

## A APPENDIX

### A.1 JUSTIFYING THE INTRODUCTION OF A META-HEAD

*Proof.* To arrive at Equation 7 we start with the closed form solution for  $\nabla_{w_i} \mathcal{L}_{T^*}^{val}(\theta^*)$  and then introduce approximations in order to produce Equation 7. First, note that :

$$\frac{\partial \mathcal{L}_{T^*}^{val}(\theta^*(\mathbf{w}))}{\partial w_i} = \left( \nabla_{\theta} \mathcal{L}_{T^*}^{val}(\theta^*(\mathbf{w})) \right)^T \left( \nabla_{w_i} \theta^*(\mathbf{w}) \right) \quad [\text{Chain rule}] \quad (8)$$

To get  $\nabla_{w_i} \theta^*(\mathbf{w})$  we invoke the Cauchy Implicit Function Theorem (IFT) as with Lorraine et al. (2020); Navon et al. (2020); Liao et al. (2018):

$$\begin{aligned} \nabla_{w_i} \theta^*(\mathbf{w}) &= \left[ \nabla_{\theta}^2 \mathcal{L}_{\text{total}}(\theta^*(\mathbf{w})) \right]^{-1} \left[ \nabla_{w_i} \nabla_{\theta} \mathcal{L}_{\text{total}}(\theta^*(\mathbf{w})) \right] \quad [\text{IFT}] \\ &= \left[ \nabla_{\theta}^2 \mathcal{L}_{\text{total}}(\theta^*(\mathbf{w})) \right]^{-1} \left[ \nabla_{w_i} \nabla_{\theta} \left( w^* \mathcal{L}_{T^*}(\theta^*(\mathbf{w})) + \sum_{T_i \in \mathbb{T}_{\text{aux}}} w_i \mathcal{L}_{T_i}(\theta^*(\mathbf{w})) \right) \right] \\ &= \left[ \nabla_{\theta}^2 \mathcal{L}_{\text{total}}(\theta^*(\mathbf{w})) \right]^{-1} \left[ \nabla_{\theta} \mathcal{L}_{T_i}(\theta^*(\mathbf{w})) \right] \quad [\text{Only terms with } w_i \text{ survive}] \end{aligned}$$

Bringing it all together, we get :

$$\frac{\partial \mathcal{L}_{T^*}^{val}(\theta^*(\mathbf{w}))}{\partial w_i} = \left( \nabla_{\theta} \mathcal{L}_{T^*}^{val}(\theta^*(\mathbf{w})) \right)^T \left( \left[ \nabla_{\theta}^2 \mathcal{L}_{\text{total}}(\theta^*(\mathbf{w})) \right]^{-1} \left[ \nabla_{\theta} \mathcal{L}_{T_i}(\theta^*(\mathbf{w})) \right] \right) \quad (9)$$

□

Computing  $\nabla_{w_i} \mathcal{L}_{T^*}^{val}(\theta^*)$  from Equation 9 is computationally unwieldy since we would not only have to optimize  $\theta$  to convergence for every step of  $w_i$  but we would also have to invert the Hessian of a typically large model. Our middle ground between Equations 9 and 6 (Equation 7) makes use of the following approximations:

- We approximate the inverse Hessian with the identity. This approximation is not new; we follow previous work like Lorraine et al. (2020) (Table 3) who explore the use of this approximation because of computational efficiency.

$$\left[ \nabla_{\theta}^2 \mathcal{L}_{\text{total}}(\theta^*(\mathbf{w})) \right]^{-1} = \lim_{i \rightarrow \infty} \sum_{j=0}^i \left( \mathbf{I} - \nabla_{\theta}^2 \mathcal{L}_{\text{total}}(\theta^*(\mathbf{w})) \right)^j \approx \mathbf{I}$$

We are assuming the contribution of terms with  $i > 0$  are negligible.

