \|\mathbf{u}_l\|_2^2), \quad (5)$$

which is shown to have the same form as the objective in Eq. (1) with appropriate hyperparameters for the deep GP. Thus, the update of  $\mathbf{W}_l$  is also consistent with the update in CigL. This suggests that our CigL can be viewed as an approximation to the deep GP using hierarchical variational inference. The final weight & mask averaging procedure can be incorporated into the Bayesian paradigm as an approximation to the posterior distribution (Srivastava et al., 2014; Maddox et al., 2019).

### 4.3 CONNECTION TO DROPOUT

Our CigL can be seen as a new version of Dropout, and our random mask  $\mathbf{Z}$  is related to the Dropout mask. Dropout is a widely used method to overcome the overfitting problem (Hinton et al., 2012; Wan et al., 2013; Srivastava et al., 2014). Two widely used types are unit dropout and weight dropout, which randomly discard units (neurons) and individual weights at each training step, respectively. Both methods use dropouts only in the training phase and remove them in the testing phase, which is equivalent to discarding  $\mathbf{Z}$  and only using  $\mathbf{W}$  for prediction. However, simply dropping  $\mathbf{Z}$  can be detrimental to the fit of the posterior. Thus, MC dropout collects multiple models by randomly selecting multiple dropout masks, which is equivalent to extracting multiple  $\mathbf{Z}$  and using one  $\mathbf{W}$  for prediction. However, only using one  $\mathbf{W}$  neither fully expresses the posterior landscape nor captures the correlation between  $\mathbf{Z}$  and  $\mathbf{W}$ . In contrast, our CigL uses multiple pairings of  $\mathbf{Z}$  and  $\mathbf{W}$ , which can better approximate the posterior under sparsity constraints.

### 4.4 CONNECTION TO WEIGHT AVERAGING

Our weight & mask averaging can be seen as an extension of weight averaging (WA), which averages the weights of multiple model samples to produce a single output model (Izmailov et al., 2018; Wortsman et al., 2022). Compared to deep ensembles (Ashukha et al., 2020), WA outputs only one model, which reduces the forward FLOPs and speeds up prediction. When these model samples are located in one low error basin, it usually leads to wider optima and better generalization. However, although WA can produce better generalization, it does not improve the confidence calibration (Wortsman et al., 2022). In contrast to WA, our weight & mask averaging uses masks for weighted averaging and improves the confidence calibration with similar FLOPs in the prediction.

## 5 EXPERIMENTS

We perform a comprehensive empirical evaluation of CigL, comparing it with the popular baseline method RigL (Evci et al., 2020). RigL is a popular sparse training method that uses weights magnitudes to prune and gradient magnitudes to grow connections.

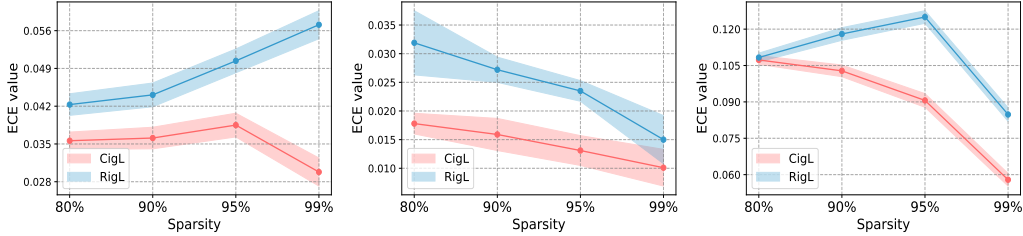
**Datasets & Model Architectures:** We follow the settings in Evci et al. (2020) for a comprehensive comparison. Our experiments are based on three benchmark datasets: CIFAR-10 and CIFAR-100 (Krizhevsky et al., 2009) and ImageNet-2012 (Russakovsky et al., 2015). For model architectures, we used ResNet-50 (He et al., 2016) and Wide-ResNet-22-2 (Zagoruyko & Komodakis, 2016). We repeat all experiments 3 times and report the mean and standard deviation.

**Sparse Training Settings:** For sparse training, we check multiple sparsities, including 80%, 90%, 95%, and 99%, which can sufficiently reduce the memory requirement and is of more interest.

**Implementations:** We follow the settings in (Evci et al., 2020; Sundar & Dwaraknath, 2021). The parameters are optimized by SGD with momentum. For the learning rate, we use piecewise constant decay scheduler. For CIFAR-10 and CIFAR-100, we train all the models for 250 epochs with a batch size of 128. For ImageNet, we train all the models for 100 epochs with a batch size of 64.

### 5.1 COMPARISON BETWEEN POPULAR SPARSE TRAINING METHOD

**Results on CIFAR-10 and CIFAR-100.** We first compare our CigL and RigL by the expected calibration error (ECE) (Guo et al., 2017), a popular measure of the discrepancy between a model’s



(a) CIFAR-10, ResNet-50 (b) CIFAR-10, Wide-ResNet-22-2 (c) CIFAR-100, ResNet-50

Figure 2: ECE value comparison between CigL and RigL at different sparsities (80%, 90%, 95%, 99%). Compared to RigL, CigL produces sparse models with smaller ECE values.

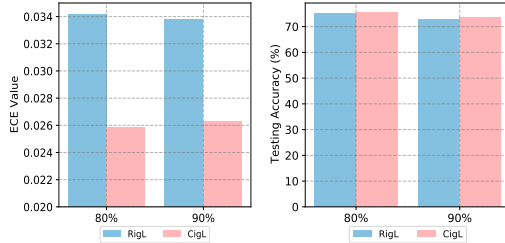
Table 1: Testing accuracy (%) comparison between CigL and RigL at different sparsities (80%, 90%, 95%, 99%). Compared to RigL, CigL maintains comparable or higher test accuracy.

	CIFAR-10		CIFAR-100	
	RiGL	CiGL	RiGL	CiGL
80% SPARSITY	94.02 (0.115)	<b>94.75</b> (0.107)	72.08 (0.109)	<b>76.84</b> (0.089)
90% SPARSITY	93.84 (0.184)	<b>94.56</b> (0.189)	71.90 (0.172)	<b>76.24</b> (0.181)
95% SPARSITY	93.19 (0.198)	<b>94.20</b> (0.202)	70.90 (0.210)	<b>74.71</b> (0.197)
99% SPARSITY	91.31 (0.205)	<b>92.42</b> (0.196)	65.57 (0.208)	<b>66.42</b> (0.206)

confidence and true accuracy, with a lower ECE indicating better confidence calibration and higher reliability. In Figure 2, the pink and blue curves represent CigL and RigL, respectively, where the colored are represent the 95% confidence intervals. We can see that the pink curves are usually lower than the blue curves for different sparsities (80%, 90%, 95%, 99%), which implies that our CigL can reduce the ECE and improve the confidence calibration of the produced sparse models.

Apart from ECE value, we also compare our CigL and RigL by the testing accuracy for multiple sparsities (80%, 90%, 95%, 99%). We summarize the results for sparse ResNet-50 in Table 1. It is observed that CigL tends to bring comparable or higher accuracy, which demonstrates that CigL can simultaneously maintain or improve the accuracy.

**Results on ImageNet-2012.** We also compare the ECE values and test accuracy of our CigL and RigL on a larger dataset, ImageNet-2012, where the sparsity of ResNet-50 is 80% and 90%. As shown in Figure 3, the pink and blue bars represent our CigL and RigL, respectively. For the comparison of ECE values in (a), the pink bars are shorter than the blue bars, indicating an improved reliability of the sparse model produced by CigL. For the test accuracy comparison in (b), the pink and blue bars are very similar in height, implying that the accuracy of CigL is comparable to that of RigL.



(a) ECE value (RM) (b) Test accuracy (RM)

Figure 3: ECE value and test accuracy(%) of CigL and RigL at 80% & 90% sparsities on ImageNet-2012. Compared with RigL, CigL has smaller ECE values and comparable test accuracies.

## 5.2 COMPARISON BETWEEN DIFFERENT DROPOUT METHODS

In this section, since our CigL is related to dropout methods, we compare our CigL with RigL using existing popular dropout methods, namely weight dropout (W-DP) and MC dropout (MC-DP). **The comparison of test accuracy** is shown in Table 2. Our CigL usually provides a comparable or higher accuracy compared to RigL. However, using weight dropout and MC dropout in RigL usually result in a decrease in accuracy. We also summarize **the comparison of the ECE value** between CigL and different dropout methods in Table 3. For each sparsity and architecture, we have marked in bold those cases where the ECE value is significantly reduced ( $\geq 15\%$  reduction compared to RigL). Our CigL are always bolded, indicating its ability to reduce ECE value and increase reliability in sparse training. But RigL + weight dropout does not significantly reduce ECE values in almost all cases and RigL + MC dropout also does not improve the calibration in highly sparse cases (99% sparsity).

Table 2: Testing accuracy (%) comparison between CigL, RigL + weight dropout (W-DP), and RigL + MC dropout (MC-DP) at different sparsities (80%, 90%, 95%, 99%). Compared to RigL, RigL + W-DP, and RigL + MC-DP, CigL maintains comparable or higher test accuracy.

