

# Supplementary Material

## Model Fusion for Personalized Learning

### A. Additional Experiment Results

This section provides extra experiment results and visualization plots that could not be included in the main text due to limited space. Additional information regarding our experiment setup is also included for better clarity. Our experimental code is also released at the GitHub repository below:

[https://github.com/zevergreenz/model\\_fusion\\_for\\_personalized\\_learning](https://github.com/zevergreenz/model_fusion_for_personalized_learning)

#### A.1. Sine Prediction Experiments

As mentioned briefly in the main text, the personalized performance of MAML (Finn et al., 2017) in problem domains with highly polarized task distributions deteriorates and varies substantially across different sub-domains of tasks. This can be observed in Table 5 below which tabulates the reported performance of all baseline methods across different sub-domains and few-shot settings. In Table 4 and Table 5, we have also extended the benchmark to include reported performance for the latest personalized federated learning work (PerFedAvg) of Fallah et al. (2020). The extended benchmark shown that on average over the entire (polarized) sine domain, PerFedAvg achieves significantly better performance than MAML (Finn et al., 2017) but still incurs higher normalized RMSE than the variant of our model (with embedding alignment). The performance difference is again most pronouncing in the 0-shot setting when there is no data for further fine-tuning. On a closer look, it appears that PerFedAvg slightly performs better than ours in sub-domain 1 under the richer-data settings of 20-shot and 40-shot but in the remaining scenarios across the remaining 4 sub-domains, ours is substantially better than PerFedAvg. This agrees with and supports our above observation that on average, PerFedAvg performs better than MAML but still noticeably under-perform when compared to our model. For more visualized performance plot showcasing the fitting of various  $m$ -shot personalized neural net on some (unseen) sine functions, see Appendix F.

	0-SHOT	10-SHOT	20-SHOT	30-SHOT	40-SHOT
OURS [1-100-1]	<b>0.135 ± 0.09</b>	<b>0.114 ± 0.03</b>	<b>0.029 ± 0.01</b>	<b>0.021 ± 0.01</b>	<b>0.017 ± 0.01</b>
OURS [1-100-1] (NO ALIGN)	0.357 ± 0.01	0.125 ± 0.01	0.043 ± 0.01	0.036 ± 0.01	0.028 ± 0.01
COLD-START [1-100-1]	0.445 ± 0.26	0.149 ± 0.03	0.085 ± 0.05	0.076 ± 0.05	0.052 ± 0.04
PERFEDAVG [1-100-1]	0.356 ± 0.02	0.146 ± 0.03	0.046 ± 0.03	0.045 ± 0.04	0.032 ± 0.03
MAML [1-100-1]	0.446 ± 0.12	0.281 ± 0.18	0.161 ± 0.07	0.149 ± 0.06	0.142 ± 0.06
MAML [1-40-40-1]	<b>0.258 ± 0.11</b>	<b>0.022 ± 0.02</b>	<b>0.014 ± 0.01</b>	<b>0.014 ± 0.01</b>	<b>0.012 ± 0.01</b>

Table 4. Reported performance of our method (with and without embedding alignment), cold-start, PerFedAvg (Fallah et al., 2020) and MAML (Finn et al., 2017) on the sine synthetic dataset. The results is averaged over all domains and independent runs.

To further demonstrate the necessity of incorporating an embedding alignment component while learning to embed the existing pre-trained neural nets, we provide below in Fig. 3 another set of visual excerpts demonstrating the defection of the no-alignment variant of our method in zero-shot scenario. These plots in particular demonstrate a few scenarios where the no-alignment version fails to capture the generic trend of the target sine function. In contrast, all plots show that our version with embedding alignment is able to capture the trend and fit the target sine function much better, which highlights the importance of having the embedding alignment in Section 3.3.

#### A.2. Meta-Model Parameterization Specification

Due to limited space, we only describe the neural parameterization of many components of our meta model representation in high-level ideas. Here, we provide a more comprehensive specification of those neural parameterizations to improve the clarity of our work.

**Neural Parameterization  $\alpha$  for the Base Module.** This module comprises two network segments which are responsible for computing the mean vector and diagonal covariance matrix of the outer multivariate Gaussian that distributes  $\mathbf{w}_\emptyset$ . In our experiment,  $\mathbf{w}_\emptyset$  is a 100-dim vector. This is a deep generative prior which probabilistically generates  $\mathbf{w}_\emptyset$  from a noise vector. In particular, a noise vector of 100-dim is passed through a dense layer comprising 100 hidden neurons with no

**Model Fusion for Personalized Learning**

	SUB-DOMAIN 1: $a \sim U[4, 5], b = 2.94$			SUB-DOMAIN 2: $a \sim U[3, 4], b = 2.44$		
	10-SHOT	20-SHOT	40-SHOT	10-SHOT	20-SHOT	40-SHOT
OURS [1-100-1]	<b>0.16 ± 0.01</b>	0.06 ± 0.01	0.04 ± 0.01	<b>0.12 ± 0.01</b>	<b>0.02 ± 0.01</b>	<b>0.01 ± 0.01</b>
COLD-START [1-100-1]	0.31 ± 0.01	0.15 ± 0.01	0.08 ± 0.01	0.26 ± 0.01	0.19 ± 0.01	0.12 ± 0.01
PERFEDAVG [1-100-1]	<b>0.16 ± 0.05</b>	<b>0.02 ± 0.01</b>	<b>0.02 ± 0.01</b>	0.14 ± 0.01	0.03 ± 0.01	<b>0.01 ± 0.01</b>
MAML [1-100-1]	0.28 ± 0.25	0.12 ± 0.02	0.11 ± 0.03	0.25 ± 0.18	0.16 ± 0.05	0.13 ± 0.02
	SUB-DOMAIN 3: $a \sim U[2, 3], b = 1.50$			SUB-DOMAIN 4: $a \sim U[1, 2], b = 0.91$		
	10-SHOT	20-SHOT	40-SHOT	10-SHOT	20-SHOT	40-SHOT
OURS [1-100-1]	<b>0.06 ± 0.01</b>	<b>0.02 ± 0.01</b>	<b>0.01 ± 0.01</b>	<b>0.11 ± 0.01</b>	<b>0.02 ± 0.01</b>	<b>0.01 ± 0.01</b>
COLD-START [1-100-1]	0.27 ± 0.01	0.21 ± 0.01	0.11 ± 0.01	0.23 ± 0.01	0.13 ± 0.01	0.11 ± 0.01
PERFEDAVG [1-100-1]	0.12 ± 0.01	0.03 ± 0.01	<b>0.01 ± 0.01</b>	0.14 ± 0.01	0.08 ± 0.04	0.05 ± 0.03
MAML [1-100-1]	0.31 ± 0.14	0.16 ± 0.08	0.15 ± 0.07	0.24 ± 0.10	0.16 ± 0.04	0.16 ± 0.04
	SUB-DOMAIN 5: $a \sim U[0, 1], b = 0.91$					
	10-SHOT	20-SHOT	40-SHOT			
OURS [1-100-1]	<b>0.13 ± 0.22</b>	0.03 ± 0.01	<b>0.01 ± 0.01</b>			
COLD-START [1-100-1]	0.16 ± 0.03	<b>0.02 ± 0.01</b>	0.02 ± 0.01			
PERFEDAVG [1-100-1]	0.18 ± 0.01	0.09 ± 0.01	0.08 ± 0.01			
MAML [1-100-1]	0.23 ± 0.11	0.09 ± 0.05	0.07 ± 0.01			

Table 5. Reported (average) normalized RMSE of personalized neural nets with architecture [1-100-1] generated by our method (with alignment), cold-start, PerFedAvg (Fallah et al., 2020) and MAML (Finn et al., 2017) on each sub-domain. The (unseen) sine functions in each sub-domain share the same phase  $b$  but have different magnitudes  $a$ , which was sampled uniformly from polarized value ranges.

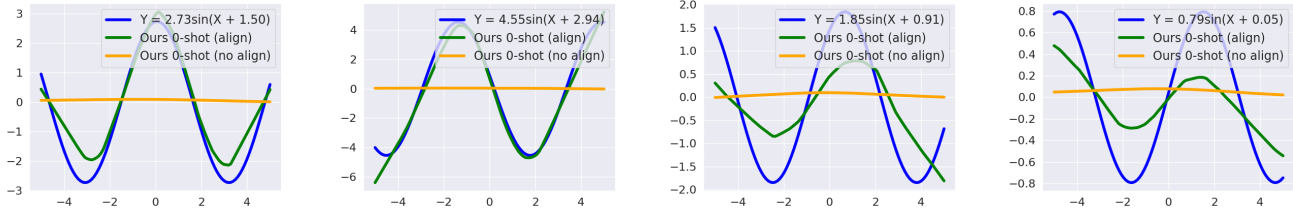


Figure 3. Visual excerpts demonstrating a few particular 0-shot scenarios where the no-alignment variant of our method generates poor/uninformed personalized neural nets. In contrast, the alignment version of our method is able to capture well the trends of the target sine functions in the same 0-shot scenarios. This highlights the importance of having an embedding alignment component when there is no learning example of the target task (i.e., 0-shot personalization) for the meta model to recognize the embedding mis-alignment via data.

activation, which produces the output of the first segment. For the second segment, the noise vector is passed through another dense layer comprising of 100 hidden neurons with tanh activation. The output is further normalized with a batch-norm layer, which finally produces the output of the second segment.

**GP-modulated Neural Parameterization  $\phi$  of the Task-Specific Module.** This module comprises a collection of  $e = 10$  independent sparse localized Gaussian processes (GPs) (Snelson, 2007) which represents the  $e$  (independent) priors over  $e$  random functions mapping from the task’s meta-data vector to a scalar value. The collective output of this  $e$  GPs is therefore an  $e$ -dim vector of independently distributed components, which is then warped into a 100-dim space of the task-specific component  $\mathbf{w}$ . The warper is designed to be a 2-layer feed-forward neural net. Its first layer comprises 256 neurons with tanh activation while its second layer comprises 100 neurons with no activation. Intuitively, this means the role of the GPs is to produce independent prediction signals which are weaved together via the 2-layer feed-forward net to generate a 100-dim task specific vector with highly correlated components. For each GP prior, the parameters comprises its kernel parameters as well as its inducing input vector. In our experiment, we use 50 inducing input vectors per GP where each inducing input is a  $\dim(\tau)$ -dim vector whose coordinates need to be learned.

