

Table 1: Starting from one auxiliary task (Exemplar-MT), we keep increasing the number of auxiliary tasks from one to four by adding one auxiliary task at a time in the order of {Rotation, VAT, EntMin}.

K	1 (+Exemplar-MT)	2 (+Rotation)	3 (+VAT)	4(+EM)
CIFAR-10	17.02	13.68	12.47	<b>11.80</b>
SVHN	9.65	5.89	5.00	<b>4.70</b>

Table 2: Effect of ARML on harmful tasks.

Exemplar-MT	+ Rotation (harmful)
17.02	<b>16.95</b> (w/ ARML)
	83.76 (w/o ARML)

- 1 We would like to thank all the reviewers for writing the insightful comments, especially during this difficult time.
- 2 **R1.1 - L110 “ $p^*(\theta)$  is ‘covered’ by  $p(\mathcal{T}_m|\theta)$ ” is not clearly explained:** Thanks! By ‘covered’, we mean  $p(\mathcal{T}_m|\theta)$  has  
3 high density both in the support set  $S$  of  $p^*(\theta)$ , and in some regions outside  $S$ . We will clarify this in the final version.
- 4 **R1.2 - Our contribution and difference from previous works:** Thanks. We proposed ARML for auxiliary task  
5 reweighting so that the least data is required to find the true parameter. Although ARML and some previous algorithms  
6 (*e.g.* CosineSim, OL\_AUX) both update task weights based on the similarity between main/auxiliary task gradient, the  
7 previous algorithms mainly focus on the **training loss**. For example, CosineSim discards the harmful tasks so that the  
8 training loss of the main task will not increase in one step. OL\_AUX updates the task weights so that the training loss  
9 decreases the fastest. However, lowering training loss may cause overfitting, especially when training data is scarce.  
10 In contrast, ARML is guaranteed to find a good prior so that the least data is required to find the parameter which  
11 generalizes the best (*i.e.* having the lowest **test error**). The superiority of ARML is verified in the experiments.
- 12 **R1.3 - Ablation on number of auxiliary tasks:** Thanks. We ablate on the number of auxiliary tasks in SSL, and  
13 observe the results in Table 1. The error rates decreases when each new task is added.
- 14 **R1.4 - Line 239-242 the text here is confusing:** Thanks. For domain generalization baseline, the model is trained only  
15 with source domain (auxiliary) data. In ‘Baseline + ARML’, for fair comparison, we stick to the same training process,  
16 *i.e.*, updating  $\theta$  using only auxiliary loss (NOT using the main task loss), and the main task data is used to adjust the  
17 task weights  $\alpha$  which brings the result improvement. We will add more elaboration on this in the final version.
- 18 **R1.5 - Non-classification experiments:** Thanks for the suggestion! We will try other tasks, *e.g.*, reinforcement learning.
- 19 **R2.1 - More explanation in the Algorithm section:** Thanks! (1) We use mini-batch sampling for main/auxiliary task  
20 data. In each mini-batch, we sample the same amount of two kinds of data. (2) Yes, during early epochs the sampling  
21 process is still warming up, and after that the estimation is more accurate. We will add more details in the section.
- 22 **R2.2 - Discussion on the efficiency of the proposed method:** Thanks. The only overhead of ARML is the update of  
23  $\alpha$ . This has little extra cost because the main/auxiliary task gradients are already calculated when updating  $\theta$ , and we  
24 only have to calculate the gradient w.r.t.  $\alpha$ . We observed a  $< 10\%$  extra cost for PyTorch, and nearly no extra cost for  
25 TensorFlow (because of the different autograd strategies of the two frameworks). The efficiency is similar to previous  
26 algorithms (*e.g.* OL\_AUX), and much higher than GridSearch. We will update this in the final version.
- 27 **R3.1 - More experiments on multi-label classification:** Thanks! We will add more experiments including 1) analyzing  
28 if the learned face attribute relationship aligns with human’s intuition, 2) varying the number of auxiliary attributes.
- 29 **R3.2 - Can ARML prevent harmful tasks?:** Thanks! Ideally, ARML can discard a harmful task by lowering its  
30 weight to 0. For verification, we turn the Rotation loss in S4L into its negative, making it a harmful task (which means  
31 when training we are actually increasing the loss). Meanwhile we keep the Exemplar-MT loss unchanged. We find that  
32 ARML lowers the weight of Rotation to 0, getting the same error rate as when only Exemplar-MT is applied (Table 2).
- 33 **R3.3 - Why does VAT + EntMin works better than S4L + ARML? Have you tried MOAM?:** Since ARML already  
34 finds the optimal task weights, the performance is mainly limited by the auxiliary tasks themselves. For example, if we  
35 further add VAT & EntMin as auxiliary tasks (which makes it MOAM), then ARML surpasses SOTA (R1.3 & Table 1).
- 36 **R4.1 - Justification behind switching sampling probability (Theorem 1):** Thanks. We think there may be a small  
37 misunderstanding, because Theorem 1 does not mean that we are minimizing a lower bound of the true objective. Our  
38 true objective is to find the optimal task weights  $\alpha^*$  so that  $D_{\text{KL}}(p^* \parallel p_{\alpha^*})$  is the smallest. Theorem 1 states that,  
39 if we choose the weights  $\hat{\alpha}$  (from the optimization of the surrogate objective), then the corresponding true objective  
40  $D_{\text{KL}}(p^* \parallel p_{\hat{\alpha}})$  is also very small, *i.e.*, upper-bounded near the optimal value  $D_{\text{KL}}(p^* \parallel p_{\alpha^*})$ .
- 41 **R4.2 - The concept of “true prior”:** This is a good point! The “true prior” is an objective prior reflecting the  
42 randomness of the environment. For example, for image classification problems with pictures taken in different places,  
43 the model parameter should stay similar to extract the same features, while bearing some slight changes (*e.g.* different  
44 image statistics in BN). The “true prior” can also be regarded as a Dirac delta if not considering such randomness.
- 45 **R4.3 - Details:** Thank you! We will correct the typos in the final version. 4) WRN-28-2 is used for SSL experiments,  
46 and ResNet18 is used for others. 5) We take three runs for each result. 6) Thanks. We will add the illustrations.