		80% SPARSITY	90% SPARSITY	95% SPARSITY	99% SPARSITY
RESNET-50	RiGL	94.02 (0.115)	93.84 (0.184)	93.19 (0.198)	91.31 (0.205)
	RiGL + W-DP	93.26 (0.114)	93.47 (0.186)	92.71 (0.193)	89.99 (0.210)
	RiGL + MC-DP	93.39 (0.105)	93.71 (0.181)	92.87 (0.205)	89.84 (0.212)
	CiGL	<b>94.75</b> (0.107)	<b>94.56</b> (0.189)	<b>94.20</b> (0.202)	<b>92.42</b> (0.196)
WRN-22-2	RiGL	93.12 (0.188)	92.26 (0.187)	91.02 (0.179)	83.82 (0.224)
	RiGL + W-DP	91.77 (0.182)	91.44 (0.191)	89.66 (0.183)	80.42 (0.215)
	RiGL + MC-DP	91.75 (0.149)	91.49 (0.187)	89.39 (0.177)	77.48 (0.198)
	CiGL	<b>93.95</b> (0.088)	<b>93.05</b> (0.219)	<b>91.34</b> (0.171)	<b>83.96</b> (0.189)

Table 3: Testing ECE comparison between CigL, RigL + weight dropout (W-DP), and RigL + MC dropout (MC-DP) at different sparsities (80%, 90%, 95%, 99%). Compared to RigL + W-DP and RigL + MC-DP, CigL more consistently achieves a significant reduction in the ECE value of RigL.

		80% SPARSITY	90% SPARSITY	95% SPARSITY	99% SPARSITY
RESNET-50	RiGL	0.0423 (0.001)	0.0441 (0.001)	0.0504 (0.001)	0.0571 (0.001)
	RiGL + W-DP	0.0504 (0.002)	0.0438 (0.001)	0.0462 (0.002)	<b>0.0315</b> (0.002)
	RiGL + MC-DP	<b>0.0322</b> (0.001)	<b>0.0200</b> (0.001)	<b>0.0121</b> (0.001)	0.0528 (0.002)
	CiGL	<b>0.0356</b> (0.001)	<b>0.0361</b> (0.001)	<b>0.0385</b> (0.001)	<b>0.0298</b> (0.001)
WRN-22-2	RiGL	0.0319 (0.003)	0.0272 (0.001)	0.0235 (0.001)	0.0150 (0.002)
	RiGL + W-DP	0.0433 (0.003)	0.0348 (0.002)	0.0256 (0.002)	0.0174 (0.003)
	RiGL + MC-DP	<b>0.0159</b> (0.001)	<b>0.0077</b> (0.002)	0.0384 (0.001)	0.1502 (0.002)
	CiGL	<b>0.0178</b> (0.001)	<b>0.0159</b> (0.001)	<b>0.0131</b> (0.001)	<b>0.0101</b> (0.002)

### 5.3 COMPARISON BETWEEN OTHER CALIBRATION METHODS

In this section, we compare our CigL with existing popular calibration methods, including mixup (Zhang et al., 2017), temperature scaling (TS) (Guo et al., 2017), and label smoothing (LS) (Szegedy et al., 2016). **The testing ECE** are depicted in Figure 4, where the pink and blue polygons represent CigL and other calibration methods, respectively. We can see that CigL usually gives smaller polygons, indicating a better confidence calibration.

### 5.4 ABLATION STUDIES

We do ablation studies to demonstrate the importance of each component in our CigL, where we train sparse networks using our CigL without random masks (CigL w/o RM) and CigL without weight & mask averaging (CigL w/o WMA), respectively. In CigL w/o RM, we search for sparse topologies using only the deterministic mask. In CigL w/o WMA, we collect multiple model samples and use prediction averaging during testing. Figures 5(a)-(b) show **the effect of random masks** on the test accuracy and ECE values, where the blue, green, and pink bars represent RigL, CigL w/o RM, and CigL, respectively. We can see that if we remove the random mask, we can still obtain an improvement in accuracy compared to RigL. However, the ECE values do not decrease as much as CigL, indicating that the CigL w/o RM is not as effective as CigL in improving the confidence calibration. Figures 5(c)-(d) further show **the effect of weight & mask averaging**. We can see that without using weight & mask averaging, the accuracy decreases and the ECE value increases in high sparsity such as 95% and 99%, demonstrating the importance of weight & mask averaging.

## 6 DISCUSSION & CONCLUSION

We for the first time identify and study the reliability problem of sparse training and find that sparse training exacerbates the over-confidence problem of DNNs. We then develop a new sparse training



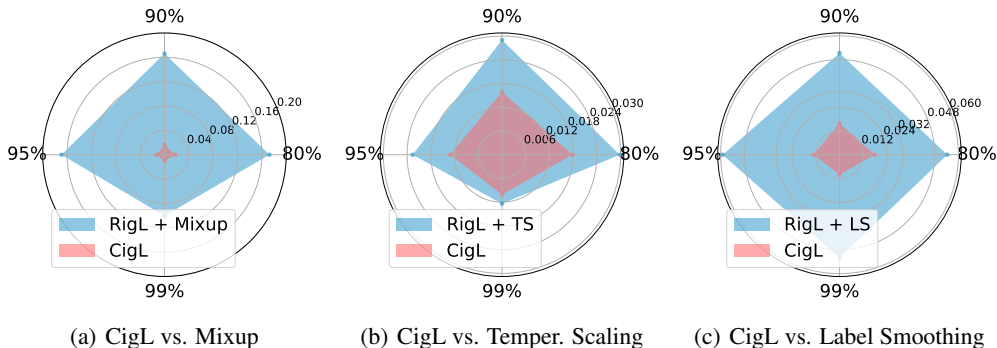


Figure 4: ECE value comparison between CigL and RigL + other calibration methods at different sparsities (80%, 90%, 95%, 99%). The pink polygons (CigL) are smaller than the blue polygons (other calibration methods), indicating a better confidence calibration using CigL compared to (a) Mixup, (b) Temperature scaling, and (c) Label smoothing.

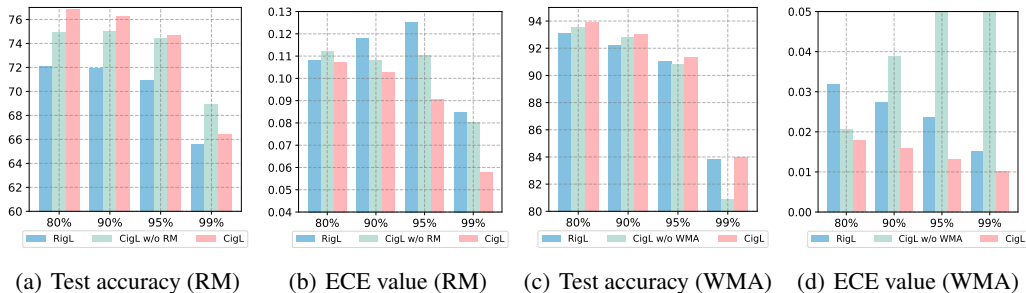


Figure 5: Ablation studies: test accuracy(%) and ECE value comparison between CigL, CigL with random mask (CigL w/o RM), and CigL without weight & mask averaging (CigL w/o WMA) at different sparsities (80%, 90%, 95%, 99%). Compared to (a)-(b) CigL w/o RM and (c)-(d) CigL w/o WMA, CigL more consistently produces sparse models with low ECE values and high accuracy.

method, CigL, to produce reliable sparse models, which can simultaneously maintain or even improve accuracy with only a slight increase in computational and storage burden. Our CigL utilizes two masks, including a deterministic mask and a random mask, which allows the sparse model to better explore the weight space. Then, we design weight & mask averaging method to combine multiple sparse weights and random masks into a single model with improved reliability. We prove that CigL can be viewed as a hierarchical variational approximation to the probabilistic deep Gaussian process. Experiments results on multiple benchmark datasets, model architectures, and sparsities show that our CigL reduces ECE values by up to **47.8%** with comparable or higher accuracy.

One phenomenon we find worth discussing is the *double descent* in reliability of sparse training. Nakkiran et al. (2021) first observed this double descent phenomenon in DNNs, where as the model size, data size, or training time increases, the performance of the model first improves, then gets worse, and then improves again. Consistent with the previous definition, we consider sparsity and reliability as the measures of model size and performance, respectively. Then, as shown in the Figure 1 (c), as the sparsity decreases (model size increases), the reliability (model performance) gets better, then gets worse, and then gets better again. To explain this phenomenon, we divided sparsity into four phases, from the left (99.9%) to the right (0%), by drawing an analogy between the phases and model accuracy and size. **(a) The sparse model starts as a poor model**, which is too sparse to learn the data well (low reliability & accuracy). **(b) It gradually becomes equivalent to a shallow model** that can learn some patterns but is not flexible enough to learn all the data well (high reliability & moderate level of accuracy). **(c) Then, it moves to a sparse deep model** that can accommodate complex patterns but suffers from poor exploration (low reliability & high accuracy). **(d) Finally, it reaches a dense deep model** with over-confidence issues (moderate level of reliability & high accuracy). It is observed that at around 95% sparsity, the sparse model can achieve comparable accuracy and high sparsity at the same time, which makes it important in practical applications. However, the ECE value is at the peak of the double-descent curve at this point, which implies that the reliability of the sparse model is at a low level. Thus, our CigL smooths the double descent curve and produce reliable models on those important high sparsity levels.

## ACKNOWLEDGMENTS

This research was partially supported by NSF Grant No. NSF CCF-1934904 (TRIPODS).

## REPRODUCIBILITY STATEMENT

The implementation code can be found in <https://github.com/StevenBoys/CigL>. All datasets and code platform (PyTorch) we use are public.