**Neural Parameterization  $\nu$  for the Crossing Module  $p(\gamma|\mathbf{w}_\emptyset, \mathbf{w})$ .** This module comprises two network segments which are responsible for producing the mean vector and diagonal covariance matrix of the outer multivariate Gaussian that distributes  $\gamma$ . The two segments share a common sub-segment that maps the concatenated vector pf  $\mathbf{w}_\emptyset$  and  $\mathbf{w}$  to a latent space. In our experiment setting, both  $\mathbf{w}_\emptyset$  and  $\mathbf{w}$  are 100-dim vector. Their concatenation vector is first fed into a dense

layer of 64 hidden neurons with ReLU activation. The output of this layer is passed through another layer of  $d = 32$  hidden neurons with ReLU activation, which produces the 32-dim latent embedding for the concatenation of  $\mathbf{w}_\emptyset$  and  $\mathbf{w}$ . Next, to produce the mean vector of  $\gamma$ , the latent embedding is passed through a dense layer comprising  $|\gamma|$  hidden neurons with no activation, which constitutes the first network segment. As for the other, the same latent embedding (i.e., the output of the common sub-segment) is passed through a separate dense layer comprising  $\dim(\gamma) = 100$  hidden neurons with tanh activation, which generates the output of the second network segment.

Note the the above parameterization of  $\nu$  describes the processing of generative process of  $\gamma$  from  $\mathbf{w}$  and  $\mathbf{w}_\emptyset$  for a single task instance  $\tau$ . However, as we pointed out in the main text, pre-trained models  $\{\gamma_i\}_{i=1}^n$  were generated in isolated and hence, their induced model vector  $\{\gamma_i\}_{i=1}^n$  do not necessarily align component by component (e.g., component 1 of  $\gamma_i$  might correspond to component 3 of  $\gamma_j$ ) and we need to learn a matching permutation to align their components. To this end, we aim to learn for each model  $\gamma_i$  a separate permutation that maps it to a universal canonical order of components. These permutations  $\{\mathbf{A}_i\}_{i=1}^n$  are applied on the above 32-dim latent space of the embedding vector of the concatenation of  $\mathbf{w}$  and  $\mathbf{w}_\emptyset$  (i.e., the output of the common sub-segment described above). In our meta-model, the permutation matrices  $\{\mathbf{A}_i\}_{i=1}^n$  are treated as the discrete part of  $\nu$ , and are optimized alternatingly with the rest as detailed previously in Section 3.2.

### A.3. Computational Complexity

In addition, we provide below a concise analysis of the computational complexity of our meta model training. Plots showing its training losses against the no. of training iterations are also provided to demonstrate its convergence empirically.

First, to understand the computation complexity of our meta-model, we first note that it comprises the 3 deep generative modules as detailed above. Two of which (the base and crossing modules) are parametric networks in nature so the cost of their feed-forward and back-propagation computation is a linear function of their parameters, input dimension and batch sizes. As both are deep generative models, their parameters are optimized numerically via re-parameterized sampling (Kingma & Welling, 2013) which incurs another linear factor of the sampling size  $h$  on their feed-forward and back-propagation complexities (since we need to repeat the same computation routine for each sample).

Second, the other component (the task-specific module) on the other hand is non-parametric due to its modulation with a sparse GP prior layer whose feed-forward (prediction) scale linearly in the total number  $n$  of observed tasks (those with pre-trained models provided) and the dimension  $|\tau|$  of their corresponding meta-data vectors. Each sparse GP also scales quadratically in the number  $m$  of inducing inputs. For back-propagation however, it scales quadratically in both the number  $m$  of inducing inputs and the meta-data dimension  $\tau$  while remaining linear in  $n$ .

Furthermore, this module is also generative in nature which means both its forward and backward computation also incur a routine of forward sampling from the posterior GP. This includes first sampling the  $h$  samples of the inducing set comprising  $m$  vectors, which scales cubically in  $m$  since (per property of GPs) these vectors are marginally distributed by a  $m$ -dimensional Gaussian with full-rank covariance matrix. Thus, summarily, the total cost of feed-forward and back-propagation via a GP component incurs cubic cost in  $m$ , quadratic cost in  $|\tau|$  and linear in  $n$  and  $h$ . Also, as there are  $e$  separate GPs and there is another 2-layer feed-forward net that warps them into the space of  $\gamma$  (pre-trained model). The total processing cost of this component has to be further scaled by linear factors of  $e$  and  $\gamma$ .

Thus, putting the above together, we arrive at a complexity cost of  $\mathbf{O}(h \cdot |\alpha| \cdot n \cdot \dim(\mathbf{w}_\emptyset) + h \cdot |\nu| \cdot n \cdot (\dim(\mathbf{w}_\emptyset) + \dim(\mathbf{w})) + h \cdot e \cdot \dim(\gamma) \cdot n \cdot \dim(\tau)^2 \cdot m^3)$  for each learning epoch. In our experiment, we choose  $h = 100$  samples and use 100 learning epoches to optimize our meta-model. Additionally, per learning epoch, we also need to alternate between optimizing the continuous parameters and discrete permutation matrices for embedding alignment across different pre-trained models. The cost for the former has been given above while the cost for the latter boils down to the cost of solving a linear sum assignment task concerning  $d \times d$  bipartite matrix (see Section 3.2). Since the weight matrix for candidate assignment can be computed and cached while optimizing the continuous parameters, the cost of solving this linear sum assignment is cubic in  $d$ , which is the latent embedding dimension of the concatenation of  $\mathbf{w}_\emptyset$  and  $\mathbf{w}$  in the specification of the crossing module above. It also involves a linear factor of  $n$  since we need to solve one linear sum assignment per pre-trained model so in total, the cost is  $\mathbf{O}(n \cdot d^3)$  which is likely subsumed by  $\mathbf{O}(|\nu| \cdot n)$  as  $d$  only accounts for a small part of the entire crossing model’s parameterization  $|\nu|$ . As such, we can ignore it.

To further provide a empirical grasp of the above complexity analysis, the running time plots (along with corresponding plots of converging training losses) with respect to settings with  $e = 5, 10, 15$  and 20 sparse GP priors are given in Fig. 4 above. All our experiments were run on an Intel(R) Xeon(R) with 32 E5-2686 v4 CPUs (2.30GHz) and a V100 GPU.

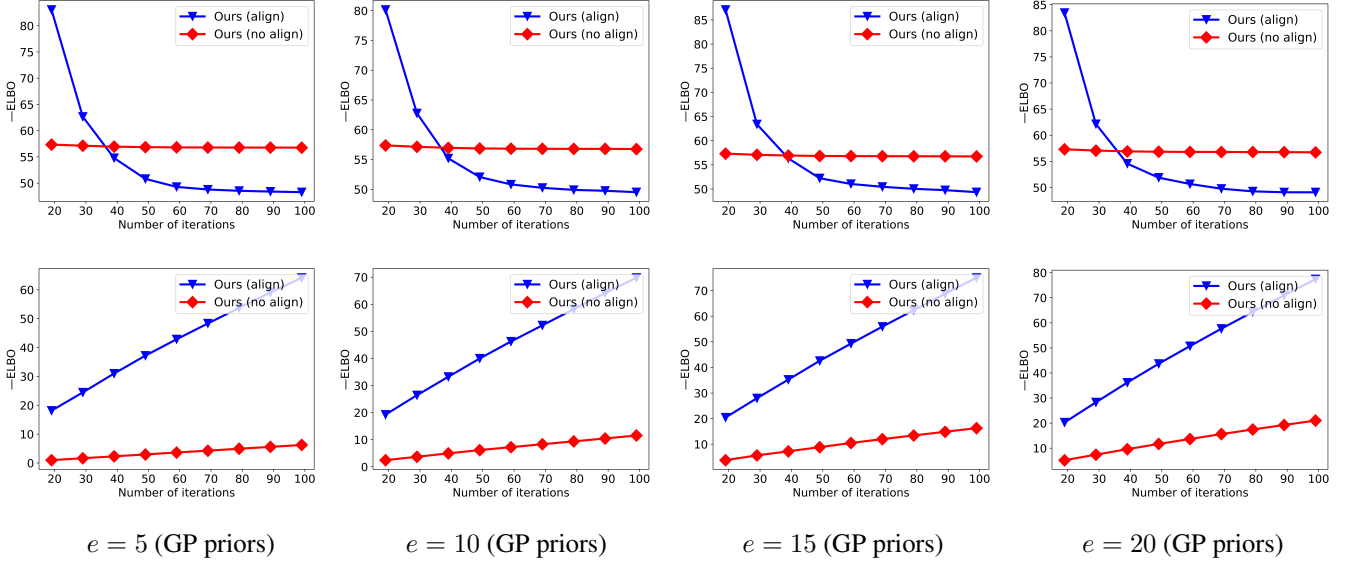


Figure 4. Plots show (top) the negative meta-model evidence lower-bound (ELBO) being reduced as we increase the no. of training iterations. This is the negation of the objective function in Eq. (16) that we sought to maximize (hence, its negative version would decrease as expected in the above plots); and the accumulated processing time (in minute) incurred by the alignment and no-alignment versions of our method. It can be seen that the negative ELBO decreases significantly when we allow for alignment during training as compared to without alignment. As a trade-off, the version with alignment incurs more processing time (up to eight-fold on average) as compared the version without alignment. The plots are provided for different parameter configurations with  $e = 5, 10, 15$  and  $20$  GP priors.

