

## A Appendix

### A.1 Additional Results

More detailed numbers on the T0 [Sanh et al. \[2022\]](#) and SuperNI [Wang et al. \[2022a\]](#) datasets using different backbones, and different adapter layouts over the base model are found in Table 4. **Multi-Task params** is the number of additional parameters that must be conserved

Model	Multi-Task Params	Adaptation Params	Avg. Test
<b>T0 Dataset</b>			
<i>Backbone T5-XL-LM</i>			
Multi-Task Full Finetuning + LoRA	2.8B	2.2M	68.9 <sub>x.x</sub>
(IA) <sup>3</sup>	540K	540K	62.4 <sub>0.4</sub>
AdapterSoup	84M	2.2M	62.1 <sub>1.0</sub>
LoRA	2.2M	2.2M	66.0 <sub>1.6</sub>
LoRA-big	35M	35M	65.4 <sub>0.9</sub>
Poly-z	17M	3.5K	66.4 <sub>0.3</sub>
Poly	17M	2.2M	68.0 <sub>1.0</sub>
MHR-z (64 h)	17M	220K	68.3 <sub>0.8</sub>
MHR (64 h)	17M	2.2M	69.1 <sub>1.0</sub>
<i>Backbone T0-3B</i>			
T-Few <a href="#">Liu et al. [2022]</a>	540K	540K	66.2 <sub>0.5</sub>
AdapterSoup	84M	2.2M	66.1 <sub>0.6</sub>
LoRA	2.2M	2.2M	67.4 <sub>0.8</sub>
LoRA-big	35M	35M	68.0 <sub>0.8</sub>
Poly-z	17M	3.5K	65.3 <sub>1.0</sub>
Poly	17M	2.2M	69.0 <sub>0.8</sub>
MHRz (64 h)	17M	220K	68.4 <sub>1.2</sub>
MHR (8 h)	17M	2.2M	69.3 <sub>1.2</sub>
<i>Backbone T0-3B light version : (k, v, ff layers only)</i>			
l-LoRA (rank 1)	934K	934K	66.2 <sub>0.9</sub>
l-LoRA (rank 16)	15M	15M	67.6 <sub>1.1</sub>
AdapterSoup (l-LoRA)	35M	934K	64.9 <sub>1.0</sub>
l-Poly-z	7.5M	2.1K	62.9 <sub>1.2</sub>
l-Poly	7.5M	934K	68.0 <sub>0.5</sub>
l-MHRz (32 h)	7.5M	74K	66.8 <sub>1.1</sub>
l-MHR (8 h)	7.5M	934K	68.5 <sub>0.7</sub>
<b>SuperNI Dataset</b>			Rouge-L
<i>Backbone T5-XL-LM light version : (k, v, ff layers only)</i>			
l-LoRA	934K	934K	67.6 <sub>0.8</sub>
l-LoRA-big	18M	18M	67.2 <sub>0.7</sub>
l-Poly-z	7.5M	2.1K	64.6 <sub>0.3</sub>
l-Poly	7.5M	934K	67.8 <sub>0.8</sub>
l-MHRz (64 h)	7.5M	147K	68.0 <sub>0.2</sub>
l-MHR (8 h)	7.5M	934K	68.5 <sub>0.3</sub>

Table 4: (top) Results on T0 dataset [Sanh et al. \[2022\]](#), we report the mean of the best validation accuracy for each test task, when evaluated every 50 train epochs. T-Few is our reproduction of the results in [Liu et al. \[2022\]](#). LoRA-big means a LoRA adapter with a larger rank. (bottom) Results on SuperNatural Instructions dataset.

after multi-task pretraining to enable transfer to a downstream task. **Adaptation Params** refer to the number of parameters required to learn a new downstream task. For e.g. Poly and MHR, the multi-task parameters includes the learned modules, but not the routing over the training tasks, as these are not required for transfer on a new task. Moreover, variants which average the learned modules prior to fine-tuning (MHR- $\mu$  and Poly- $\mu$ ) will have both multi-task and adaptation parameters equal to that of a single shared adapter, since after multi-task pretraining one can average the modules.