## REFERENCES

- Arsenii Ashukha, Alexander Lyzhov, Dmitry Molchanov, and Dmitry Vetrov. Pitfalls of in-domain uncertainty estimation and ensembling in deep learning. *arXiv preprint arXiv:2002.06470*, 2020.
- Guillaume Bellec, David Kappel, Wolfgang Maass, and Robert Legenstein. Deep rewiring: Training very sparse deep networks. *International Conference on Learning Representations (ICLR)*, 2018.
- Sameer Bibikar, Haris Vikalo, Zhangyang Wang, and Xiaohan Chen. Federated dynamic sparse training: Computing less, communicating less, yet learning better. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pp. 6080–6088, 2022.
- Christopher M Bishop and Nasser M Nasrabadi. *Pattern recognition and machine learning*, volume 4. Springer, 2006.
- Mariusz Bojarski, Davide Del Testa, Daniel Dworakowski, Bernhard Firner, Beat Flepp, Praseoon Goyal, Lawrence D Jackel, Mathew Monfort, Urs Muller, Jiakai Zhang, et al. End to end learning for self-driving cars. *arXiv preprint arXiv:1604.07316*, 2016.
- Tianlong Chen, Zhenyu Zhang, Jun Wu, Randy Huang, Sijia Liu, Shiyu Chang, and Zhangyang Wang. Can you win everything with a lottery ticket? *Transactions of Machine Learning Research*, 2022.
- Andreas Damianou and Neil D Lawrence. Deep gaussian processes. In *Artificial intelligence and statistics*, pp. 207–215. PMLR, 2013.
- Tim Dettmers and Luke Zettlemoyer. Sparse networks from scratch: Faster training without losing performance. *arXiv preprint arXiv:1907.04840*, 2019.
- Anastasia Dietrich, Frithjof Gressmann, Douglas Orr, Ivan Chelombiev, Daniel Justus, and Carlo Luschi. Towards structured dynamic sparse pre-training of bert. *arXiv preprint arXiv:2108.06277*, 2021.
- Utku Evci, Fabian Pedregosa, Aidan Gomez, and Erich Elsen. The difficulty of training sparse neural networks. *arXiv preprint arXiv:1906.10732*, 2019.
- Utku Evci, Trevor Gale, Jacob Menick, Pablo Samuel Castro, and Erich Elsen. Rigging the lottery: Making all tickets winners. In *International Conference on Machine Learning*, pp. 2943–2952. PMLR, 2020.
- Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. *International Conference on Learning Representations (ICLR)*, 2019.
- Yarin Gal and Zoubin Ghahramani. Dropout as a bayesian approximation: Representing model uncertainty in deep learning. In *international conference on machine learning*, pp. 1050–1059. PMLR, 2016.
- Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.
- Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 770–778, 2016.

- Zheng He, Zeke Xie, Quanzhi Zhu, and Zengchang Qin. Sparse double descent: Where network pruning aggravates overfitting. In *International Conference on Machine Learning*, pp. 8635–8659. PMLR, 2022.
- Geoffrey E Hinton, Nitish Srivastava, Alex Krizhevsky, Ilya Sutskever, and Ruslan R Salakhutdinov. Improving neural networks by preventing co-adaptation of feature detectors. *arXiv preprint arXiv:1207.0580*, 2012.
- Shaoyi Huang, Bowen Lei, Dongkuan Xu, Hongwu Peng, Yue Sun, Mimi Xie, and Caiwen Ding. Dynamic sparse training via balancing the exploration-exploitation trade-off. *arXiv preprint arXiv:2211.16667*, 2022.
- Pavel Izmailov, Dmitrii Podoprikin, Timur Garipov, Dmitry Vetrov, and Andrew Gordon Wilson. Averaging weights leads to wider optima and better generalization. *arXiv preprint arXiv:1803.05407*, 2018.
- Siddhant Jayakumar, Razvan Pascanu, Jack Rae, Simon Osindero, and Erich Elsen. Top-kast: Top-k always sparse training. *Advances in Neural Information Processing Systems*, 33:20744–20754, 2020.
- Xiaoqian Jiang, Melanie Osl, Jihoon Kim, and Lucila Ohno-Machado. Calibrating predictive model estimates to support personalized medicine. *Journal of the American Medical Informatics Association*, 19(2):263–274, 2012.
- Alex Krizhevsky, Geoffrey Hinton, et al. Learning multiple layers of features from tiny images. *Technical report, University of Toronto*, 2009.
- Shiwei Liu, Lu Yin, Decebal Constantin Mocanu, and Mykola Pechenizkiy. Do we actually need dense over-parameterization? in-time over-parameterization in sparse training. In *International Conference on Machine Learning*, pp. 6989–7000. PMLR, 2021.
- Wesley J Maddox, Pavel Izmailov, Timur Garipov, Dmitry P Vetrov, and Andrew Gordon Wilson. A simple baseline for bayesian uncertainty in deep learning. *Advances in Neural Information Processing Systems*, 32, 2019.
- Decebal Constantin Mocanu, Elena Mocanu, Peter Stone, Phuong H Nguyen, Madeleine Gibescu, and Antonio Liotta. Scalable training of artificial neural networks with adaptive sparse connectivity inspired by network science. *Nature communications*, 9(1):1–12, 2018.
- Hesham Mostafa and Xin Wang. Parameter efficient training of deep convolutional neural networks by dynamic sparse reparameterization. In *International Conference on Machine Learning*, pp. 4646–4655. PMLR, 2019.
- Preetum Nakkiran, Gal Kaplun, Yamini Bansal, Tristan Yang, Boaz Barak, and Ilya Sutskever. Deep double descent: Where bigger models and more data hurt. *Journal of Statistical Mechanics: Theory and Experiment*, 2021(12):124003, 2021.
- Jeremy Nixon, Michael W Dusenberry, Linchuan Zhang, Ghassen Jerfel, and Dustin Tran. Measuring calibration in deep learning. In *CVPR Workshops*, volume 2, 2019.
- Ozan Özdenizci and Robert Legenstein. Training adversarially robust sparse networks via bayesian connectivity sampling. In *International Conference on Machine Learning*, pp. 8314–8324. PMLR, 2021.
- Kanil Patel, William Beluch, Kilian Rambach, Michael Pfeiffer, and Bin Yang. Improving uncertainty of deep learning-based object classification on radar spectra using label smoothing. In *2022 IEEE Radar Conference (RadarConf22)*, pp. 1–6. IEEE, 2022.
- Rahul Rahaman et al. Uncertainty quantification and deep ensembles. *Advances in Neural Information Processing Systems*, 34:20063–20075, 2021.
- Rajesh Ranganath, Dustin Tran, and David Blei. Hierarchical variational models. In *International conference on machine learning*, pp. 324–333. PMLR, 2016.

- Olga Russakovsky, Jia Deng, Hao Su, Jonathan Krause, Sanjeev Satheesh, Sean Ma, Zhiheng Huang, Andrej Karpathy, Aditya Khosla, Michael Bernstein, et al. Imagenet large scale visual recognition challenge. *International journal of computer vision*, 115(3):211–252, 2015.
- Jonathan Schwarz, Siddhant Jayakumar, Razvan Pascanu, Peter E Latham, and Yee Teh. Power-propagation: A sparsity inducing weight reparameterisation. *Advances in Neural Information Processing Systems*, 34:28889–28903, 2021.
- Gowthami Somepalli, Liam Fowl, Arpit Bansal, Ping Yeh-Chiang, Yehuda Dar, Richard Baraniuk, Micah Goldblum, and Tom Goldstein. Can neural nets learn the same model twice? investigating reproducibility and double descent from the decision boundary perspective. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 13699–13708, 2022.
- Nitish Srivastava, Geoffrey Hinton, Alex Krizhevsky, Ilya Sutskever, and Ruslan Salakhutdinov. Dropout: a simple way to prevent neural networks from overfitting. *The journal of machine learning research*, 15(1):1929–1958, 2014.
- Yiyou Sun and Yixuan Li. On the effectiveness of sparsification for detecting the deep unknowns. *arXiv preprint arXiv:2111.09805*, 2021.
- Varun Sundar and Rajat Vadiraj Dwaraknath. [reproducibility report] rigging the lottery: Making all tickets winners. *arXiv preprint arXiv:2103.15767*, 2021.
- Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 2818–2826, 2016.
- Bindya Venkatesh, Jayaraman J Thiagarajan, Kowshik Thopalli, and Prasanna Sattigeri. Calibrate and prune: Improving reliability of lottery tickets through prediction calibration. *arXiv preprint arXiv:2002.03875*, 2020.
- Li Wan, Matthew Zeiler, Sixin Zhang, Yann Le Cun, and Rob Fergus. Regularization of neural networks using dropconnect. In *International conference on machine learning*, pp. 1058–1066. PMLR, 2013.
- Yezhen Wang, Bo Li, Tong Che, Kaiyang Zhou, Ziwei Liu, and Dongsheng Li. Energy-based open-world uncertainty modeling for confidence calibration. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 9302–9311, 2021.
- Mitchell Wortsman, Gabriel Ilharco, Samir Ya Gadre, Rebecca Roelofs, Raphael Gontijo-Lopes, Ari S Morcos, Hongseok Namkoong, Ali Farhadi, Yair Carmon, Simon Kornblith, et al. Model soups: averaging weights of multiple fine-tuned models improves accuracy without increasing inference time. In *International Conference on Machine Learning*, pp. 23965–23998. PMLR, 2022.
- Mengzhou Xia, Zexuan Zhong, and Danqi Chen. Structured pruning learns compact and accurate models. *arXiv preprint arXiv:2204.00408*, 2022.
- Lu Yin, Vlado Menkovski, Meng Fang, Tianjin Huang, Yulong Pei, Mykola Pechenizkiy, Decbal Constantin Mocanu, and Shiwei Liu. Superposing many tickets into one: A performance booster for sparse neural network training. *arXiv preprint arXiv:2205.15322*, 2022.
- Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *British Machine Vision Conference*, 2016.
- Hongyi Zhang, Moustapha Cisse, Yann N Dauphin, and David Lopez-Paz. mixup: Beyond empirical risk minimization. *arXiv preprint arXiv:1710.09412*, 2017.
- Jize Zhang, Bhavya Kailkhura, and T Yong-Jin Han. Mix-n-match: Ensemble and compositional methods for uncertainty calibration in deep learning. In *International conference on machine learning*, pp. 11117–11128. PMLR, 2020.
- Xiao Zhou, Weizhong Zhang, Hang Xu, and Tong Zhang. Effective sparsification of neural networks with global sparsity constraint. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 3599–3608, 2021.