## B. Derivation of the Meta-Model Likelihood in Eq. (8)

The full distribution of our meta-model is given by

$$p(\mathbf{w}_\emptyset, \mathbf{w}, \gamma, \mathbf{s} | \tau) = p(\mathbf{w}_\emptyset) p(\mathbf{s}) \prod_{i=1}^n \left[ p(\gamma_i | \mathbf{w}_\emptyset, \mathbf{w}_i) p(\mathbf{w}_i | \mathbf{s}) \right]. \quad (20)$$

Marginalizing out  $\mathbf{w}_\emptyset$  and  $\mathbf{w}$ , we obtain the meta-model evidence in Eq. (8), Eq. (9) and Eq. (10)

$$\begin{aligned}
 p(\gamma | \tau) &= \int_{\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}} p(\mathbf{w}_\emptyset, \mathbf{w}, \gamma, \mathbf{s} | \tau) d\mathbf{s} d\mathbf{w} d\mathbf{w}_\emptyset \\
 &= \int_{\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}} p(\mathbf{w}_\emptyset) p(\mathbf{s}) \prod_{i=1}^n \left[ p(\gamma_i | \mathbf{w}_\emptyset, \mathbf{w}_i) p(\mathbf{w}_i | \mathbf{s}) \right] d\mathbf{s} d\mathbf{w} d\mathbf{w}_\emptyset \\
 &= \int_{\mathbf{w}_\emptyset} p(\mathbf{w}_\emptyset) \left[ \int_{\mathbf{w}, \mathbf{s}} p(\mathbf{s}) \prod_{i=1}^n \left[ p(\gamma_i | \mathbf{w}_\emptyset, \mathbf{w}_i) p(\mathbf{w}_i | \mathbf{s}) \right] d\mathbf{s} d\mathbf{w} \right] d\mathbf{w}_\emptyset = \mathbb{E}_{\mathbf{w}_\emptyset \sim p(\mathbf{w}_\emptyset)} \left[ h(\mathbf{w}_\emptyset) \right], \quad (21)
 \end{aligned}$$

where we define the auxiliary function  $h(\mathbf{w}_\emptyset)$  as

$$h(\mathbf{w}_\emptyset) = \int_{\mathbf{w}, \mathbf{s}} p(\mathbf{s}) \prod_{i=1}^n \left[ p(\gamma_i | \mathbf{w}_\emptyset, \mathbf{w}_i) p(\mathbf{w}_i | \mathbf{s}) \right] d\mathbf{s} d\mathbf{w}. \quad (22)$$

Now, substituting  $\mathbf{w} = \mathbf{g}(\boldsymbol{\tau})$  or  $\mathbf{g}$  for short

$$\begin{aligned}
 h(\mathbf{w}_\emptyset) &= \int_{\mathbf{g}(\boldsymbol{\tau}), \mathbf{s}} p(\mathbf{s}) \prod_{i=1}^n p(\mathbf{g}(\boldsymbol{\tau}_i) | \mathbf{s}) \prod_{i=1}^n \left[ p(\boldsymbol{\gamma}_i | \mathbf{w}_\emptyset, \mathbf{g}(\boldsymbol{\tau}_i)) \right] d\mathbf{s} d\mathbf{w} \\
 &= \int_{\mathbf{g}(\boldsymbol{\tau})} \underbrace{\left[ \int_{\mathbf{s}} p(\mathbf{s}) \prod_{i=1}^n \left[ p(\mathbf{g}(\boldsymbol{\tau}_i) | \mathbf{s}) \right] d\mathbf{s} \right]}_{p(\mathbf{g}(\boldsymbol{\tau}))} \prod_{i=1}^n \left[ p(\boldsymbol{\gamma}_i | \mathbf{w}_\emptyset, \mathbf{g}(\boldsymbol{\tau}_i)) \right] d\mathbf{w} \\
 &= \mathbb{E}_{\mathbf{g} \sim p(\mathbf{g})} \left[ \prod_{i=1}^n p(\boldsymbol{\gamma}_i | \mathbf{w}_\emptyset, \mathbf{g}(\boldsymbol{\tau}_i)) \right] \text{ where } p(\mathbf{g}) = \mathbb{E}_{\mathbf{s} \sim p(\mathbf{s})} \left[ \prod_{i=1}^n p(\mathbf{g}(\boldsymbol{\tau}_i) | \mathbf{s}) \right]. \tag{23}
 \end{aligned}$$

### C. Derivation of the Variational Inequality in Eq. (11)

From the chain rule (also called the product rule) of probability, suppose that  $p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\gamma}, \boldsymbol{\tau}) \neq 0$ , we have

$$p(\boldsymbol{\gamma} | \boldsymbol{\tau}) = \frac{p(\boldsymbol{\gamma}, \mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\tau})}{p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\gamma}, \boldsymbol{\tau})}. \tag{24}$$

Assuming that  $p(\boldsymbol{\gamma} | \boldsymbol{\tau}) \neq 0$  (consequently,  $p(\boldsymbol{\gamma}, \mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\tau}) \neq 0$ ), taking log on both sides of the above equality

$$\log p(\boldsymbol{\gamma} | \boldsymbol{\tau}) = \log \left( \frac{p(\boldsymbol{\gamma}, \mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\tau})}{p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\gamma}, \boldsymbol{\tau})} \right). \tag{25}$$

For any distribution  $q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)$  parameterized by  $\omega$  such that  $q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega) \ll p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\gamma}, \boldsymbol{\tau})$ , i.e.,  $q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)$  is absolutely continuous with respect to  $p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\gamma}, \boldsymbol{\tau})$ , we have

$$\mathbb{E}_{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \left[ \log p(\boldsymbol{\gamma} | \boldsymbol{\tau}) \right] = \mathbb{E}_{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \left[ \log \left( \frac{p(\boldsymbol{\gamma}, \mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\tau})}{p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\gamma}, \boldsymbol{\tau})} \right) \right]. \tag{26}$$

Since  $p(\boldsymbol{\gamma} | \boldsymbol{\tau})$  is a constant with respect to  $q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)$ , we can pull that out of the expectation

$$\begin{aligned}
 \log p(\boldsymbol{\gamma} | \boldsymbol{\tau}) &= \mathbb{E}_{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \left[ \log \left( \frac{p(\boldsymbol{\gamma}, \mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\tau})}{p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\gamma}, \boldsymbol{\tau})} \right) \right] \\
 &= \mathbb{E}_{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \left[ \log \left( \frac{p(\boldsymbol{\gamma}, \mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\tau}) q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)}{p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\gamma}, \boldsymbol{\tau}) q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \right) \right] \\
 &= \mathbb{E}_{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \left[ \log \left( \frac{p(\boldsymbol{\gamma}, \mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\tau})}{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \right) \right] + \mathbb{E}_{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \left[ \log \left( \frac{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)}{p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\gamma}, \boldsymbol{\tau})} \right) \right] \\
 &= \mathbb{E}_{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \left[ \log \left( \frac{p(\boldsymbol{\gamma}, \mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\tau})}{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \right) \right] + \mathbb{D}_{\text{KL}} \left[ q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega) \parallel p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\gamma}, \boldsymbol{\tau}) \right] \\
 &\geq \mathbb{E}_{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \left[ \log \left( \frac{p(\boldsymbol{\gamma}, \mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\tau})}{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \right) \right] \\
 &= \mathbb{E}_{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \left[ \log \left( p(\boldsymbol{\gamma}, \mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \boldsymbol{\tau}) \right) \right] - \mathbb{E}_{q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega)} \left[ \log \left( q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega) \right) \right], \tag{27}
 \end{aligned}$$

where the second-last inequality is due to the non-negativity of the KL divergence.

### D. Derivation of Eq. (18)

From Eq. (17), the probability  $p(\boldsymbol{\gamma}_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu)$  is given as

$$p(\boldsymbol{\gamma}_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) = \prod_{u=1}^d \mathbf{N} \left( [\boldsymbol{\gamma}_i]_u \mid \sum_{v=1}^d \mathbf{A}_i^{uv} [\mathbf{m}_\lambda^i]_v, \sum_{v=1}^d \mathbf{A}_i^{uv} [\mathbf{v}_\lambda^i]_v \right). \tag{28}$$

Assuming that  $p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \neq 0$ , taking log on both sides, we have

$$\log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) = \sum_{u=1}^d \log \mathbf{N} \left( [\gamma_i]_u \mid \sum_{v=1}^d \mathbf{A}_i^{uv} [\mathbf{m}_\lambda^i]_v, \sum_{v=1}^d \mathbf{A}_i^{uv} [\mathbf{v}_\lambda^i]_v \right). \quad (29)$$

Since  $\mathbf{A}_i$  is a permutation matrix, for every row  $u$ , there exists exactly one column  $v_*(u)$  – here, we make the dependence on the row  $u$  explicit – such that

$$\mathbf{A}_i^{uv} = \begin{cases} 1 & \text{if } v = v_*(u) \\ 0 & \text{otherwise} \end{cases}.$$

Using the above property, the following three equalities hold for any  $u = 1, \dots, d$

$$\sum_{v=1}^d \mathbf{A}_i^{uv} [\mathbf{m}_\lambda^i]_v = [\mathbf{m}_\lambda^i]_{v_*(u)}, \quad (30)$$

$$\sum_{v=1}^d \mathbf{A}_i^{uv} [\mathbf{v}_\lambda^i]_v = [\mathbf{v}_\lambda^i]_{v_*(u)}, \quad (31)$$

$$\sum_{v=1}^d \mathbf{A}_i^{uv} \log \mathbf{N}([\gamma_i]_u \mid [\mathbf{m}_\lambda^i]_v, [\mathbf{v}_\lambda^i]_v) = \log \mathbf{N}([\gamma_i]_u \mid [\mathbf{m}_\lambda^i]_{v_*(u)}, [\mathbf{v}_\lambda^i]_{v_*(u)}), \quad (32)$$

Continuing from Eq. (29), we have

$$\begin{aligned} \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) &= \sum_{u=1}^d \log \mathbf{N} \left( [\gamma_i]_u \mid \sum_{v=1}^d \mathbf{A}_i^{uv} [\mathbf{m}_\lambda^i]_v, \sum_{v=1}^d \mathbf{A}_i^{uv} [\mathbf{v}_\lambda^i]_v \right) \\ &= \sum_{u=1}^d \log \mathbf{N} \left( [\gamma_i]_u \mid [\mathbf{m}_\lambda^i]_{v_*(u)}, [\mathbf{v}_\lambda^i]_{v_*(u)} \right) = \sum_{u=1}^d \sum_{v=1}^d \mathbf{A}_i^{uv} \log \mathbf{N}([\gamma_i]_u \mid [\mathbf{m}_\lambda^i]_v, [\mathbf{v}_\lambda^i]_v), \end{aligned} \quad (33)$$

where the second equality is due to Eq. (30) and Eq. (31). The third equality is due to Eq. (32).