- Instead of training the whole network to convergence, at each time-step, we fix the body of the network and train a special head  $\phi^*$  to convergence on a small batch of end-task training data. We then use  $[\theta_{\text{body}}; \phi^*]$  as a proxy for  $\theta^*$ . This is a computationally feasible work-around to training all of  $\theta$  to convergence to get a single step gradient estimate. Especially in the continued pre-training setting where a pre-trained *generalist* model like BERT is used as  $\theta_{\text{body}}$ , this approximation is reasonable. To our knowledge, we are the first to suggest this approximation.

$$\nabla_{\theta} \mathcal{L}_{T^*}^{val}(\theta^*) \rightarrow \nabla_{\theta} \mathcal{L}_{T^*}^{val}([\theta_{\text{body}}; \phi^*])$$

- Above, we have approximated  $\theta^* = [\theta_{\text{body}}; \phi^*]$ . Since  $\phi^*$  is only used to evaluate end-task ( $T^*$ ) validation data, it means  $\theta$  remains unchanged with respect to the training data for task  $T_i$ . Thus  $\nabla_{\theta} \mathcal{L}_{T_i}([\theta_{\text{body}}; (\phi^*, \dots, \phi^t)]) = \nabla_{\theta} \mathcal{L}_{T_i}([\theta_{\text{body}}; \phi^t]) = \nabla_{\theta} \mathcal{L}_{T_i}(\theta)$

Bringing it all together, we get Equation 7 repeated here:

$$\frac{\partial \mathcal{L}_{T^*}^{val}(\theta^*(\mathbf{w}))}{\partial w_i} \approx (\nabla_{\theta} \mathcal{L}_{T_i})^T (\nabla_{\theta} \mathcal{L}_{T^*}^{val}([\theta_{\text{body}}; \phi^*]_t))$$

## A.2 CALCULATING P-VALUES FROM PERMUTATION TEST

We used the permutation test (Good, 2005; Dror et al., 2018) to test for statistical significance. For each test, we generate 10000 permutations to calculate significance level. This is sufficient to converge to a stable p-value without being a computational burden. We chose this over the common student t-test because :

1. We have only 10 runs per algorithm and permutation tests are more robust at low sample size
2. Permutation test is assumption free. Student t-tests assume that the samples are normally distributed
3. Permutation test is robust to variance in the samples, so even though error-bars can overlap, we still establish significant differences in the samples. Variance in our results is expected due to small dataset sizes of end-tasks.

## A.3 ALGORITHM FOR META-TARTAN

---

### Algorithm 1: End-task Aware Training via Meta-learning (META-TARTAN)

---

**Require:**  $T^*$ ,  $\mathbf{T}_{\text{aux}}$ : End-task, Set of auxiliary pre-training tasks

**Require:**  $\eta, \beta_1, \beta_2$ : Step size hyper-parameters

**Initialize :**

Pre-trained RoBERTa as shared network body,  $\theta_{\text{body}}$

Task weightings:  $w^*, w_i = \frac{1}{|\mathbf{T}_{\text{aux}}|+1}$

**Randomly initialize :**

end-task head as  $\phi'$

meta head for end-task as  $\phi^*$

task head,  $\phi^i$ , for each  $T_i \in \mathbf{T}_{\text{aux}}$

**while not done do**

$B_{\text{tr}}^* \sim T_{\text{train}}^*$  // Sample a batch from end-task

$g_{\theta}^*, g_{\phi}^* \leftarrow [\nabla_{\theta}, \nabla_{\phi'}] \left( \mathcal{L}_{T^*}(\theta, \phi', B_{\text{tr}}^*) \right)$  // Get end-task grads

$g_{\theta}^i, g_{\phi}^i \leftarrow [\nabla_{\theta}, \nabla_{\phi^i}] \left( \mathcal{L}_{T_i}(\theta, \phi^i, B_i) \right)$  // Get task grads.  $\forall i \in [n], B_i \sim T_i$

// Learn a new meta head

$\phi^* \leftarrow \text{estimate\_meta\_head}(B_{\text{tr}}^*, \beta_2, \theta, \phi^*)$  //  $B_{\text{tr}}^* \sim T_{\text{train}}^*$

$g_{\text{meta}}^* \leftarrow \nabla_{\theta} \mathcal{L}_{T^*}(\theta, \phi^*, B_{\text{val}}^*)$  //  $B_{\text{val}}^* \sim T_{\text{val}}^*$