## A APPENDIX: BACKGROUND

In this section, we briefly summarize CigL, Gaussian processes, and hierarchical variational inference, which will be used to support the main theoretical analysis of this work.

### A.1 SPARSE TRAINING: CIGL

We first review our CigL method for the case of a single hidden layer neural network (NN). This is done for ease of notation, and it is straightforward to generalise to multiple layers (Gal & Ghahramani, 2016). Denote by  $\mathbf{W}_1, \mathbf{W}_2$  the sparse weight matrices connecting the first layer to the hidden layer and connecting the hidden layer to the output layer respectively. For the sparse mask controlling the sparse topology, we use  $\mathbf{M}_1, \mathbf{M}_2$  to denote the corresponding deterministic masks for  $\mathbf{W}_1$  and  $\mathbf{W}_2$ , respectively. And we use  $\mathbf{Z}_1, \mathbf{Z}_2$  to denote the corresponding random masks for  $\mathbf{W}_1$  and  $\mathbf{W}_2$ , respectively. These linearly transform the layers' inputs before applying some element-wise non-linearity  $\sigma(\cdot)$ . Denote by  $\mathbf{b}$  the biases by which we shift the input of the non-linearity. We assume the model to output  $D$  dimensional vectors while its input is  $Q$  dimensional vectors, with  $K$  hidden units. Thus  $\mathbf{W}_1, \mathbf{M}_1$ , and  $\mathbf{Z}_1$  are  $Q \times K$  matrices,  $\mathbf{W}_2, \mathbf{M}_2$ , and  $\mathbf{Z}_2$  are  $K \times D$  matrices, and  $\mathbf{b}$  is a  $K$  dimensional vector. A sparse NN model with the two masks would output  $\hat{\mathbf{y}} = \sigma(\mathbf{x}\mathbf{Z}_1 \odot \mathbf{W}_1 + \mathbf{b})\mathbf{Z}_2 \odot \mathbf{W}_2$  given some input  $\mathbf{x}$ .

For the update of masks, the deterministic masks is updated by exploiting the magnitude of weights and gradients. The random mask is randomly sampled. For the update of weights, we use  $E$  to denote the loss function, which is the euclidean loss for regression problem

$$E = \frac{1}{2N} \sum_{n=1}^N N \|y_n - \hat{y}_n\|_2^2, \quad (6)$$

where  $y_n$  is the observed response,  $\hat{y}_n$  is the prediction based on input  $\mathbf{x}_n$  for  $n = 1, \dots, N$ .

For classification task with  $D$  classes, we use softmax function to map the output  $\hat{y}_n$  to the probability score for each class  $\hat{p}_{nd} = \exp(\hat{y}_{nd}) / (\sum_k \exp(\hat{y}_{nk}))$ , and the loss function will be

$$E = -\frac{1}{N} \sum_{n=1}^N N \log(\hat{p}_{n,c_n}), \quad (7)$$

where  $c_n \in [1, 2, \dots, D]$  is the true class label for  $\mathbf{x}_n$ .

During NN optimization process, apart from the loss function mentioned above,  $l_2$  regularisation is often used to improve the performance, leading to a minimisation objective:

$$\mathcal{L}_{\text{CigL}} := \frac{1}{N} \sum_{i=1}^N E(y_i, \hat{y}_i) + \lambda_1 \|\mathbf{W}_1\|_2^2 + \lambda_2 \|\mathbf{W}_2\|_2^2 + \lambda_3 \|\mathbf{b}\|_2^2. \quad (8)$$

Therefore, during the training process of CigL, we iteratively update the sparse weights  $\mathbf{W}$  by minimising Eq. (8) and update deterministic mask based on the magnitude of weights and gradients.

### A.2 GAUSSIAN PROCESS

The Gaussian process (GP) is a popular non-parametric Bayesian method to model distributions over functions, which can be applied to both regression and classification tasks. It has good performance in various fields, but it will bring huge computation burden when faced with a large number of data. The use of variational inference for GP can make it scalable to large data.

Given a data  $\{\mathbf{x}_n, \mathbf{y}_n\}, n = 1, \dots, N$ , the task is to estimate an unknown function  $\mathbf{y} = f(\mathbf{x})$ , where  $\mathbf{X} \in \mathbb{R}^{N \times Q}$  and  $\mathbf{Y} \in \mathbb{R}^{N \times D}$ . GP usually put a prior over the function space and we want to fit the posterior distribution over the function space:

$$p(\mathbf{f}|\mathbf{X}, \mathbf{Y}) \propto p(\mathbf{Y}|\mathbf{X}, \mathbf{f})p(\mathbf{f}).$$

Within Gaussian process, we usually place a Gaussian prior over the function space, and it equivalent to placing a joint Gaussian distribution over all function values

$$\begin{aligned} \mathbf{F}|\mathbf{X} &\sim \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{X})) \\ \mathbf{Y}|\mathbf{F} &\sim \mathcal{N}(\mathbf{F}, \tau^{-1}) \end{aligned} \quad (9)$$

where  $\tau$  is a precision parameter and  $\mathbf{I}_N$  is the identity matrix with dimensions  $N \times N$ .

For classification tasks, we can formulate the model as

$$\mathbf{F}|\mathbf{X} \sim \mathcal{N}(\mathbf{0}, \mathbf{K}(\mathbf{X}, \mathbf{X})) \quad (10)$$

$$\mathbf{Y}|\mathbf{F} \sim \mathcal{N}(\mathbf{F}, \tau^{-1}),$$

$$c_n|\mathbf{Y} \sim \text{Categorical}\left(\frac{\exp(\hat{y}_{nd})}{\sum_k \exp(\hat{y}_{nk})}\right) \quad (11)$$

An important aspect of Gaussian process is the choice of the covariance function, which reflects how we believe the similarity between each pair of inputs  $\mathbf{x}_i$  and  $\mathbf{x}_j$ . One widely-used covariance function is stationary squared exponential covariance function. In addition, some non-stationary covariance function are proposed, such as dot-product kernels and more flexible deep network kernels.

### A.3 HIERARCHICAL VARIATIONAL INFERENCE

Variational inference (VI) is a broadly-used technique to approximate intractable integrals in Bayesian modeling, which sets up a parameterized family of tractable distributions over the latent variables and then optimizes the parameters to be close to the posterior. More specifically, suppose  $\mathbf{w}$  is the set of random variables defining our model. Then, the predictive distribution will be formulated as

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{w})p(\mathbf{w}|\mathbf{X}, \mathbf{Y})d\mathbf{w},$$

where the posterior  $p(\mathbf{w}|\mathbf{X}, \mathbf{Y})$  is usually intractable. Thus, we define a family of tractable approximating variational distributions  $q(\mathbf{w})$  to approach the posterior.

To find the closest approximating distribution among the family of  $q(\mathbf{w})$ , we minimise the Kullback–Leibler (KL) divergence between  $q(\mathbf{w})$  and posterior  $p(\mathbf{w}|\mathbf{X}, \mathbf{Y})$ , which is equivalent to maximising the log evidence lower bound (ELBO) with respect to  $q(\mathbf{w})$ :

$$\mathcal{L}_{\text{VI}} := \int q(\mathbf{w}) \log p(\mathbf{Y}|\mathbf{X}, \mathbf{w})d\mathbf{w} - D_{\text{KL}}(q(\mathbf{w})||p(\mathbf{w})).$$

After obtaining an good approximation  $q(\mathbf{w})$ , we can update the predictive distribution to

$$p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{X}, \mathbf{Y}) = \int p(\mathbf{y}^*|\mathbf{x}^*, \mathbf{w})q(\mathbf{w})d\mathbf{w}.$$

However, when faced with posterior difficult to fit,  $q(\mathbf{w})$  can be limited and not flexible enough to approach the posterior. In this case, VI cannot capture the posterior dependencies between latent variables that both improve the fidelity of the approximation and are sometimes intrinsically meaningful. To solve this limitation of VI, hierarchical variational inference (HVI) is proposed, which can capture both posterior dependencies between the latent variables and more complex marginal distributions (Ranganath et al., 2016). More specifically about HVI, we extend the limited family of VI distribution hierarchically, i.e., by placing a prior on the parameters of the likelihood. Suppose VI uses  $q(\mathbf{w}; \lambda)$  to approximate the posterior where  $\lambda$  is the variational parameters to optimise. HVI will add prior on  $\lambda$  and uses  $q(\mathbf{w}|\lambda)q(\lambda; \theta)$ . The ELBO equivalently will be as:

$$\mathcal{L}_{\text{HVI}} := \mathbb{E}_{q_{\text{HVI}}(\mathbf{w}; \lambda)}[\log p(\mathbf{x}, \mathbf{w}) - \log q_{\text{HVI}}(\mathbf{w}; \theta)].$$

This ELBO can be further bounded by

$$\mathcal{L}_{\text{HVI}} \leq \mathbb{E}_{q_{\text{HVI}}(\mathbf{w}; \lambda)}[\log p(\mathbf{x}, \mathbf{w})] - \mathbb{E}_{q(\mathbf{w}, \lambda)}[\log q(\lambda) + \log q(\mathbf{w}|\lambda) - \log r(\lambda|\mathbf{w}; \theta)].$$

where  $r(\lambda|\mathbf{w}; \theta)$  is introduced to apply the variational principle.