## E. Theoretical Analysis

This section presents the theoretical analysis of our personalized learning framework in Section 3. The analysis comprises two parts. In **Part I**, we will derive an auxiliary result that analyzes the behavior of the distribution<sup>7</sup>  $q(\gamma_* | \gamma)$  of the personalized solution model  $\gamma_*$  for target task  $\tau_*$ , which is derived from the optimized meta-model representation (Section 3.1). Here, the optimized representation corresponds to the universal maximizer of Eq. (16) or equivalently Eq. (19), which we sought to maximize previously in Section 3.4. This result will then be leveraged next in **Part II** to derive a non-trivial bound on the generalized performance gap between using a personalized model  $\gamma_* \sim q(\gamma_* | \gamma)$ ; and an oracle model  $\gamma_*^\circ$  which is defined to be a solution model that has best performance on  $\tau_*$ , and is unknown to us.

### E.1. Part I

First, recall that in part **B** of Section 3.4, the exact posterior<sup>8</sup>  $p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s} | \gamma)$  of our generative model (Fig. 1) is approximated with  $q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}) = q(\mathbf{w}_\emptyset) p(\mathbf{w} | \mathbf{s}) \prod_{\kappa=1}^e p(\mathbf{s}_\kappa)$  where  $q(\mathbf{w}_\emptyset)$  is parameterized separately from  $p(\mathbf{s}_\kappa)$  and  $p(\mathbf{w} | \mathbf{s})$  which are parts of our generative model. We can then show below (Lemma 1) that when the optimization of Eq. (16) converges to a universal maximizer, we can represent  $q(\mathbf{w}_\emptyset)$  explicitly in terms of the generative model’s components.

**Lemma 1.** *Suppose that the model evidence lower-bound  $\mathbf{L}(\alpha, \phi, \nu, \omega)$  defined in Eq. (16) achieves its (global) maximum,*

<sup>7</sup> $\gamma_*$  is also conditioned on the target task’s meta data  $\tau_*$  as well as the related tasks’ meta data  $\tau = \{\tau_i\}_{i=1}^n$  but we omitted these terms to avoid cluttering the notations. For the rest of the analysis, conditioning on  $\tau_*$  and  $\tau$  should therefore be implicitly understood.

<sup>8</sup>The posterior is over the base  $\mathbf{w}_\emptyset$  and specific  $\mathbf{w} = \{\mathbf{w}_i\}_{i=1}^n$  components, as well as the latter’s representative sets  $\mathbf{s} = \{\mathbf{s}_\kappa\}_{\kappa=1}^e$  (one set per encoding dimension  $\kappa$ ) as detailed previously in Section 3.2.

then the following is true:

$$q(\mathbf{w}_\emptyset; \omega) = \frac{1}{\mathbf{Z}} p(\mathbf{w}_\emptyset; \alpha) \prod_{i=1}^n \exp \left( \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \right] \right) \text{ almost everywhere,} \quad (34)$$

where  $\mathbf{Z} = \mathbb{E}_{p(\mathbf{w}_\emptyset; \alpha)} \left[ \prod_{i=1}^n \exp \left( \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \right] \right) \right]$  is the normalization term.

*Proof.* Let us first recall the expression of Eq. (16) below:

$$\begin{aligned} \mathbf{L}(\alpha, \phi, \nu, \omega) &= - \left( \mathbb{E}_{q(\mathbf{w}_\emptyset; \omega)} \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ - \sum_{i=1}^n \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \right] + \mathbb{D}_{\text{KL}} \left( q(\mathbf{w}_\emptyset; \omega) \parallel p(\mathbf{w}_\emptyset; \alpha) \right) \right) \\ &= - \mathbb{E}_{q(\mathbf{w}_\emptyset; \omega)} \left[ \log \exp \left( \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ - \sum_{i=1}^n \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \right] \right) + \log \frac{q(\mathbf{w}_\emptyset; \omega)}{p(\mathbf{w}_\emptyset; \alpha)} \right] \\ &= - \mathbb{E}_{q(\mathbf{w}_\emptyset; \omega)} \left[ \log q(\mathbf{w}_\emptyset; \omega) - \log \left( p(\mathbf{w}_\emptyset; \alpha) \prod_{i=1}^n \exp \left( \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \right] \right) \right) \right] \\ &= - \mathbb{D}_{\text{KL}} \left( q(\mathbf{w}_\emptyset; \omega) \parallel \frac{1}{\mathbf{Z}} p(\mathbf{w}_\emptyset; \alpha) \prod_{i=1}^n \exp \left( \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \right] \right) \right) + \log \mathbf{Z} \end{aligned} \quad (35)$$

where  $\mathbf{Z} = \mathbb{E}_{p(\mathbf{w}_\emptyset; \alpha)} \left[ \prod_{i=1}^n \exp \left( \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \right] \right) \right]$  is the normalization term, which is a constant with respect to  $\omega$ . Then, due to the non-negativity property of KL divergence, any choices of  $\omega$  that result in a positive divergence would not be optimal. Therefore, the model evidence lower-bound  $\mathbf{L}(\alpha, \phi, \nu, \omega)$  achieves its (global) maximum value if and only if the KL divergence is 0, which occurs if and only if:

$$\begin{aligned} q(\mathbf{w}_\emptyset; \omega) &= \frac{1}{\mathbf{Z}} p(\mathbf{w}_\emptyset; \alpha) \prod_{i=1}^n \exp \left( \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \right] \right) \text{ almost everywhere,} \\ q(\mathbf{w}_\emptyset; \omega) &\propto p(\mathbf{w}_\emptyset; \alpha) \left[ \prod_{i=1}^n \exp \left( \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \right] \right) \right] \text{ almost everywhere,} \end{aligned} \quad (36)$$

which completes our proof of Lemma 1.  $\square$

Leveraging this result, one can derive an immediate expression for  $q(\gamma_* | \gamma)$  as detailed next in Lemma 2 below.

**Lemma 2.** *Suppose that the model evidence lower-bound  $\mathbf{L}(\alpha, \phi, \nu, \omega)$  defined in Eq. (16) achieves its (global) maximum. Then, its induced predictive distribution  $q(\gamma_* | \gamma)$  of our framework coincides with  $\pi(\gamma_* | \gamma)$  where*

$$\pi(\gamma_* | \gamma) = \frac{\mathbb{E}_{p(\mathbf{w}_\emptyset; \alpha)} \left[ \exp \left( \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ \sum_{i=1}^n \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \right] \right) \int_{\mathbf{w}_*} p(\gamma_* | \mathbf{w}_*, \mathbf{w}_\emptyset; \nu) \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ p(\mathbf{w}_* | \mathbf{w}, \mathbf{s}; \phi) \right] d\mathbf{w}_* \right]}{\mathbb{E}_{p(\mathbf{w}_\emptyset; \alpha)} \left[ \exp \left( \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ \sum_{i=1}^n \log p(\gamma_i | \mathbf{w}_i, \mathbf{w}_\emptyset; \nu) \right] \right) \right]}$$

which does not depend on the parameterization  $\omega$  of  $q(\mathbf{w}_\emptyset; \omega)$ .

*Proof.* Let us first recall our framework's approximate expression of the predictive distribution  $q(\gamma_* | \gamma)$  for the target model

$\gamma_*$  given observations of existing, relevant models  $\gamma$ :

$$\begin{aligned}
 q(\gamma_*|\gamma) &= \int p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}|\gamma) \left[ p(\mathbf{w}_*|\mathbf{w}, \mathbf{s}; \phi) p(\gamma_*|\mathbf{w}_*, \mathbf{w}_\emptyset; \nu) \right] d\mathbf{w}_* d\mathbf{w} d\mathbf{s} d\mathbf{w}_\emptyset \\
 &\simeq \int q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega) \left[ p(\mathbf{w}_*|\mathbf{w}, \mathbf{s}; \phi) p(\gamma_*|\mathbf{w}_*, \mathbf{w}_\emptyset; \nu) \right] d\mathbf{w}_* d\mathbf{w} d\mathbf{s} d\mathbf{w}_\emptyset \\
 &= \int q(\mathbf{w}_\emptyset; \omega) \left[ p(\mathbf{w}, \mathbf{s}; \phi) p(\mathbf{w}_*|\mathbf{w}, \mathbf{s}; \phi) p(\gamma_*|\mathbf{w}_*, \mathbf{w}_\emptyset; \nu) d\mathbf{w}_* d\mathbf{w} d\mathbf{s} \right] d\mathbf{w}_\emptyset \\
 &= \int q(\mathbf{w}_\emptyset; \omega) \mathbb{E}_{p(\mathbf{w}, \mathbf{s}; \phi)} \left[ \int p(\mathbf{w}_*|\mathbf{w}, \mathbf{s}; \phi) p(\gamma_*|\mathbf{w}_\emptyset, \mathbf{w}_*; \nu) d\mathbf{w}_* \right] d\mathbf{w}_\emptyset, \tag{37}
 \end{aligned}$$

where the first step is by our variational approximation scheme, the second step follows from the surrogate approximation  $p(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}|\gamma) \simeq q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s})$ , and the third step follows from our choice of the surrogate  $q(\mathbf{w}_\emptyset, \mathbf{w}, \mathbf{s}; \omega) = q(\mathbf{w}_\emptyset; \omega) p(\mathbf{w}|\mathbf{s}; \phi) \prod_{\kappa=1}^e p(\mathbf{s}_\kappa; \phi) = q(\mathbf{w}_\emptyset; \omega) p(\mathbf{w}, \mathbf{s}; \phi)$ . Finally, plugging the expression of  $q(\mathbf{w}_\emptyset)$  in Lemma 1 (at a global maximizer) into Eq. 37 yields the desired form of  $\pi(\gamma_*|\gamma)$ .  $\square$

Lemma 2 thus establishes the form of the predictive distribution  $q(\gamma_*|\gamma)$  at optimality. We will exploit this result in **Part II** below to establish a non-trivial bound on the expected performance gap between an oracle model and a sampled model from  $q(\gamma_*|\gamma)$  on target task  $\tau_*$ , which confirms the key result that we stated previously in Theorem 1 (Section 4) of the main text.

