

1 Appendix

2 The following manuscript provides the supplementary materials of the main paper: Regularizing
3 Neural Networks with Meta-Learning Generative Models.

4 A Algorithm of Meta Generative Regularization

Algorithm 1 Meta Generative Regularization

Require: Training dataset \mathcal{D} , validation dataset \mathcal{D}_{val} , main model f_θ , generator G , finder F_ϕ , training batchsize B , pseudo batchsize B_p , validation batchsize B_{val} , step size η and ξ , hyperparameter λ and λ_{KL} .

Ensure: Trained main model f_θ

```
1: while not converged do  
2:    $\{(x^i, y^i)\}_{i=1}^B \sim \mathcal{D}$   
3:    $\{z^i\}_{i=1}^{B_p} \sim \mathcal{N}(0, I)$   
4:   // Updating  $\phi$  for MPS  
5:    $\{(x_{\text{val}}^i, y_{\text{val}}^i)\}_{i=1}^{B_{\text{val}}} \sim \mathcal{D}$   
6:    $\{x_p^i\}_{i=1}^{B_p} = \{G_\Phi(F_\phi(z^i), y_p^i)\}_{i=1}^{B_p}$   
7:    $\theta' \leftarrow \theta - \eta \nabla_\theta (\frac{1}{B} \ell(f_\theta(x^i), y^i) + \frac{\lambda}{B_p} \ell_{\text{PCR}}(x_p^i; \psi))$   
8:    $\phi \leftarrow \phi - \xi \nabla_\phi (\frac{1}{B_{\text{val}}} \ell(f_{\theta'}(x_{\text{val}}), y_{\text{val}}) + \lambda_{\text{KL}}(D_{\text{KL}}(p_\phi(z) \| p(z))))$   
9:   // Updating  $\theta$  with PCR  
10:   $\{x_p^i\}_{i=1}^{B_p} = \{G_\Phi(F_\phi(z^i), y_p^i)\}_{i=1}^{B_p}$   
11:   $\theta \leftarrow \theta - \eta \nabla_\theta (\frac{1}{B} \ell(f_\theta(x^i), y^i) + \frac{\lambda}{B_p} \ell_{\text{PCR}}(x_p^i; \psi))$   
12: end while
```

5 B Additional Experiments

6 B.1 Evaluation of Gradient Approximation

7 Here, we evaluate the gradient approximation by Eq. (9). As shown in Table 1, the 1st-order
8 approximation by Eq. (9) well approximated the second-order gradients in speeding up over 10%
with 0.08 of the accuracy drop.

Table 1: Performance comparison between MPS with 2nd-order gradients and 1st-order approximated gradients (ResNet-18, Cars).

Method	Top-1 Acc. (%)	Wall Clock Time (hours)
2nd-Order	87.30 \pm .39	6.55
1st-Order Approx.	87.22 \pm .15	5.79

9

10 B.2 Ablation study of F_ϕ

11 In Section 3.2, we introduce F_ϕ for meta-optimized parameters and the residual architectures with
12 MLP defined by Eq. (10). We performed an ablation study of MPS with respect to the meta-optimized
13 parameters and the architectures of F_ϕ . We compared MPS with a variant of MPS optimizing G_Φ
14 instead of F_ϕ . We also attempted other architectures for F_ϕ including **Linear**: $W_\phi(z) + b$, **MLP**:
15 $\text{MLP}_\phi(z)$, and **Residual+Shallow**: $z + \tanh(W_\phi(z) + b)$. The results of these variations are shown
16 in Table 2. We observed that MPS with G_Φ caused failures of training f_θ and degrades the accuracy.
17 On the other hand, all variants of MPS with F_ϕ succeeded in boosting the models without MPS.
18 Thus, restricting the number of optimized parameters is important, and determining an optimal z with
19 the finder F_ϕ is effective on the optimization problems of MPS. For the variants of MPS with F_ϕ ,
20 we observed that the residual architectures and regularization by $D_{\text{KL}}(p_\phi(z) \| p(z))$ contributed to
21 the successes. Interestingly, MPS with Linear F_ϕ outperformed MPS with MLP F_ϕ , i.e., significantly
22 transforming the input $z \sim p(z)$ by complex functions results in low accuracy. These results suggest

23 that better latent vectors in \mathcal{Z} to train f_θ can exist near the uniformly sampled input z . Thus, limiting
 24 the search range by \tanh in the residual architectures can help in finding better latent vectors.

Table 2: Ablation study of MPS (ResNet-18, Cars).1

Method	Top-1 Acc. (%)
Without MPS (PCR)	86.32 \pm .07
MPS	87.22\pm.15
MPS with G_ϕ	84.47 \pm .05
MPS with Linear F_ϕ	86.51 \pm .09
MPS with MLP F_ϕ	86.35 \pm .13
MPS with Residual+Shallow F_ϕ	86.88 \pm .16
MPS w/o $D_{\text{KL}}(p_\phi(z) p(z))$	86.92 \pm .22

25 B.3 MGR with Diffusion Models

26 We tested our method on EDM [1], a recent diffusion model. Due to the computation cost, we used a
 27 10% reduced CIFAR-10 as the dataset. We optimized F_ϕ to search the first step noise of the diffusion
 28 process. Table 3 shows that our method with EDM improves Base Model. However, the overhead
 29 of incorporating diffusion models was significant; it takes more than ten times longer training than
 30 GANs. In future work, we will investigate lighter-weight methods using the diffusion model.

Table 3: Performance studies on Diffusion Model (ResNet-18 on Cars)

Method	Top-1 Acc. (%)
Base Model	86.49 \pm .48
GDA (EDM)	85.80 \pm .30
MGR	88.49\pm.12

31 B.4 Updating F_ϕ without Meta-optimization

32 MPS consists of meta-learning on validation losses requiring bi-level optimization, which is a
 33 relatively heavy computation. One can consider if F_ϕ could be trained without meta-optimization.
 34 Here, we try alternative methods other than meta-learning to update F_ϕ . Instead of meta-learning, we
 35 used a strategy of choosing hard examples via optimizing F_ϕ . That is, we optimize F_ϕ by maximizing
 36 the training cross-entropy (CE) loss and the PCR loss on synthetic samples. Note that, in both cases,
 37 we used the PCR loss for synthetic samples when training classifiers. Table 4 shows the results.
 38 Optimizing F_ϕ with CE and PCR slightly improved the baselines but significantly underperformed
 39 our method (MGR). This result can justify using meta-optimizing F_ϕ to generate useful samples for
 40 classifiers. Nevertheless, this idea could inspire a sampling method that does not require bi-level
 41 optimization in future work.

Table 4: Performance comparison of updating strategies for F_ϕ (ResNet-18 on Cars)

Method	Top-1 Acc. (%)
Base Model	85.50 \pm .10
PCR	86.36 \pm .08
Optimizing F_ϕ w/ CE	86.52 \pm .21
Optimizing F_ϕ w/ PCR	86.44 \pm .68
MGR	87.22\pm.15

42 **References**

- 43 [1] Tero Karras, Miika Aittala, Timo Aila, and Samuli Laine. Elucidating the design space of
44 diffusion-based generative models. In *Advances in Neural Information Processing Systems*, 2022.