## B APPENDIX: CIGL AS A HIERARCHICAL BAYESIAN APPROXIMATION

We show that sparse deep NNs trained with CigL are mathematically equivalent to approximate hierarchical variational inference in the deep Gaussian process (marginalised over its covariance

function parameters). For this, we build on previous work (Gal & Ghahramani, 2016) that proved unit dropout applied before every weight layer are mathematically equivalent to approximate variational inference in the deep Gaussian process. Starting with the full Gaussian process we will develop an approximation that will be shown to be equivalent to the sparse NN optimisation objective with CigL (eq. (8)) with either the Euclidean loss in the case of regression or softmax loss in the case of classification. Our derivation takes regression as an example, which can be extended to classification by Section 4 of the Appendix of Gal & Ghahramani (2016). This view of CigL will allow us to derive new probabilistic results in sparse training.

### B.1 A GAUSSIAN PROCESS APPROXIMATION

In this section, we will re-parameterise the deep GP model and marginalise over the additional auxiliary random variables, which is built on Gal & Ghahramani (2016). To define our covariance function, let  $\sigma(\cdot)$  be some non-linear activation function and  $\mathbf{K}(\mathbf{x}, \mathbf{y})$  can be formulated as:

$$\mathbf{K}(\mathbf{x}, \mathbf{y}) = \int p(\mathbf{w})p(\mathbf{b})\sigma(\mathbf{w}^\top \mathbf{x} + \mathbf{b})\sigma(\mathbf{w}^\top \mathbf{y} + \mathbf{b})d\mathbf{w}d\mathbf{b},$$

where  $p(\mathbf{w})$  is a standard multivariate normal distribution in dimension  $Q$ .

We use Monte Carlo integration with  $K$  samples to approximate the integral above and get the finite rank covariance function

$$\widehat{\mathbf{K}}(\mathbf{x}, \mathbf{y}) = \frac{1}{K} \sum_{k=1}^K \sigma(\mathbf{w}_k^\top \mathbf{x} + \mathbf{b}_k)\sigma(\mathbf{w}_k^\top \mathbf{y} + \mathbf{b}_k),$$

where  $\mathbf{w}_k \sim p(\mathbf{w})$  and  $\mathbf{b}_k \sim p(\mathbf{b})$ .  $K$  is the number of hidden units in our single hidden layer sparse NN approximation. The generative model will be as follow when we use  $\widehat{\mathbf{K}}$  instead of  $\mathbf{K}$ :

$$\begin{aligned} \mathbf{w} &\sim p(\mathbf{w}), \quad \mathbf{b}_k \sim p(\mathbf{b}), \\ \mathbf{W}_1 &= [\mathbf{w}_{qk}]_{q=1, k=1}^Q, \quad \mathbf{b} = [\mathbf{b}_k]_{k=1}^K, \\ \widehat{\mathbf{K}}(\mathbf{x}, \mathbf{y}) &= \frac{1}{K} \sum_{k=1}^K \sigma(\mathbf{w}_k^\top \mathbf{x} + \mathbf{b}_k)\sigma(\mathbf{w}_k^\top \mathbf{y} + \mathbf{b}_k), \\ \mathbf{F}|\mathbf{X}, \mathbf{W}_1, \mathbf{b} &\sim \mathcal{N}(0, \widehat{\mathbf{K}}(\mathbf{x}, \mathbf{y})), \\ \mathbf{Y}|\mathbf{F} &\sim \mathcal{N}(\mathbf{F}, \tau^{-1}\mathbf{I}_n), \end{aligned}$$

where  $\mathbf{W}_1 \in \mathbb{R}^{Q \times K}$  which parameterise the covariance function.

We can get the predictive distribution by integrating over the  $\mathbf{F}$ ,  $\mathbf{W}_1$ , and  $\mathbf{b}$

$$p(\mathbf{Y}|\mathbf{X}) = \int p(\mathbf{Y}|\mathbf{F})p(\mathbf{F}|\mathbf{W}_1, \mathbf{b}, \mathbf{X})p(\mathbf{W}_1)p(\mathbf{b})d\mathbf{W}_1d\mathbf{b}.$$

Denoting a  $1 \times K$  row vector

$$\phi(\mathbf{x}, \mathbf{W}_1, \mathbf{b}) = \sqrt{\frac{1}{K}}\sigma(\mathbf{W}_1^\top \mathbf{x} + \mathbf{b})$$

and a  $N \times K$  feature matrix  $\Phi = [\phi(\mathbf{x}_n, \mathbf{W}_1, \mathbf{b})]_{n=1}^N$ . Then, we can get  $\widehat{\mathbf{K}}(\mathbf{X}, \mathbf{X}) = \Phi\Phi^\top$  and the predictive distribution can be rewritten as

$$p(\mathbf{Y}|\mathbf{X}) = \int \mathcal{N}(\mathbf{Y}; 0, \Phi\Phi^\top + \tau^{-1}\mathbf{I}_N)p(\mathbf{W}_1)p(\mathbf{b})d\mathbf{W}_1d\mathbf{b}$$

The normal distribution of  $\mathbf{Y}$  inside the integral above can be written as a joint normal distribution over  $\mathbf{y}_d$  which denoting the  $d$ -th columns of the  $N \times D$  matrix  $\mathbf{Y}$  ( $d = 1, \dots, D$ ). For each term in the joint distribution, following Bishop & Nasrabadi (2006), we introduce a  $K \times 1$  auxiliary random variable  $\mathbf{w}_d \sim \mathcal{N}(0, \mathbf{I}_K)$ ,

$$\mathcal{N}(\mathbf{y}_d; 0, \Phi\Phi^\top + \tau^{-1}\mathbf{I}_N) = \int \mathcal{N}(\mathbf{y}_d; \Phi\mathbf{w}_d, \tau^{-1}\mathbf{I}_N)\mathcal{N}(\mathbf{w}_d; 0, \mathbf{I}_K)d\mathbf{w}_d.$$

We use  $\mathbf{W}_2 = [\mathbf{w}_d]_{d=1}^D \in \mathbb{R}^{K \times D}$  and we get the predictive distribution as

$$p(\mathbf{Y}|\mathbf{X}) = \int p(\mathbf{Y}|\mathbf{X}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{b})p(\mathbf{W}_1)p(\mathbf{W}_2)p(\mathbf{b})d\mathbf{W}_1d\mathbf{W}_2d\mathbf{b}.$$

## B.2 HIERARCHICAL VARIATIONAL INFERENCE IN THE APPROXIMATE MODEL

We next approximate the posterior over these variables with appropriate hierarchical approximating variational distributions. We define a hierarchical variational distribution as:

$$q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}) := q(\mathbf{W}_1)q(\mathbf{W}_2)q(\mathbf{b}) = \int q(\mathbf{W}_1|\mathbf{U}_1)q(\mathbf{W}_2|\mathbf{U}_2)q(\mathbf{U}_1)q(\mathbf{U}_2)q(\mathbf{b})d\mathbf{U}_1d\mathbf{U}_2,$$

where we define  $q(\mathbf{W}_1)$  to be a Gaussian mixture distribution with two components, which is factorised over  $Q$  and  $K$ :

$$\begin{aligned} q(\mathbf{W}_1|\mathbf{U}_1) &= \prod_{q=1}^Q \prod_{k=1}^K q(\mathbf{w}_{qk}|\mathbf{u}_{qk}), \\ q(\mathbf{w}_{qk}|\mathbf{u}_{qk}) &= p_1\mathcal{N}(\mathbf{u}_{qk}, \sigma^2) + (1 - p_1)\mathcal{N}(0, \sigma^2), \\ q(\mathbf{u}_{qk}) &= \mathcal{N}(\mathbf{v}_{qk}, \sigma^2) \end{aligned}$$

where  $p_1 \in [0, 1]$ , and  $\sigma > 0$ . Similarly, we can define a hierarchical variational distribution over  $\mathbf{W}_2$

$$\begin{aligned} q(\mathbf{W}_2|\mathbf{U}_2) &= \prod_{k=1}^K \prod_{d=1}^D q(\mathbf{w}_{kd}|\mathbf{u}_{kd}), \\ q(\mathbf{w}_{kd}|\mathbf{u}_{kd}) &= p_2\mathcal{N}(\mathbf{u}_{kd}, \sigma^2) + (1 - p_2)\mathcal{N}(0, \sigma^2), \\ q(\mathbf{u}_{kd}) &= \mathcal{N}(\mathbf{v}_{kd}, \sigma^2) \end{aligned}$$

For  $\mathbf{b}$ , we use a simple Gaussian distribution

$$q(\mathbf{b}) = \mathcal{N}(\mathbf{u}, \sigma^2 \mathbf{I}_K).$$

## B.3 EVALUATING THE LOG EVIDENCE LOWER BOUND FOR REGRESSION

Next we evaluate the log evidence lower bound for the task of regression. The log evidence lower bound is as below

$$\mathcal{L}_{\text{GP-VI}} := \int q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}) \log p(\mathbf{Y}|\mathbf{X}, \mathbf{W}_1, \mathbf{W}_2, \mathbf{b}) - D_{\text{KL}}(q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b})||p(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b})).$$

where the integration is with respect to  $\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}$ .