## E.2. Part II

Let  $\gamma_* \in \mathcal{H}$  denotes a model for  $\tau_*$  in the hypothesis space  $\mathcal{H}$  and  $\mathbf{D}_*$  denotes the (unknown) data distribution of  $\tau_*$ . Let  $\ell_\gamma(\mathbf{x}, y)$  denotes the non-negative evaluation loss of predicting  $\gamma(\mathbf{x})$  when the ground-truth is  $y$ . Examples of such evaluation loss function includes the squared error  $\ell_\gamma(\mathbf{x}, y) = (\gamma(\mathbf{x}) - y)^2$  (for regression) and the 0-1 loss  $\ell_\gamma(\mathbf{x}, y) = \mathbb{I}(\gamma(\mathbf{x}), y)$  (for classification). The generalized loss of a model  $\gamma_*$  on  $\tau_*$  is defined as,

$$\mathbf{G}(\gamma_*) \triangleq \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \ell_{\gamma_*}(\mathbf{x}, y) \right]. \tag{38}$$

We assume there exists an oracle model  $\gamma_*^\circ \in \mathcal{H}$  for  $\tau_*$  that achieves the minimum generalized loss, which is defined as:

$$\gamma_*^\circ \triangleq \arg \min_{\gamma_* \in \mathcal{H}} \mathbf{G}(\gamma_*) \text{ or equivalently, } \mathbf{G}(\gamma_*^\circ) = \min_{\gamma_*} \mathbf{G}(\gamma_*). \tag{39}$$

We also assume that the evaluation loss admits the triangle inequality such that  $\ell_\gamma(\mathbf{x}, y) \leq \ell_\gamma(\mathbf{x}, \gamma'(\mathbf{x})) + \ell_{\gamma'}(\mathbf{x}, y)$  and  $\ell_\gamma(\mathbf{x}, \gamma'(\mathbf{x})) \leq \ell_\gamma(\mathbf{x}, y) + \ell_{\gamma'}(\mathbf{x}, y)$ . It is easy to see that both examples above of the squared and 0-1 losses admit such triangle inequalities. Then, let  $q$  be the abbreviation for  $q(\gamma_*|\gamma)$  (as defined above). The generalized loss of sampling solution model from  $q$  for target task  $\tau_*$  with distribution  $\mathbf{D}_*$  is defined as

$$\mathbf{G}(q) = \mathbb{E}_{\gamma_* \sim q} \left[ \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \ell_{\gamma_*}(\mathbf{x}, y) \right] \right]. \tag{40}$$

To bound the gap between  $\mathbf{G}(\gamma_*^\circ)$  and  $\mathbf{G}(q)$ , we first establish the following immediate results in Lemma 3 and Lemma 4.

**Intuition.** In lay terms, Lemma 3 re-iterates a well-known result that can be used in Lemma 4 to bound the expected generalized loss  $\mathbf{G}(q)$  of a model distribution in terms of auxiliary functions. This in turn can be bound by a small fraction of the oracle's generalized loss  $\mathbf{G}(\gamma_*^\circ)$  under a mild assumption regarding the performance of the oracle model on a random input. This is formally stated in Lemma 5 – see its premise assumption and its discussion that follows – which is the stepping stone to establish our key result in Theorem 1 (Section 4).

**Lemma 3.** For any function  $g(\mathbf{x})$  and any two distributions  $p(\mathbf{x})$  and  $\pi(\mathbf{x})$  on a certain measurable space, we have

$$\mathbb{E}_{\mathbf{x} \sim p} \left[ g(\mathbf{x}) \right] \leq \mathbb{D}_{\text{KL}} \left( p \parallel \pi \right) + \log \mathbb{E}_{\mathbf{x} \sim \pi} \left[ \exp \left( g(\mathbf{x}) \right) \right]. \tag{41}$$



*Proof.* This is a well-known result often referred to as the **Change of Measure** inequality (Seldin et al., 2015). For the sake of self-containment, we concisely reproduce the its proof below:

$$\begin{aligned}
 \mathbb{E}_{\mathbf{x} \sim p} \left[ g(\mathbf{x}) \right] &= \mathbb{E}_{\mathbf{x} \sim p} \left[ \log \left( \frac{p(\mathbf{x})}{\pi(\mathbf{x})} \times \frac{\pi(\mathbf{x})}{p(\mathbf{x})} \times \exp(g(\mathbf{x})) \right) \right] \\
 &= \mathbb{E}_{\mathbf{x} \sim p} \left[ \log \frac{p(\mathbf{x})}{\pi(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p} \left[ \log \left( \frac{\pi(\mathbf{x})}{p(\mathbf{x})} \times \exp(g(\mathbf{x})) \right) \right] \\
 &= \mathbb{D}_{\text{KL}}(p \parallel \pi) + \mathbb{E}_{\mathbf{x} \sim p} \left[ \log \left( \frac{\pi(\mathbf{x})}{p(\mathbf{x})} \times \exp(g(\mathbf{x})) \right) \right] \\
 &\leq \mathbb{D}_{\text{KL}}(p \parallel \pi) + \log \mathbb{E}_{\mathbf{x} \sim p} \left[ \frac{\pi(\mathbf{x})}{p(\mathbf{x})} \times \exp(g(\mathbf{x})) \right] = \mathbb{D}_{\text{KL}}(p \parallel \pi) + \log \mathbb{E}_{\mathbf{x} \sim \pi} \left[ \exp(g(\mathbf{x})) \right], \quad (42)
 \end{aligned}$$

where the inequality follows from applying Jensen inequality, which completes our proof of Lemma 3.  $\square$

**Lemma 4.** For any reference distribution  $\pi(\gamma_* | \gamma)$ , let  $\mathbf{F}_{\gamma_*^\circ}(\pi) \triangleq \mathbb{E}_{\mathbf{x} \sim \mathbf{D}_*} \mathbb{E}_{\gamma_* \sim \pi} [\exp(\ell_{\gamma_*}(\mathbf{x}, \gamma_*^\circ(\mathbf{x})))]$ . We have

$$\left| \mathbf{G}(q) - \mathbf{G}(\gamma_*^\circ) \right| \leq \mathbb{D}_{\text{KL}}(q \parallel \pi) + \log \mathbf{F}_{\gamma_*^\circ}(\pi). \quad (43)$$

*Proof.* By definition,  $\mathbf{G}(\gamma_*^\circ) \leq \mathbf{G}(\gamma_*)$  for any choice of  $\gamma_*$ . Hence,  $\mathbf{G}(\gamma_*^\circ) = \mathbb{E}_{\gamma_* \sim q} [\mathbf{G}(\gamma_*^\circ)] \leq \mathbb{E}_{\gamma_* \sim q} [\mathbf{G}(\gamma_*)] = \mathbf{G}(q)$ . On the other hand, we also have

$$\begin{aligned}
 \mathbf{G}(q) &= \mathbb{E}_{\gamma_* \sim q} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \ell_{\gamma_*}(\mathbf{x}, y) \right] \\
 &\leq \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim q} \left[ \ell_{\gamma_*}(\mathbf{x}, y) + \ell_{\gamma_*}(\mathbf{x}, \gamma_*^\circ(\mathbf{x})) \right] \right] \\
 &= \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \ell_{\gamma_*^\circ}(\mathbf{x}, y) \right] + \mathbb{E}_{\mathbf{x} \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim q} \left[ \ell_{\gamma_*}(\mathbf{x}, \gamma_*^\circ(\mathbf{x})) \right] \right] \\
 &= \mathbf{G}(\gamma_*^\circ) + \mathbb{E}_{\mathbf{x} \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim q} \left[ \ell_{\gamma_*}(\mathbf{x}, \gamma_*^\circ(\mathbf{x})) \right] \right] \\
 &\leq \mathbf{G}(\gamma_*^\circ) + \mathbb{E}_{\mathbf{x} \sim \mathbf{D}_*} \left[ \mathbb{D}_{\text{KL}}(q \parallel \pi) + \log \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, \gamma_*^\circ(\mathbf{x}))) \right] \right] \\
 &= \mathbf{G}(\gamma_*^\circ) + \mathbb{D}_{\text{KL}}(q \parallel \pi) + \mathbb{E}_{\mathbf{x} \sim \mathbf{D}_*} \log \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, \gamma_*^\circ(\mathbf{x}))) \right] \\
 &\leq \mathbf{G}(\gamma_*^\circ) + \mathbb{D}_{\text{KL}}(q \parallel \pi) + \log \mathbb{E}_{\mathbf{x} \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, \gamma_*^\circ(\mathbf{x}))) \right] \right] \\
 &= \mathbf{G}(\gamma_*^\circ) + \mathbb{D}_{\text{KL}}(q \parallel \pi) + \log \mathbf{F}_{\gamma_*^\circ}(\pi), \quad (44)
 \end{aligned}$$

where the first inequality is due to the triangle inequality of the evaluation loss, the second inequality follows from the change of measure inequality in Lemma 3, and the last inequality follows from Jensen inequality. Thus, combining this with the previously proved statement, we have

$$\mathbf{G}(\gamma_*^\circ) \leq \mathbf{G}(q) \leq \mathbf{G}(\gamma_*^\circ) + \mathbb{D}_{\text{KL}}(q \parallel \pi) + \log \mathbf{F}_{\gamma_*^\circ}(\pi), \quad (45)$$

which implies  $|\mathbf{G}(q) - \mathbf{G}(\gamma_*^\circ)| \leq \mathbb{D}_{\text{KL}}(q \parallel \pi) + \log \mathbf{F}_{\gamma_*^\circ}(\pi)$  as desired. This completes our proof of Lemma 4.  $\square$

**Intuition.** Next, we show that when the oracle model is highly accurate in the sense that the chance for it to incur a loss worse than a fraction of its generalized loss is vanishingly small – see Eq. (46),  $\log \mathbf{F}_{\gamma_*^\circ}(\pi)$  can be upper-bounded by a small fraction of  $\mathbf{G}(\gamma_*^\circ)$  and a log data term. This is formally stated in Lemma 5 below.