// Update task weightings

$w^* \leftarrow w^* + \eta \cos(g_{\text{meta}}^*, g_{\theta}^*)$

$w_i \leftarrow w_i + \eta \cos(g_{\text{meta}}^*, g_{\theta}^i)$

// Update task parameters

$\alpha^*, \alpha_1, \dots, \alpha_{|\mathbf{T}_{\text{aux}}|} = \text{softmax}(w^*, w_1, \dots, w_{|\mathbf{T}_{\text{aux}}|})$

Update  $\theta_{\text{body}} \leftarrow \theta_{\text{body}} - \beta_1 (\alpha^* g_{\theta}^* + \sum_i \alpha_i g_{\theta}^i)$

Update  $\left( \phi_i \leftarrow \phi_i - \beta_2 g_{\phi}^i \right), \left( \phi' \leftarrow \phi' - \beta_2 g_{\phi}^* \right)$

**end**

**Result :**  $\theta, \phi'$

---

## A.4 VISION EXPERIMENTS

We validate that the gains from end-task Aware Training are not siloed to only learning from text. We conduct an experiment comparing end-task aware training on images to its end-task agnostic variant. We use the Cifar100 dataset (Krizhevsky et al., 2009). We use the Medium-Sized Mammals superclass (one of the 20 coarse labels) as our main task whilst the other 19 super classes are used as auxiliary data. Our primary task is thus a 5-way classification task of images different types of

Method	Medium-Sized Mammals
Regular (Task-Agnostic) Pre-training	46.7 <sub>2.2</sub>
MT-TARTAN	51.3 <sub>1.2</sub>
META-TARTAN	<b>52.3<sub>3.8</sub></b>

Table 5: We report averages across 3 random seeds. Best average task accuracy is bolded. medium-sized mammals whilst whilst the remaining 95 classes are grouped into a single auxiliary task.

As can be seen from Table 5, being end-task aware improves over task agnostic pre-training. We find that, again, when our auxiliary task consist of solely domain data and no task data, META-TARTAN performs better than MT-TARTAN (as measured by averaged performance).

#### A.5 FULL TAPT TABLE WITH SIGNIFICANCE LEVELS

We repeat Table 1 and provide details about levels of statistical signifiance.

Task	TAPT	MT-TARTAN	$p$ -values	META-TARTAN	$p$ -values
ACL-ARC	67.74 <sub>3.68</sub>	<b>70.48</b> <sub>4.42</sub>	0.040	70.08 <sub>4.70</sub>	0.069
SCIERC	79.53 <sub>1.93</sub>	<b>80.81</b> <sub>0.74</sub>	0.038	<b>81.48</b> <sub>0.82</sub>	0.005
CHEMPROT	82.17 <sub>0.065</sub>	<b>84.29</b> <sub>0.63</sub>	0.000	<b>84.49</b> <sub>0.50</sub>	0.000

Table 6: Duplicate of Table 1. Significance levels as computed from the permutation test. All  $p$ -values are relative to the TAPT column. Statistically significant performance( $p$ -value from permutation test  $< 0.05$ ), is boldfaced

#### A.6 FULL DAPT/DAPT+TAPT TABLE

We repeat Table 3 and provide details about levels of statistical signifiance.

Task	DAPT	DAPT+TAPT	MT-TARTAN	$p$ -values	META-TARTAN	$p$ -values
ACL-ARC	68.60 <sub>2.62</sub>	69.12 <sub>5.76</sub>	71.58 <sub>1.65</sub>	0.110	71.05 <sub>2.37</sub>	0.174
SCIERC	76.44 <sub>1.19</sub>	77.62 <sub>1.38</sub>	<b>81.02</b> <sub>1.24</sub>	0.000	<b>81.41</b> <sub>1.70</sub>	0.000
CHEMPROT	80.76 <sub>0.54</sub>	78.22 <sub>0.74</sub>	<b>83.77</b> <sub>0.60</sub>	0.000	<b>83.38</b> <sub>0.89</sub>	0.000

Table 7: Duplicate of Table 2. Significance levels as computed from the permutation test. All  $p$ -values are relative to  $\max(\text{DAPT}, \text{DAPT} + \text{TAPT})$ . Statistically significant performance( $p$ -value from permutation test  $< 0.05$ ), is boldfaced