During regression, we can rewrite the integrand as a sum:

$$\begin{aligned} \log p(\mathbf{Y}|\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}) &= \sum_{d=1}^D \log \mathcal{N}(\mathbf{y}_d; \Phi \mathbf{w}_d, \tau^{-1} \mathbf{I}_N), \\ &= -\frac{ND}{2} \log(2\pi) + \frac{ND}{2} \log(\tau) - \sum_{d=1}^D \frac{\tau}{2} \|\mathbf{y}_d - \Phi \mathbf{w}_d\|_2^2. \end{aligned}$$

as the output dimensions of a multi-output Gaussian process are assumed to be independent. Denote  $\hat{\mathbf{Y}} = \Phi \mathbf{W}_2$ . We can then sum over the rows instead of the columns of  $\hat{\mathbf{Y}}$  and write

$$\sum_{d=1}^D \frac{\tau}{2} \|\mathbf{y}_d - \hat{\mathbf{y}}_d\|_2^2 = \sum_{n=1}^N \frac{\tau}{2} \|\mathbf{y}_n - \hat{\mathbf{y}}_n\|_2^2.$$

Here we have  $\hat{\mathbf{y}}_n = \phi(\mathbf{x}, \mathbf{W}_1, \mathbf{b}) \mathbf{W}_2 = \sqrt{\frac{1}{K}} \sigma (\mathbf{x}_n \mathbf{W}_1 + \mathbf{b}) \mathbf{W}_2$ , leading to

$$\log p(\mathbf{Y}|\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}) = \sum_{n=1}^N \log \mathcal{N}(\mathbf{y}_n; \phi(\mathbf{x}_n, \mathbf{W}_1, \mathbf{b}) \mathbf{W}_2, \tau^{-1} \mathbf{I}_D).$$



Therefore, we can update the log evidence lower bound as

$$\sum_{n=1}^N \int q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}) \log p(\mathbf{Y}|\mathbf{x}_n, \mathbf{W}_1, \mathbf{W}_2, \mathbf{b}) - D_{\text{KL}}(q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b})\|p(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b})).$$

We re-parametrise the integrands in the sum to not depend on  $\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}$  directly, but instead on the standard normal distribution and the Bernoulli distribution. Let  $q(\epsilon_1) = \mathcal{N}(0, \mathbf{I}_{Q \times K})$ ,  $q(\mathbf{z}_{1,q,k}) = \text{Bernoulli}(p_1)$ ,  $q(\epsilon_2) = \mathcal{N}(0, \mathbf{I}_{K \times D})$ ,  $q(\mathbf{z}_{2,k,d}) = \text{Bernoulli}(p_2)$ ,  $q(\epsilon) = \mathcal{N}(0, \mathbf{I}_K)$ ,  $q(\epsilon_3) = \mathcal{N}(0, \mathbf{I}_{Q \times K})$ , and  $q(\epsilon_4) = \mathcal{N}(0, \mathbf{I}_{K \times D})$ . Then, we can have

$$\begin{aligned} \mathbf{W}_1 &= \mathbf{Z}_1 \odot (\mathbf{U}_1 + \sigma \epsilon_1) + (1 - \mathbf{Z}_1) \odot \sigma \epsilon_1, \\ \mathbf{W}_2 &= \mathbf{Z}_2 \odot (\mathbf{U}_2 + \sigma \epsilon_2) + (1 - \mathbf{Z}_2) \odot \sigma \epsilon_2, \\ \mathbf{b} &= \mathbf{u} + \sigma \epsilon, \quad \mathbf{U}_1 = \sigma \epsilon_3, \quad \mathbf{U}_2 = \sigma \epsilon_4 \end{aligned}$$

where  $\odot$  means element-wise multiplication. Thus, we can update the above the sum over the integrals

$$\begin{aligned} &\sum_{n=1}^N \int q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b}) \log p(\mathbf{y}_d|\mathbf{x}_n, \mathbf{W}_1, \mathbf{W}_2, \mathbf{b}), \\ &= \sum_{n=1}^N \int q(\mathbf{Z}_1, \epsilon_1, \mathbf{Z}_2, \epsilon_2, \epsilon, \epsilon_3, \epsilon_4) \log p(\mathbf{y}_d|\mathbf{x}_n, \mathbf{W}_1(\mathbf{Z}_1, \epsilon_1, \epsilon_3), \mathbf{W}_2(\mathbf{Z}_2, \epsilon_2, \epsilon_4), \mathbf{b}(\epsilon)). \end{aligned}$$

For the first term in  $\mathcal{L}_{\text{GP-MC}}$ , we can estimate the integrals using Monte Carlo integration with a distinct single sample to obtain

$$\mathcal{L}_{\text{GP-MC}} := \sum_{n=1}^N \int \log p(\mathbf{y}_d|\mathbf{x}_n, \widehat{\mathbf{W}}_1^n, \widehat{\mathbf{W}}_2^n, \widehat{\mathbf{b}}^n) - D_{\text{KL}}(q(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b})\|p(\mathbf{W}_1, \mathbf{W}_2, \mathbf{b})).$$

Following Gal & Ghahramani (2016), by optimising the stochastic objective  $\mathcal{L}_{\text{GP-MC}}$ , we can converge to the same limit as  $\mathcal{L}_{\text{GP-VI}}$ , which justifies this stochastic approximation.

Moving to the second term in  $\mathcal{L}_{\text{GP-MC}}$ , we use  $\mathbf{w}$  and  $\mathbf{u}$  to denote certain component in the weights ( $\mathbf{W}_1$ ) and variational parameters ( $\mathbf{U}_1$ ), respectively. Then, we can have

$$\begin{aligned} &-D_{\text{KL}}\left(q(\mathbf{w})\|p(\mathbf{w})\right) = -\int q(\mathbf{w}) \log \frac{q(\mathbf{w})}{p(\mathbf{w})} d\mathbf{w} \\ &= \int \int q(\mathbf{w}, \mathbf{u}) d\mathbf{u} \log \frac{p(\mathbf{w})}{\int q(\mathbf{w}, \mathbf{u}) d\mathbf{u}} d\mathbf{w} \\ &= \int q(\mathbf{u}) \left\{ \int q(\mathbf{w}|\mathbf{u}) \log \frac{p(\mathbf{w})}{\int q(\mathbf{w}, \mathbf{u}) d\mathbf{u}} d\mathbf{w} \right\} d\mathbf{u} \\ &= \int q(\mathbf{u}) \left[ \int q(\mathbf{w}|\mathbf{u}) \log p(\mathbf{w}) d\mathbf{w} \right] d\mathbf{u} - \int q(\mathbf{u}) \left\{ \int q(\mathbf{w}|\mathbf{u}) \log \left[ \int q(\mathbf{w}, \mathbf{u}) d\mathbf{w} \right] d\mathbf{w} \right\} d\mathbf{u} \\ &= \int q(\mathbf{u}) \left[ \int q(\mathbf{w}|\mathbf{u}) \log p(\mathbf{w}) d\mathbf{w} \right] d\mathbf{u} - \int \left[ \int q(\mathbf{w}, \mathbf{u}) d\mathbf{u} \right] \log \left[ \int q(\mathbf{w}, \mathbf{u}) d\mathbf{u} \right] d\mathbf{w} \end{aligned} \quad (12)$$

For the first term in Eq. (12), we can first follow the Proposition 1 in (Gal & Ghahramani, 2016) to approximate  $\int q(\mathbf{w}|\mathbf{u}) \log p(\mathbf{w}) d\mathbf{w}$  as below:

$$\int q(\mathbf{w}|\mathbf{u}) \log p(\mathbf{w}) d\mathbf{w} \approx -\frac{1}{2}(p \cdot \mathbf{u}^2 + \sigma^2)$$

Then, we can approximate the first term in Eq. (12) as

$$\begin{aligned} &\int q(\mathbf{u}) \left[ \int q(\mathbf{w}|\mathbf{u}) \log p(\mathbf{w}) d\mathbf{w} \right] d\mathbf{u} \approx \int q(\mathbf{u}) \left[ -\frac{1}{2}(p \cdot \mathbf{u}^2 + \sigma^2) \right] d\mathbf{u} \\ &= -\frac{1}{2}\sigma^2 - \frac{p}{2} \int q(\mathbf{u}) \mathbf{u}^2 d\mathbf{u} = -\frac{p+1}{2}\sigma^2 - \frac{p}{2}v^2 \end{aligned}$$

For the second term in Eq. (12), we can estimate the integral  $\int q(\mathbf{w}, \mathbf{u})d\mathbf{u}$  using Monte Carlo integration with a distinct single sample to obtain

$$\int [\int q(\mathbf{w}, \mathbf{u})d\mathbf{u}] \log[\int q(\mathbf{w}, \mathbf{u})d\mathbf{u}]d\mathbf{w} \approx \int q(\mathbf{w}|\hat{\mathbf{u}}) \log[q(\mathbf{w}|\hat{\mathbf{u}})]d\mathbf{w}$$

Then we can follow the Proposition 1 in (Gal & Ghahramani, 2016) to approximate it as

$$\int q(\mathbf{w}|\hat{\mathbf{u}}) \log[q(\mathbf{w}|\hat{\mathbf{u}})]d\mathbf{w} \approx \frac{1}{2}(\log \sigma^2 + 1 + 2 \log \pi) + C$$

Therefore, we can get the approximation for Eq. (12) as Eq. (13):

$$-D_{\text{KL}}(q(\mathbf{w})\|p(\mathbf{w})) = -\frac{p+1}{2}\sigma^2 - \frac{p}{2}\mathbf{v}^2 + \frac{1}{2}(\log \sigma^2 + K(1 + 2 \log \pi)) + C \quad (13)$$

Then, we can have the following equation based on Eq. (13):

$$D_{\text{KL}}(q(\mathbf{W}_1)\|p(\mathbf{W}_1)) \approx \frac{QK(p+1)}{2}\sigma^2 - \frac{QK}{2}(\log(\sigma^2) + 1) + \frac{p_1}{2} \sum_{q=1}^Q \sum_{k=1}^K \mathbf{v}_{qk}^2 + C,$$

where  $C$  is a constant and  $D_{\text{KL}}(q(\mathbf{W}_2)\|p(\mathbf{W}_2))$  can be approximated in a similar way. For  $D_{\text{KL}}(q(\mathbf{b})\|p(\mathbf{b}))$ , it can be written as

$$D_{\text{KL}}(q(\mathbf{b})\|p(\mathbf{b})) = \frac{1}{2}(\mathbf{u}^\top \mathbf{u} + K(\sigma^2 - \log(\sigma^2) - 1)) + C.$$