**Lemma 5.** *Assume there exist constants  $c > 0$  and  $r > 1$  such that*

$$\Pr_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left( \ell_{\gamma_*^\circ}(\mathbf{x}, y) > \frac{1}{r} \mathbf{G}(\gamma_*^\circ) \right) \leq c \times \exp \left( \frac{1}{r} \mathbf{G}(\gamma_*^\circ) - 2 \sup_{\gamma, \mathbf{x}, y} [\ell_\gamma(\mathbf{x}, y)] \right). \quad (46)$$

Then, it can be shown that for  $\mathbf{F}_{\gamma_*^\circ}(\pi)$  as defined in Lemma 4 above,

$$\log \mathbf{F}_{\gamma_*^\circ}(\pi) \leq \frac{1}{r} \mathbf{G}(\gamma_*^\circ) + \log \left( \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, y)) \right] \right] + c \right). \quad (47)$$

*Proof.* To prove this, we first expand the expression of  $\mathbf{F}_{\gamma_*^\circ}(\pi)$  using the triangle inequality of the evaluation loss,

$$\begin{aligned} \log \mathbf{F}_{\gamma_*^\circ}(\pi) &= \log \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp \left( \ell_{\gamma_*}(\mathbf{x}, \gamma_*^\circ(\mathbf{x})) \right) \right] \right] \\ &\leq \log \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp \left( \ell_{\gamma_*^\circ}(\mathbf{x}, y) + \ell_{\gamma_*}(\mathbf{x}, y) \right) \right] \right] \\ &= \log \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \exp \left( \ell_{\gamma_*^\circ}(\mathbf{x}, y) \right) \times \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp \left( \ell_{\gamma_*}(\mathbf{x}, y) \right) \right] \right] \triangleq \log \mathbf{F}_{\gamma_*^\circ}^\dagger(\pi). \end{aligned} \quad (48)$$

Now, let  $\psi \triangleq \sup_{\gamma, \mathbf{x}, y} \ell_\gamma(\mathbf{x}, y)$  and  $\delta = \Pr(\ell_{\gamma_*^\circ}(\mathbf{x}, y) > (1/r)\mathbf{G}(\gamma_*^\circ))$ . We can now split the space of  $(\mathbf{x}, y)$  into two parts: (a)  $\mathbb{D}_*^{(a)} = \{(\mathbf{x}, y) \mid \ell_{\gamma_*^\circ}(\mathbf{x}, y) \leq (1/r)\mathbf{G}(\gamma_*^\circ)\}$ ; and (b)  $\mathbb{D}_*^{(b)} = \{(\mathbf{x}, y) \mid \ell_{\gamma_*^\circ}(\mathbf{x}, y) > (1/r)\mathbf{G}(\gamma_*^\circ)\}$ . As such, if  $(\mathbf{x}, y)$  is to be sampled uniformly in part (a), its sampling distribution is  $\mathbf{D}_*^{(a)}(\mathbf{x}, y) = \mathbb{I}((\mathbf{x}, y) \in \mathbb{D}_*^{(a)}) \times \mathbf{D}_*(\mathbf{x}, y) / (1 - \delta)$ . Likewise, if  $(\mathbf{x}, y)$  is in part (b), its sampling distribution is  $\mathbf{D}_*^{(b)}(\mathbf{x}, y) = \mathbb{I}((\mathbf{x}, y) \in \mathbb{D}_*^{(b)}) \times \mathbf{D}_*(\mathbf{x}, y) / \delta$ . Thus,

$$\begin{aligned} \mathbf{F}_{\gamma_*^\circ}^\dagger(\pi) &\leq (1 - \delta) \exp \left( \frac{1}{r} \mathbf{G}(\gamma_*^\circ) \right) \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*^{(a)}} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, y)) \right] \right] \\ &\quad + \delta \exp(\psi) \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*^{(b)}} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, y)) \right] \right] \\ &\leq (1 - \delta) \exp \left( \frac{1}{r} \mathbf{G}(\gamma_*^\circ) \right) \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*^{(a)}} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, y)) \right] \right] + \delta \exp(2\psi) \\ &\leq \exp \left( \frac{1}{r} \mathbf{G}(\gamma_*^\circ) \right) \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, y)) \right] \right] + \delta \exp(2\psi) \\ &\leq \exp \left( \frac{1}{r} \mathbf{G}(\gamma_*^\circ) \right) \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, y)) \right] \right] + c \times \exp \left( \frac{1}{r} \mathbf{G}(\gamma_*^\circ) - 2\psi + 2\psi \right) \\ &= \exp \left( \frac{1}{r} \mathbf{G}(\gamma_*^\circ) \right) \left[ \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, y)) \right] \right] + c \right], \end{aligned} \quad (49)$$

where the last inequality follows from our assumption above in Eq. (46). Then, taking logarithm on both sides of Eq. (49) and plugging the result in Eq. (48) yield Eq. (47), which completes our proof of Lemma 5.  $\square$

**Key Result.** Finally, chaining up the above results, we can obtain the following key result which is the main thrust of our theoretical contribution here. It also reveals important insights that explain the rationalities behind the (somewhat artificial) formulation of our personalized learning framework, as discussed below.

**Theorem 1.** Assume there exist constants  $c > 0$  and  $r > 1$  such that (same assumption from Lemma 5)

$$\Pr_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left( \ell_{\gamma_*^\circ}(\mathbf{x}, y) > \frac{1}{r} \mathbf{G}(\gamma_*^\circ) \right) \leq c \times \exp \left( \frac{1}{r} \mathbf{G}(\gamma_*^\circ) - 2 \sup_{\gamma, \mathbf{x}, y} [\ell_\gamma(\mathbf{x}, y)] \right). \quad (50)$$

Then, for any reference distribution  $\pi(\gamma_*|\gamma)$ , we have with probability at least  $1 - \delta$  over the sampling of an  $m$ -shot dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$  where  $(\mathbf{x}_i, y_i) \sim \mathbf{D}_*$ , the following inequality holds simultaneously for all choices of  $q(\gamma_*|\gamma)$ :

$$\left| \mathbf{G}(q) - \mathbf{G}(\gamma_*^\circ) \right| \leq \mathbb{D}_{\text{KL}}(q|\pi) + \frac{1}{r} \mathbf{G}(\gamma_*^\circ) + \log \left( \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp[\ell_{\gamma_*}(\mathbf{x}_i, y_i)] \right] + \mathbf{O} \left( \left[ \frac{\log \left( \frac{4m}{\delta} \right)}{2m-1} \right]^{\frac{1}{2}} \right) \right)$$

where (as defined previously)  $\gamma_*^\circ$  denotes the oracle solution model of task  $\tau_*$ .

**Proof.** First, we will chain up the results of Lemma 4 and Lemma 5 to obtain the following:

$$\begin{aligned} \left| \mathbf{G}(q) - \mathbf{G}(\gamma_*^\circ) \right| &\leq \mathbb{D}_{\text{KL}}(q|\pi) + \log \mathbf{F}_{\gamma_*^\circ}(\pi) \\ &\leq \mathbb{D}_{\text{KL}}(q|\pi) + \frac{1}{r} \mathbf{G}(\gamma_*^\circ) + \log \left( \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, y)) \right] \right] + c \right). \end{aligned} \quad (51)$$

Now, to bound the last term on the RHS of the above inequality, we define the following auxiliary loss  $h_{\gamma_*}(\mathbf{x}, y) = \exp(\ell_{\gamma_*}(\mathbf{x}, y) - \sup \ell_{\gamma_*}(\mathbf{x}, y))$  where the supremum is over the choice of  $\gamma$  and  $(\mathbf{x}, y)$ . As such,  $h_{\gamma_*}(\mathbf{x}, y) \in (0, 1)$ .

Then, applying the PAC-Bayes result in (McAllester, 1999) to relate the generalized and empirical loss of sampling solution model  $\gamma_* \sim \pi$  for task  $\tau_*$  with data distribution  $\mathbf{D}_*$  using instance loss  $h_{\gamma_*}(\mathbf{x}, y)$ , we have with probability at least  $1 - \delta$  over the choice of  $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$ ,

$$\begin{aligned} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim \pi} [h_{\gamma_*}(\mathbf{x}, y)] \right] &\leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\gamma_* \sim \pi} [h_{\gamma_*}(\mathbf{x}_i, y_i)] + \left( \frac{\mathbb{D}_{\text{KL}}(\pi|\pi) + \log \left( \frac{4m}{\delta} \right)}{2m-1} \right)^{\frac{1}{2}} \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\gamma_* \sim \pi} [h_{\gamma_*}(\mathbf{x}_i, y_i)] + \left( \frac{\log \left( \frac{4m}{\delta} \right)}{2m-1} \right)^{\frac{1}{2}}. \end{aligned} \quad (52)$$

Note that we use a special case of the PAC-Bayesian bound here where the target and reference distributions are the same. Thus, multiplying both sides of the above inequality with  $\exp(\sup \ell_{\gamma_*}(\mathbf{x}, y))$  yields

$$\begin{aligned} \mathbb{E}_{(\mathbf{x}, y) \sim \mathbf{D}_*} \left[ \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}, y)) \right] \right] &\leq \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}_i, y_i)) \right] \\ &\quad + \exp \left( \sup \ell_{\gamma_*}(\mathbf{x}_i, y_i) \right) \left( \frac{\log \left( \frac{4m}{\delta} \right)}{2m-1} \right)^{\frac{1}{2}} \\ &= \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp(\ell_{\gamma_*}(\mathbf{x}_i, y_i)) \right] + \mathbf{O} \left( \left[ \frac{\log \left( \frac{4m}{\delta} \right)}{2m-1} \right]^{\frac{1}{2}} \right). \end{aligned} \quad (53)$$

Finally, taking the logarithm of both sides, plugging the result into Eq. (51), and absorbing the constant  $c$  into the big- $\mathbf{O}$  notation complete our proof of Theorem 1. Thus, setting  $C_m = \mathbf{O}((\log(4m/\delta)/(2m-1))^{1/2})$  reproduces the key result that we mentioned in Theorem 1 in the main text. Now, combining this result with our previously established result in **Part I** reveals an interesting insight into the fine-tuning part of our personalized learning framework as discussed in Corollary 1 below. We note that an informal discussion of this had been provided previously Section 4 above.