## 520 A.2 Navigating the parameter efficiency / performance trade-off of tuning only 521 the routing

522 Here we provide additional results on how different routing based methods can be more  
523 expressive when only learning a new routing function (over *frozen* modules) to adapt to a  
524 new task.

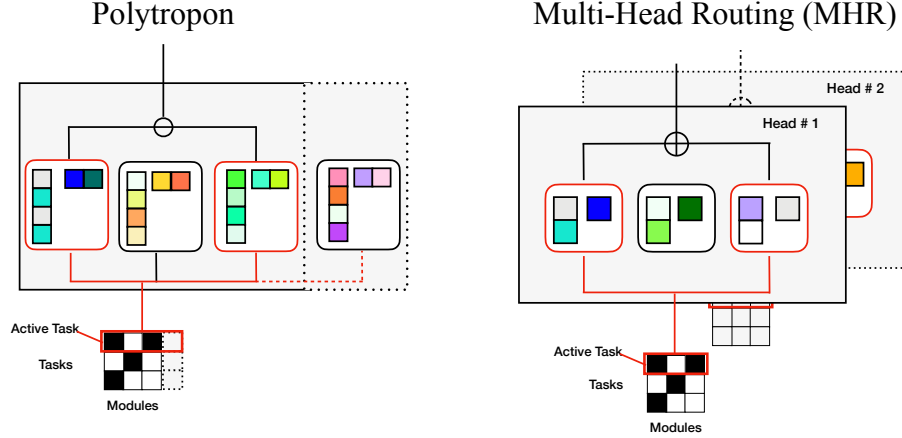


Figure 4: Different ways to control the expressivity of routing based methods. *Left* : In Polytropon, one can only add additional modules, resulting in a linear parameter increase. *Right* : In MHR, additional heads only introduce routing matrices  $\mathbf{Z}$ , resulting in a negligible parameter increase.

525 In Fig. 4 (left), we see that in order to build more expressive routing functions  $\mathbf{Z}$ , in Poly  
526 one can only do so by increasing the number of skills at each layer. However, this has a  
527 significant impact on the number of multi-task parameters which much be kept in order to  
528 perform few-shot transfer. MHR on the other hand, can increase routing capacity in a much  
529 more parameter efficient way.

### 530 A.2.1 On the granularity of routing tensor in MHR

531 Here we provide additional results when modifying the granularity of  $\mathbf{Z}$  for MHR. We see that  
532 one can easily trade-off more parameters for better performance.

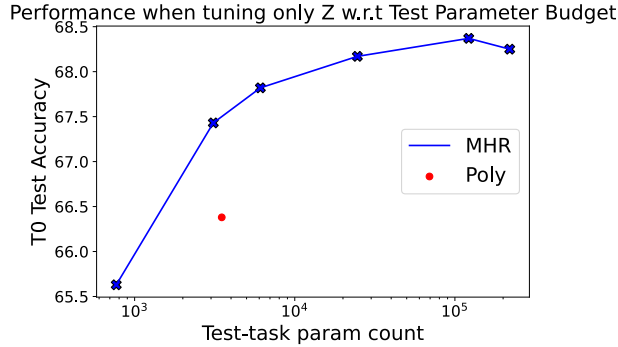


Figure 5: Routing-Only Fine-Tuning (MHR- $z$ )

## 533 **B Broader Impact**

534 In our work, we focus on advancing parameter-efficient fine-tuning methods for cross-task  
535 generalization. While our research primarily addresses technical challenges and performance  
536 improvements, when applying such methods, it is crucial to consider the potential negative  
537 societal impacts. Specifically, we believe that prior to applying our proposed adaptation  
538 method, critically examining the potential biases and ethical implications of the underlying  
539 large language model, and the data itself must be properly addressed. This includes issues  
540 related to fairness, privacy, and the spread of misinformation.