#### B.4 LOG EVIDENCE LOWER BOUND OPTIMISATION FOR CIGL

Next we explain the relation between the above equations with equations for CigL. Ignoring the constant terms  $\tau, \sigma$  we obtain the maximisation objective

$$\mathcal{L}_{\text{GP-MC}} \propto -\frac{\tau}{2} \sum_{n=1}^N \|\mathbf{y}_n - \hat{\mathbf{y}}_n\|_2^2 - \frac{p_1}{2} \|\mathbf{V}_1\|_2^2 - \frac{p_2}{2} \|\mathbf{V}_2\|_2^2 - \frac{1}{2} \|\mathbf{u}\|_2^2. \quad (14)$$

We will show the equivalence between the iterative update of  $\mathbf{M}_1, \mathbf{Z}_1$  and  $\mathbf{U}_1$  in CigL and the hierarchical variational inference for deep GP. The update for  $\mathbf{M}_2, \mathbf{Z}_2$  and  $\mathbf{U}_2$  will be similar. In the hierarchical variational distribution, the distribution of sparse weights  $\mathbf{W}_1$  depends on three random variables  $\mathbf{M}_1, \mathbf{Z}_1$ , and  $\mathbf{U}_1$ . Given  $\mathbf{Z}_1$  and  $\mathbf{U}_1$ , we can know the variational distribution for  $\mathbf{M}_1$  as:

$$q(\mathbf{M}_1|\mathbf{U}_1) \propto \exp(\mathbf{M}_1 \odot |\mathbf{U}_1|) \quad (15)$$

where  $\mathbf{M}_1$  is under certain sparsity constraint. Thus, we can update  $\mathbf{M}_1$  by choosing  $\widehat{\mathbf{M}}_1$  that maximising Eq. (15), which is aligned with the  $\mathbf{M}_1$  update procedure in CigL.

Given  $\mathbf{M}_1$ , we can use  $\widehat{\mathbf{V}}_1$  to approximate  $\widehat{\mathbf{U}}_1$ . Then, we can let  $\sigma$  tend to zero (Gal & Ghahramani, 2016), and the random variable realisations  $\widehat{\mathbf{W}}_1^n, \widehat{\mathbf{W}}_2^n, \hat{\mathbf{b}}^n$  can be

$$\widehat{\mathbf{W}}_1^n \approx \widehat{\mathbf{Z}}_1 \odot \widehat{\mathbf{U}}_1, \quad \widehat{\mathbf{W}}_2^n \approx \widehat{\mathbf{Z}}_2 \odot \widehat{\mathbf{U}}_2, \quad \hat{\mathbf{b}}^n \approx \hat{\mathbf{u}}$$

Then, we can get

$$\hat{\mathbf{y}}_n \approx \sqrt{\frac{1}{K}} \sigma (\mathbf{x}_n(\widehat{\mathbf{Z}}_1 \odot \widehat{\mathbf{U}}_1) + \hat{\mathbf{u}})(\widehat{\mathbf{Z}}_2 \odot \widehat{\mathbf{U}}_2).$$

We scale the optimisation objective by a positive constant  $\frac{1}{\tau N}$  and get the objective:

$$\mathcal{L}_{\text{GP-MC}} \propto -\frac{1}{2N} \sum_{n=1}^N \|\mathbf{y}_n - \hat{\mathbf{y}}_n\|_2^2 - \frac{p_1}{2\tau N} \|\mathbf{U}_1\|_2^2 - \frac{p_2}{2\tau N} \|\mathbf{U}_2\|_2^2 - \frac{1}{2\tau N} \|\mathbf{u}\|_2^2.$$

So, we recover equation for CigL. With correct stochastic optimisation scheduling, both will converge to the same limit.

## B.5 MASK & WEIGHT AVERAGING FOR PREDICTION

For prediction, we design mask & weight averaging (WMA) to produce the final output model. Specifically, we collect samples of  $\{\widehat{\mathbf{Z}}_1, \widehat{\mathbf{U}}_1\}$  during the optimization process, where we use  $\widehat{\mathbf{V}}_1$  at different epochs to approximate  $\widehat{\mathbf{U}}_1$ . By using weight & mask averging and letting  $\sigma$  tend to zero (Gal & Ghahramani, 2016), the random variable  $\widehat{\mathbf{W}}_1, \widehat{\mathbf{W}}_2, \widehat{\mathbf{b}}$  can be approximated as

$$\widehat{\mathbf{W}}_1 \approx \frac{1}{S} \sum_{s=1}^S \widehat{\mathbf{Z}}_1^{(s)} \odot \widehat{\mathbf{U}}_1^{(s)}, \quad \widehat{\mathbf{W}}_2 \approx \frac{1}{S} \sum_{s=1}^S \widehat{\mathbf{Z}}_2^{(s)} \odot \widehat{\mathbf{U}}_2^{(s)}, \quad \widehat{\mathbf{b}} \approx \frac{1}{S} \sum_{s=1}^S \widehat{\mathbf{u}}^{(s)}.$$

WMA can be seen as approximating the mean of the posterior based on samples from variational distribution  $q(\mathbf{W})$  using moment matching, which is justified as below:

(i) CigL is connected to the Bayesian approach because of using both deterministic and random masks to explore the weight space. As shown in Equation 3, the design of  $\mathbf{Z}$  and  $\mathbf{M}$  results in a hierarchical variational distribution  $q(\mathbf{W})$ , where the hierarchy expands the approximation family and leads to a better posterior approximation capability. In Equation 3, updating the mask is equivalent to updating and to bring closer to the posterior 3.

(ii) In a similar idea to weight dropout (Gal & Ghahramani, 2016), WMA can be considered as a Bayesian approximation, which is used to approximate the mean of the posterior by moment matching.

- Weight dropout has been shown to be equivalent to the Bayesian approximation method (Gal & Ghahramani, 2016). After obtaining the final  $\mathbf{W}$ , dropout approximates  $\int f(\mathbf{W}\mathbf{Z})p(\mathbf{Z})d\mathbf{Z}$  using  $f(E(\mathbf{W}\mathbf{Z}))$ , where  $f$  is the neural network and the mean  $E(\mathbf{W}\mathbf{Z}) = \int \mathbf{W}\mathbf{Z}p(\mathbf{Z})d\mathbf{Z}$  (Srivastava et al., 2014). This approximation actually approximates the whole posterior by the first moment of the posterior (i.e., the mean).
- For our WMA, since we assume a hierarchy, we collect multiple samples of  $\mathbf{W}$  and  $\mathbf{Z}$ . Then, the WMA is used to approximate the first moment of the posterior which is used as an approximation of the posterior itself (Srivastava et al., 2014).
- In addition, if we really want the second moment, it is straightforward to obtain an estimation based on samples using moment matching again, similar to Maddox et al. (2019). We do not estimate the second moment since the sparse training typically wants a single sparse model in the end to reduce both computational and memory costs, and we find that using the posterior mean already significantly improves the calibration of the sparse training.

## C APPENDIX: ADDITIONAL EXPERIMENTAL RESULTS

### C.1 STRONGER CORRELATION BETWEEN HIDDEN VARIABLES

Empirically, we find that a stronger correlation between  $\mathbf{Z}$  and  $\mathbf{W}$  in sparse training. We use CigL to train sparse Wide-ResNet-22-2 on CIFAR-10 at multiple sparsities (0%, 50%, 80%, 90%). Then, we randomly draw five random masks  $\mathbf{Z}_i, i \in 1, \dots, 5$  from Bernoulli distribution. Using the final sparse weights  $\mathbf{W}$ , we obtain several new sparse models  $\mathbf{Z}_i \odot \mathbf{W}, i \in 1, \dots, 5$ , and record their test accuracies. We compare the test accuracy of  $\mathbf{W}$  with the average accuracy of  $\mathbf{Z}_i \odot \mathbf{W}, i \in 1, \dots, 5$  to see the correlation. The larger decrease in test accuracy after multiplying by  $\mathbf{Z}_i$  implies a stronger correlation between  $\mathbf{Z}$  and  $\mathbf{W}$ .

As shown in Figure 6 (a), both the accuracy of  $\mathbf{W}$  (red curve) and the average accuracy of  $\mathbf{Z}_i \odot \mathbf{W}, i \in 1, \dots, 5$  (blue curve) are decrease with increasing sparsity, and we see a more pronounced decrease in the blue curve. Figure 6 (b) further shows the decrease in test accuracy at each sparsity. We observe that the decrease is very small in the dense or sparse model at low sparsity. However, when it shifts to high sparsity such as 90%, we observe a larger decrease, which indicates a stronger correlation between  $\mathbf{Z}$  and  $\mathbf{W}$  in sparse training.

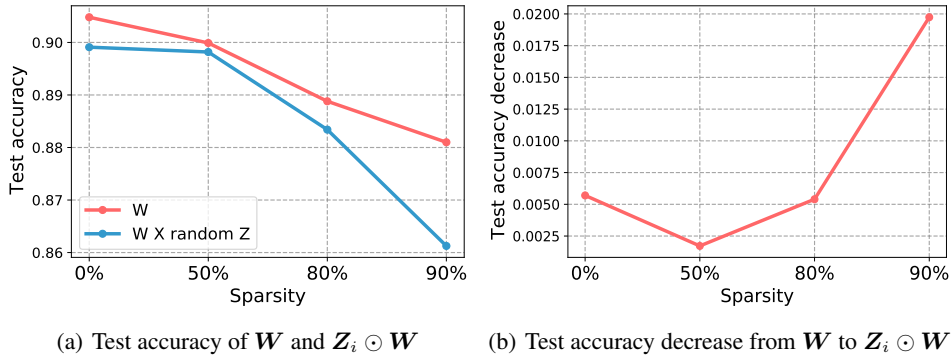


Figure 6: (a) Test accuracy of sparse model  $W$  and the newly produced sparse model  $Z_i \odot W$ . (b) decrease in test accuracy from sparse model  $W$  to newly produced sparse model  $Z_i \odot W$ . At low sparsity, the decrease of the dense or sparse models is small. At high sparsity, the decrease is larger.