**Corollary 1.** Assuming the same assumption of Theorem 1 regarding the oracle solution model for  $\tau_*$  and constructing the reference distribution  $\pi(\gamma_*|\gamma)$  following its established parameterization in Lemma 2, with arbitrarily high probability,

the gap between the generalized losses of  $\gamma_*^\circ$  and  $q(\gamma_*|\gamma)$  is at most  $\mathbf{G}(\gamma_*^\circ)/r$  when  $m \rightarrow \infty$  and the optimization of Eq. (16) reaches a global optimizer. Furthermore, when  $m$  is small, the established bound in Theorem 1 also reveals how the process of fine-tuning a personalized model using the few-shot dataset  $\{(\mathbf{x}_i, y_i)\}_{i=1}^m$  can be incorporated seamlessly into the learning of the meta-model representation to provably control the bound gap.

**Proof.** To see this, we first take a closer look at the proved bound of Theorem 1,

$$\left| \mathbf{G}(q) - \mathbf{G}(\gamma_*^\circ) \right| \leq \mathbb{D}_{\text{KL}}(q \parallel \pi) + \frac{1}{r} \mathbf{G}(\gamma_*^\circ) + \log \left( \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp \left[ \ell_{\gamma_*}(\mathbf{x}_i, y_i) \right] \right] \right) + \mathbf{O} \left( \left[ \frac{\log \left( \frac{4m}{\delta} \right)}{2m-1} \right]^{\frac{1}{2}} \right)$$

Now, using the results of **Part I** (see Lemmas 1 and 2) we know that when the optimization of Eq. (16) reaches a universal optimizer,  $q$  converges to  $\pi$  in Lemma 2. At that point,  $\mathbb{D}_{\text{KL}}(q \parallel \pi) = 0$  and the first term on the RHS of the above bound disappears. As for the remaining terms, taking limit when  $m \rightarrow \infty$  on both sides of the reduced bound further zero out the log term, thus leaving only a fraction  $1/r$  of the oracle loss.

Secondly, in case  $m$  is small, the excess loss term depends heavily on  $\log \left( \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\gamma_* \sim \pi} [\exp[\ell_{\gamma_*}(\mathbf{x}_i, y_i)]] \right)$  which immediately suggests a loss function for fine-tuning and an appropriate scale to combine with the personalization loss. In particular, it can be seen from the above arguments that optimizing  $q$  and  $\pi$  simultaneously via minimizing

$$\mathbf{H}(q, \pi) = -\mathbf{L}(q) + \log \left( \frac{1}{m} \sum_{i=1}^m \mathbb{E}_{\gamma_* \sim \pi} \left[ \exp \left[ \ell_{\gamma_*}(\mathbf{x}_i, y_i) \right] \right] \right) \quad (54)$$

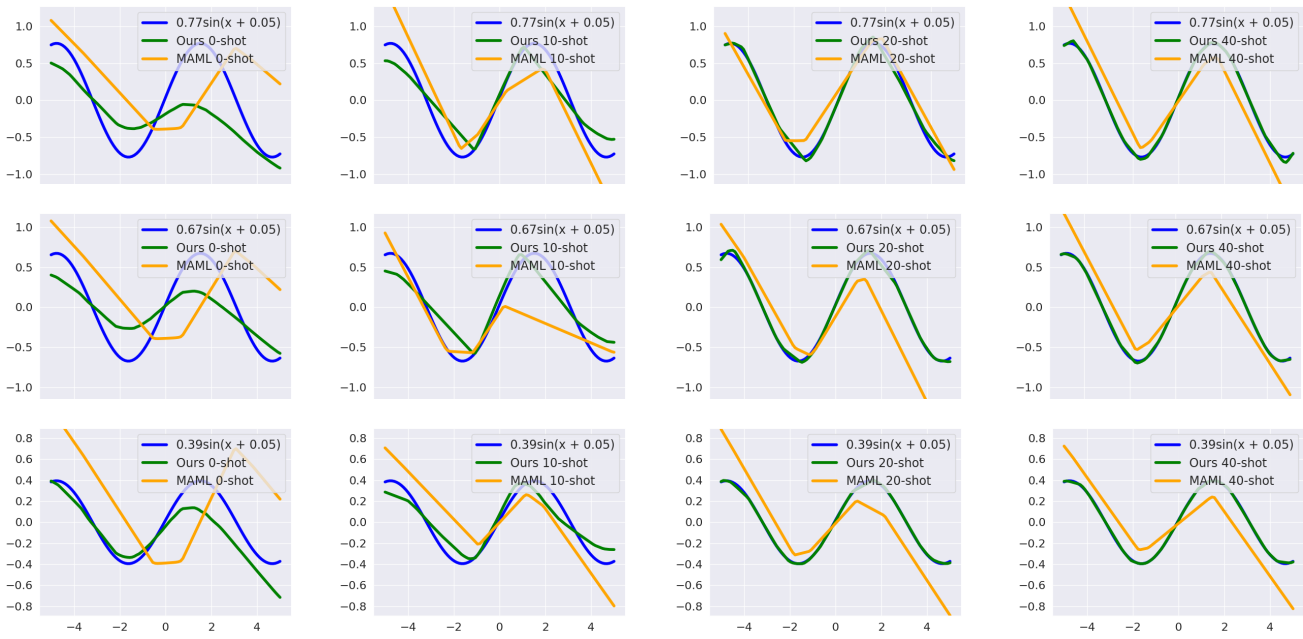
will result in a performance gap bound by  $\mathbf{G}(\gamma_*^\circ)/r$  and a minimized excess loss involving the few-shot dataset for fine-tuning. Here, we ignore the constant term that scales with double log, which is negligible. This is true because the parameterizations of  $q$  and  $\pi$  are decoupled (see Lemma 1) and the data term in the above loss does not depend on  $q$ , which effectively forces  $\mathbb{D}_{\text{KL}}(q \parallel \pi)$  to be zero if  $(q, \pi)$  is a universal minimizer of  $\mathbf{H}(q, \pi)$ . Thus, optimizing  $\mathbf{H}$  instead of the original ELBO objective in Eq. (16) allows us to not only zero out the divergence term but also reduce the data fit further on the few-shot dataset in the new task context (i.e., fine-tuning).

**Remark.** The bound above also meets our sanity check in the limit of information and quality of the oracle model. To elaborate, in the limit of information, both the divergence and data terms disappear as argued above, thus leaving a fraction  $1/r$  of the oracle loss  $\mathbf{G}(\gamma_*^\circ)$ . Now, if the target task  $\tau_*$  can be solved with high accuracy by an oracle model bounded within a model space  $\mathcal{H}$ , one would expect  $r$  (see the assumption in Lemma 5) to be large and  $c$  to be small while  $\mathbf{G}(\gamma_*^\circ)$  to be vanishingly small. Both of which contribute towards zeroing out the remaining loss terms of our bound in the limit of information and oracle quality, thus resulting in a personalized model distribution whose generalized performance is arbitrarily close to that of an oracle model.

In light of this, the bound in Theorem 1 also suggests an intuition regarding an intricate trade-off between the bound gap (in the limit of information) and the expressiveness of the model space. On one hand, if the model space is highly expressive, it is likely that our assumption for the oracle model will hold with large  $r$  which tighten the excess loss but in exchange, it will be harder for the divergence term to zero out via optimizing the meta-model loss in Eq. (16) as there are now more parameters to be optimized. On the other hand, if the model space is reduced in complexity, it is easier for the optimization of Eq. (16) to find a global optimizer and zero out the divergence term but with the reduced complexity, the oracle model (i.e., the best in the defined model space) becomes less accurate and hence, the assumption might only hold with large  $c$  and small  $r$ . This has the effect of pushing the bound towards being vacuous. When this happens, model fusion become perhaps a less well-defined task since even the oracle model cannot sufficiently solve the task with high accuracy.

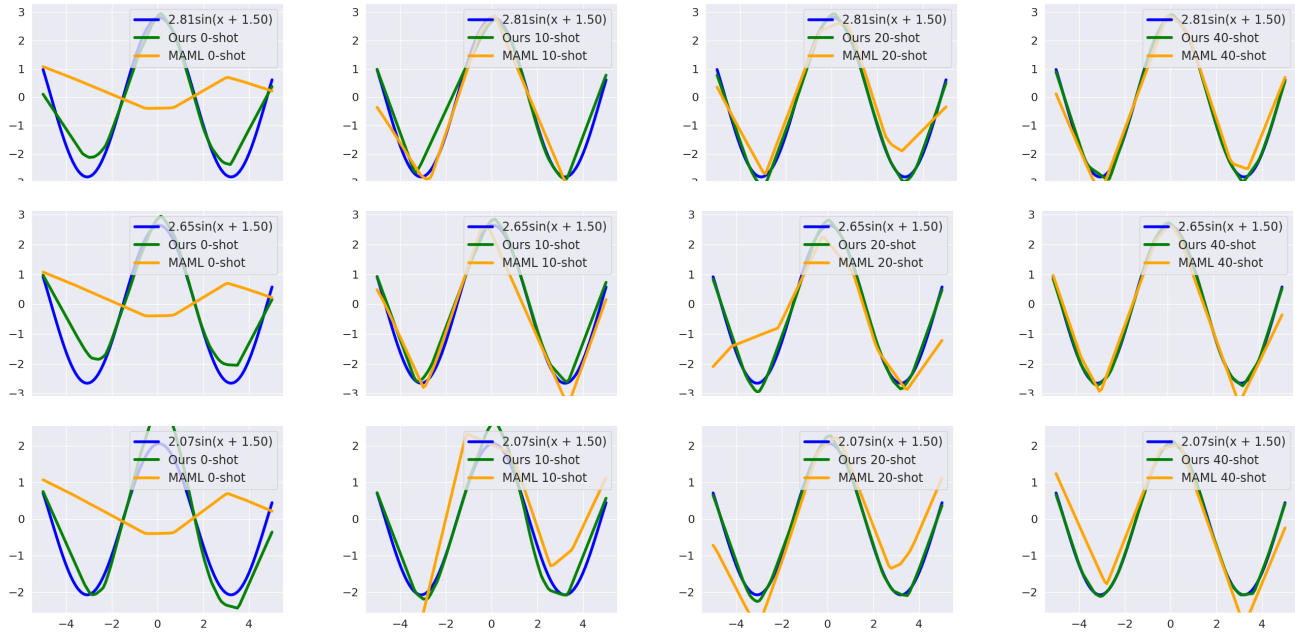
## F. Extra Visualized Performance Plots across Different Sine Tasks in Different Sub-Domains

Figures 5, 6 and 7 provide additional visualized fitting performance of our personalized model towards various sine tasks across all 5 polarized sine domains (corresponding to those reported in Appendix A).

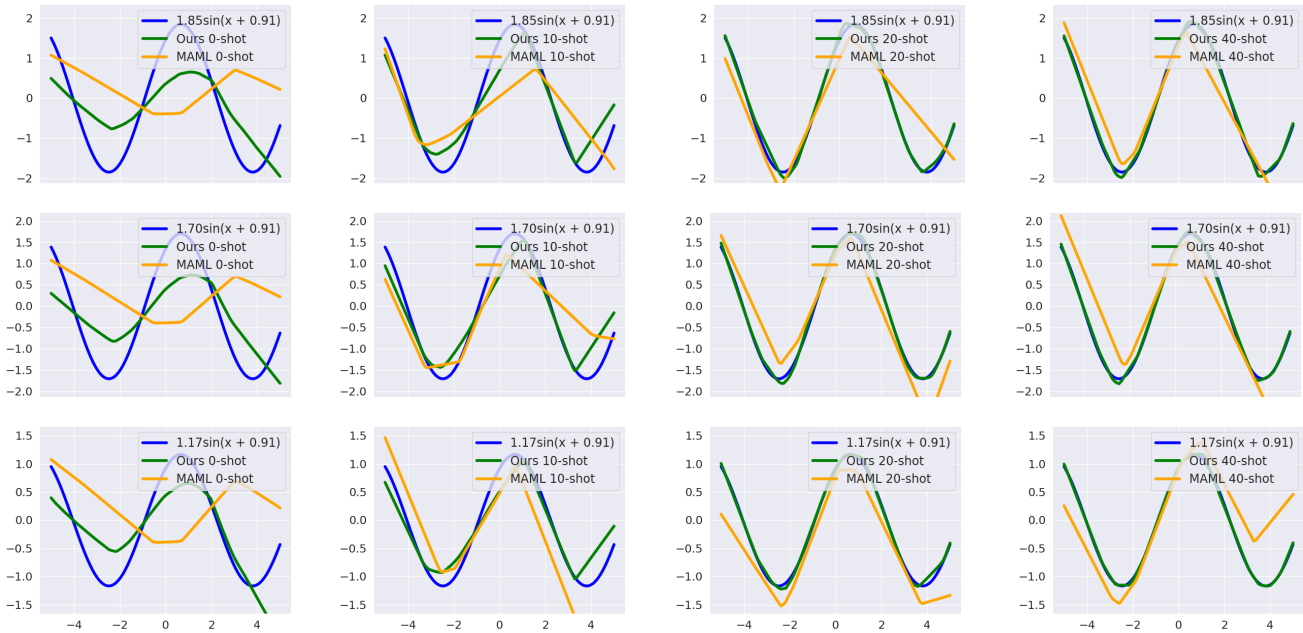


Sub-Domain 1:  $a \sim U[0, 1]$ ,  $b = 0.91$

Figure 5. (Part 1) Additional visual excerpts demonstrating how well the warm-start neural net with architecture [1-100-1] generated by MAML and our method fit 2 unseen sine functions  $y = a \sin x + b$  of sub-domain 1 in 0-, 10-, 20- and 40-shot settings.

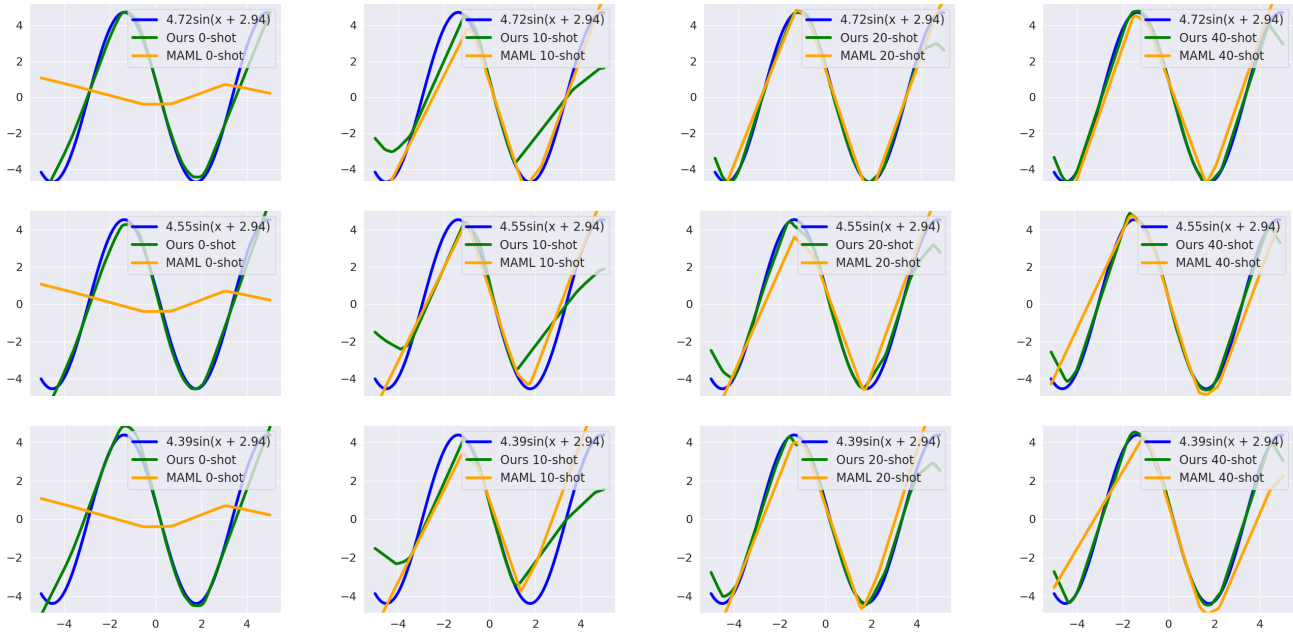


Sub-Domain 3:  $a \sim U[2, 3], b = 1.50$

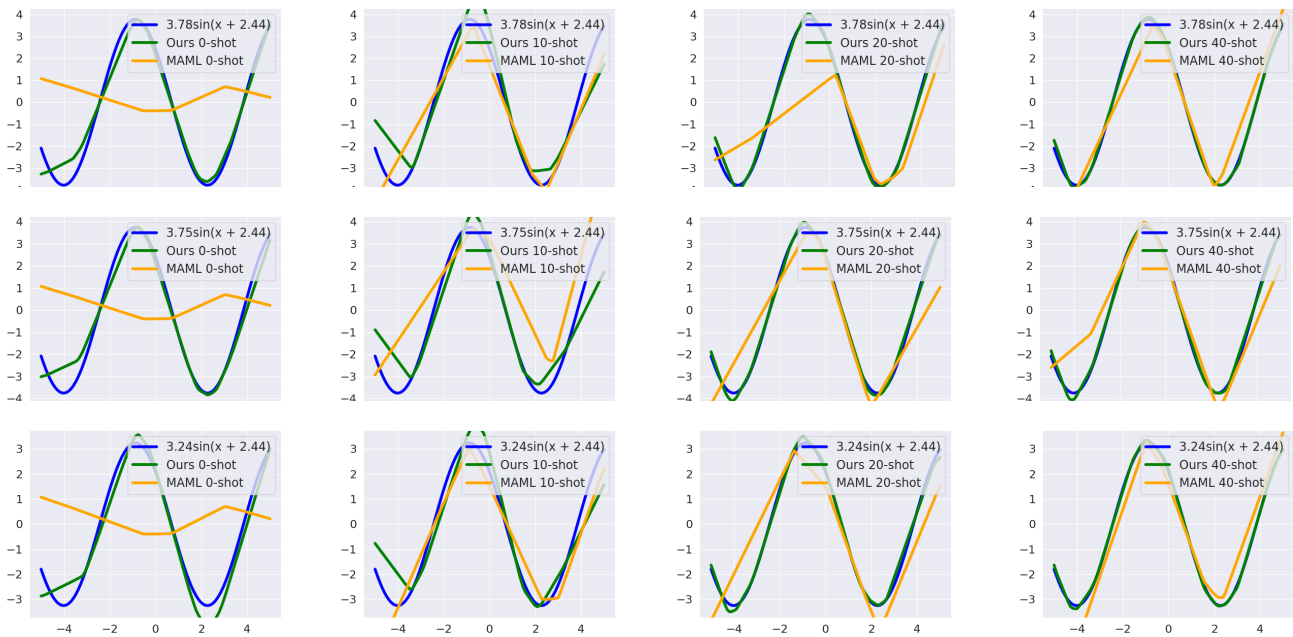


Sub-Domain 2:  $a \sim U[1, 2], b = 0.91$

Figure 6. (Part 2) Additional visual excerpts demonstrating how well the warm-start neural net with architecture [1-100-1] generated by MAML and our method fit 2 unseen sine functions  $y = a \sin x + b$  of sub-domains 2 and 3 in 0-, 10-, 20- and 40-shot settings.



Sub-Domain 5:  $a \sim U[4, 5]$ ,  $b = 2.94$



Sub-Domain 4:  $a \sim U[3, 4]$ ,  $b = 2.44$

Figure 7. (Part 3) Additional visual excerpts demonstrating how well the warm-start neural net with architecture [1-100-1] generated by MAML and our method fit 2 unseen sine functions  $y = a \sin x + b$  of sub-domains 4 and 5 in 0-, 10-, 20- and 40-shot settings.