## C.2 RELIABILITY IN SPARSE TRAINING

To get a more comprehensive understanding of the reliability issues in sparse training, we also evaluated the ECE values of the sparse models generated by SET (Mocanu et al., 2018). We find that the sparse model produced by SET is also more over-confident than the dense model. As shown in Table 4, the ECE values of dense ResNet-50 are smaller than those of sparse ResNet-50 on both CIFAR-10 and CIFAR-100. This proves the reliability issue of sparse training.

Table 4: ECE value of sparse ResNet-50 on CIFAR-10 and CIFAR-100 produced by SET at different sparsity including 0%, 50%, 80%, 90%, 95%, 99%.

SPARSITY	0%	50%	80%	90%	95%	99%
CIFAR-100	0.0381	0.0429	0.0416	0.0459	0.0460	0.0589
CIFAR-100	0.0841	0.0931	0.1058	0.1290	0.1282	0.0873

## C.3 MORE COMPARISON WITH SPARSE TRAINING BASELINE

We further compare our CigL with a recent Sparse training baseline Sup-tickets (Yin et al., 2022) to show the effectiveness of CigL in reducing ECE values. Table 5 shows the change in ECE after using Sup-tickets or our CigL. We can see that Sup-tickets brings only a limited reduction in ECE, while the reduction of our CigL is much larger than that of Sup-tickets.

Table 5: ECE value changes of Sup-tickets and CigL in ResNet-50 on CIFAR-10 and CIFAR-100 at different sparsity including 80%, 90%, 95%.

	CIFAR-10			CIFAR-100		
	80%	90%	95%	80%	90%	95%
SUP-TICKETS	-0.0012	-0.0005	-0.0007	-0.0005	-0.0010	-0.0010
CIGL	<b>-0.0067</b>	<b>-0.0080</b>	<b>-0.0119</b>	<b>-0.0141</b>	<b>-0.0113</b>	<b>-0.0104</b>

## C.4 MORE RESULTS ABOUT THE EFFECT OF WEIGHT & MASK AVERAGING

To demonstrate that weight & mask averaging (WMA) is not effective in reducing ECE alone, we add more results of using only WMA without random masking (CigL w/o RM). Table 6 shows the change in ECE after using CigL w/o RM or our CigL. We find that when only WMA is used,

the ECE value cannot be effectively reduced, either increasing or with limited reduction. On the contrary, the reduction of our CigL is much larger than that of CigL w/o RM.

Table 6: ECE value changes of CigL w/o RM and CigL in ResNet-50 (CIFAR-100) and Wide-ResNet-22-2 (CIFAR-10) at different sparsity including 80%, 90%, 95%, 99%.

	RESNET-50, CIFAR-100				WIDE-RESNET-22-2, CIFAR-10			
	80%	90%	95%	99%	80%	90%	95%	99%
CigL w/o RM	0.0038	-0.0098	-0.0144	-0.0046	0.0060	0.0129	0.0149	0.0017
CigL	<b>-0.0010</b>	<b>-0.0152</b>	<b>-0.0344</b>	<b>-0.0269</b>	<b>-0.0141</b>	<b>-0.0113</b>	<b>-0.0104</b>	<b>-0.0049</b>

## D MORE DISCUSSION

### D.1 WEIGHT SPACE EXPLORATION

ITOP (Liu et al., 2021) and DST-EE (Huang et al., 2022) study weight space exploration in sparse training and emphasized its importance. Compared to their studies, our work has two main differences that address their limitations.

On the one hand, our work has a different goal from ITOP and DST-EE with respect to encouraging exploration of the weight space. Specifically, our work aims to better explore the weight space to find more reliable models, while ITOP and DST-EE aims to build models with higher accuracy, ignoring the safety aspects.

On the other hand, the exploration of weight space has two aspects, namely “which weight is active” & “what value that weight has”. The limitation of ITOP is that, given the mask, the second aspect is not addressed and the optimization of the algorithm remains more challenging than dense training due to the pseudo-local optimization introduced by the sparsity constraints. To meet this challenge, ITOP increases the iterations between mask updates, leading to an increase in training time. For DST-EE, it mainly targets the first aspect. In contrast to their study, our work addresses this limitation, as shown in the following discussion:

The first aspect of weight space exploration is reflected by the ITOP rate, which is the percentage of all weights that have ever been selected as active weights by the mask. The second aspect of weight space exploration is reflected by the idea of “reliable exploration” in the ITOP paper. Ideally, a reliable exploration should allow a model to find the good direction and jump out of the bad local optimum. The sparsity constraints introduce some pseudo-local optima, which is difficult to jump out of. Our random mask can randomly cut off some directions and force the model to explore other directions, thus encouraging the model to better explore the weight space and avoid missing the correct direction.

### D.2 DOUBLE DESCENT IN RELIABILITY

One phenomenon we find worth discussing is the *double descent* in the reliability of sparse training. We discuss it in Section 6, where we divide the sparsity into four stages, i.e., poor model, shallow model, sparse deep model, and dense deep model.

The four stages are first supported by intuition. In the discussion, we draw analogies between model types such as “shallow models” and “poor models” in terms of model accuracy (expressiveness) and size. Consistent with the previous definition of double descent (Nakkiran et al., 2021; Somepalli et al., 2022), we consider sparsity as a measure of model size. Intuitively, as we gradually reduce the model size (increase the sparsity), we will go through four stages.

The four arguments are also supported by our sparse training experiments on ResNet-50 at CIFAR-100. As shown in Figure 1 (c), we can infer the model type by sparsity and accuracy:

- For 99.7% sparsity, the accuracy of the model is 41.7%, which is similar to a shallow model.

- For 99.9% sparsity, the accuracy of the model is 23.5%, which can be viewed as a poor model.

More detailed and quantitative support for these four arguments is beyond the main scope of this paper and could be a good direction for future research. One potential direction is the use of effective depth as a measure of stage identification.

### D.3 WEIGHT & MASK AVERAGING

Without the use of WMA, the analysis in Section 4.2 would be a non-hierarchical Bayesian method or a poor approximation to a hierarchical Bayesian approach.

(i) In the absence of WMA, the algorithm can be viewed as a non-hierarchical variational inference. As described in Section 4.3, using the final  $\mathbf{W}$  for prediction without WMA is equivalent to using weight dropout in RigL. Thus, the analysis in Section 4.2 will be updated in a similar way to Section 3 in Gal & Ghahramani (2016) which shows that weight dropout can be viewed as a non-hierarchical Bayesian approximation.

(ii) Without WMA, the algorithm can also be viewed as a poor approximation to hierarchical variational inference.

- If we continue to interpret the algorithm without WMA using the current analysis structure from Section 4.1, then how we generate the final posterior approximation will change.
- In this case, although the algorithm is still a Bayesian approximation, we only use the final  $\mathbf{W}$  to represent the posterior, which does not effectively capture all the information we explore from the weight space and the increased correlation between  $\mathbf{Z}$  and  $\mathbf{W}$ .
- Therefore, it turns out to be a bad hierarchical approximation, which limits its power.

### D.4 MULTI-MASK SPARSE DNNs

Existing multi-mask methods are not designed for improved weight space exploration. Bibikar et al. (2022) considers sparse training in federated learning and investigates the aggregation of multiple masks in edge devices. Xia et al. (2022) utilizes multiple masks with different granularities to allow greater flexibility in structured pruning and to improve accuracy. Despite the use of multiple masks, existing work (Xia et al., 2022; Bibikar et al., 2022) differ significantly from our work. They still use deterministic masks, which still suffer from the lack of exploration of the weight space, and consider only the accuracy of sparse models. In addition, their setups are federated learning and pruning, which are different from our work.

### D.5 SPARSE DNNs: PRUNING & SPARSE TRAINING

Although pruning (e.g., Lottery Tickets) and sparse training are related and both produce subnetworks with high accuracy, their goals and discovering journeys are quite different, which leads to significant differences in several important properties, including uncertainty, geometry of the loss surface, generalization ability, and so on.

(i) For the goal, Lottery Tickets mainly aim to reduce the inference cost, while the sparse training also aims to save resources during the training phase.

(ii) Lottery Tickets and sparse training are different in several important properties. As shown in Figure 11 of Chen et al. (2022), the blue and purple bars represent Lottery Tickets and sparse training with a sparsity level of 79%, respectively.

- For uncertainty, sparse training does not improve the confidence calibration compared to dense training, while Lottery Tickets allows for improved confidence calibration.
- For the geometry of the loss surface, sparse training leads to larger trace values and cannot locate flat local minima. In contrast, Lottery Tickets can still locate flat local minima.
- For generalization ability, sparse training provides higher accuracy and improved robustness compared to dense training, while Lottery Tickets provide relatively less improvement.

(iii) The main reason for the different properties is their different discovering journeys.

For Lottery Tickets:

- It retrains the weights from the initial training phase after each pruning, which significantly increases the training time but allows more time for the model to explore the weight space.
- It starts from a dense model and has low sparsity in the early stages, which reduces the difficulty of weight space exploration caused by the sparsity constraints.

For sparse training:

- It maintains a high level of sparsity throughout the training process, which does not extend the training time to enable more exploration of the weight space.
- In addition, maintaining high sparsity can cut off a large portion of the optimization route and produce more spurious local minima, thus making training very difficult.
- Chen et al. (2022) shows the differences in the properties of Lottery Tickets and sparse training at the 79% sparsity level. The difficulty of training typically increases with increasing sparsity, implying that the difference is likely to be greater at higher sparsity levels.