LST: Ladder Side-Tuning for Parameter and Memory Efficient Transfer Learning

Anonymous Author(s) Affiliation Address email

Abstract

Fine-tuning large pre-trained models on downstream tasks has been adopted in a 1 variety of domains recently. However, it is costly to update the entire parameter 2 set of large pre-trained models. Although recently proposed parameter-efficient 3 transfer learning (PETL) techniques allow updating a small subset of parameters 4 inside a pre-trained backbone network for a new task, they only reduce the training 5 memory requirement by up to 30%. This is because the gradient computation 6 for the trainable parameters still requires backpropagation through the large pre-7 trained backbone model. To address this, we propose Ladder Side-Tuning (LST), 8 a new PETL technique that can also reduce training memory requirements by 9 more substantial amounts. Unlike existing parameter-efficient methods that insert 10 additional parameters inside backbone networks, we train a ladder side network, a 11 12 small and separate network that takes intermediate activations as input via shortcut connections (ladders) from backbone networks and makes predictions. LST 13 has significantly lower memory requirements than previous methods, because it 14 does not require backpropagation through the backbone network, but instead only 15 through the side network and ladder connections. We conduct our experiments 16 on various models (T5 and CLIP-T5) and tasks (GLUE, VQA, GQA, NLVR², 17 MSCOCO). LST saves 69% of the memory costs to fine-tune the whole network, 18 while other methods only save 26% of that in similar parameter usages (hence, 19 2.7x more memory savings). Moreover, LST outperforms Adapter and LoRA in 20 21 a low-memory regime. To further show the advantage of this memory efficiency, we also apply LST to large T5 models (T5-large, T5-3B), attaining better GLUE 22 performance than full fine-tuning and other PETL methods. The trend also holds in 23 the experiments on vision-and-language tasks, where LST performs comparably to 24 other PETL methods when training a similar number of parameters but only uses 25 46% of their GPU memory. 26

27 **1** Introduction

Recently, large-scale pre-training and fine-tuning of transformers [48] have been successful in various 28 domains [8, 31, 38, 33, 47, 7, 19, 2, 10]. As the model size grows rapidly, fine-tuning the entire 29 parameter set of the large pre-trained model has become very costly. Parameter-efficient transfer 30 learning (PETL) [26, 21, 34, 35, 29, 45, 52, 16, 22, 17, 36, 54, 14, 46, 53] is a recent research 31 32 direction for the online setting, where tasks arrive in an online or multi-task setting. The goal is to build a system that performs well on all tasks without training an entire new model for every 33 34 new task. Concretely, PETL methods select a small subset of pre-trained parameters and/or insert a few parameters to a pre-trained network and update those parameters for new tasks, while freezing 35 most of the original parameters. In the natural language processing (NLP), computer vision (CV), 36 and vision-and-language (VL) domains, two types of parameters have been commonly updated 37

Submitted to 36th Conference on Neural Information Processing Systems (NeurIPS 2022). Do not distribute.



Figure 1: Comparison between full finetuning, Adapter, LoRA, BitFit, and Ladder Side-Tuning over GLUE tasks. The y-axis is the average accuracy of 8 GLUE tasks, while the x-axis is the GPU memory usage during training. Unless specially stated, we use the T5-base in the figure.



Figure 2: Comparison between transfer learning with (a) Adapters, (b) Prompt Tuning, and our (b) Ladder Side-Tuning (LST). LST reduces memory usage by removing the need of backpropgation through backbone networks.

for parameter-efficient transfer learning: (1) *Adapters* [21, 35, 34]: small modules inserted into transformer blocks; (2) *Prompt* [29, 26]: small parameters concatenated with input embeddings.

However, while parameter-efficient techniques reduce the number of parameters to update, they
 do not reduce the memory requirement during training by much (up to 30%). In other words, if a

⁴¹ about reduce the memory requirement during training by much (up to 5070). In other words, if a ⁴² large pre-trained language model does not fit on a GPU, these techniques usually do not help to fit

43 the model on the GPU. Since the updated parameters are inside the backbone language models, to

44 calculate gradients for these parameters for backpropagation, we still need to run the backward pass

through the large pre-trained language models. This prevents PETL methods from being applied to

⁴⁶ many real-world applications with limited computational resources.

To address this issue, we propose Ladder Side-Tuning (LST), a memory-efficient PETL method. 47 LST separates trainable parameters from the backbone model to construct a side network, which is 48 responsible for adapting the entire model to new tasks. Concretely, we train a *ladder* side network, 49 50 a lightweight network that takes intermediate activations via shortcut connections (ladders) from 51 the backbone networks as input and makes predictions. As shown in Fig. 2, unlike previous (a) adapters, (b) prompt tuning methods, (c) our LST does not add trainable parameters inside the 52 pre-trained model. This completely eliminates the need for expensive backpropagation of a large 53 backbone network and saves substantial memory during transfer learning. Instead of initializing the 54 side network's weights randomly, we utilize a network pruning technique to retrieve a smaller pruned 55 network and use it as the side network. In addition to our standard design of the side network, we 56 also boost the efficiency of LST by dropping layers of the side network. We empirically demonstrate 57 that the layer dropping can significantly improve both memory and parameter efficiency without 58 59 sacrificing performance. Furthermore, during inference, even though we have to propagate forward through two distinct networks, LST does not necessarily use more inference time because the same 60 61 level of the backbone network and the side network can be computed in parallel.

We conduct comprehensive studies on LST using diverse NLP and VL tasks, namely, GLUE [49], 62 VQA [15], GQA [23], NLVR² [44] and MSCOCO [6]. Overall, in GLUE experiments, LST saves 63 69% of the GPU memory that is needed for fine-tuning the entire backbone model, saving 2.7x 64 memory compared against Adapter and LoRA. Also, LST outperforms other PETL methods in 65 a low-memory regime. To take advantage of this memory efficiency, we also apply LST to large 66 language models (T5-large and T5-3B) and find that it outperforms other techniques when GPU 67 memory utilization is similar. The findings still hold in the VL experiments; LST is not only a 68 method that can fit in a 16GB GPU with 7.5% trainable parameters, but it also has comparable or 69 better performance than other PETL methods. To justify our design of LST, we conduct ablation 70 studies on initialization strategies and alternatives to add shortcut connections. The results reveal that 71 these components considerably help performance, while only adding minor computation overhead. 72

73 2 Related Work

74 2.1 Parameter-efficient Transfer Learning (PETL)

PETL for NLP. In the past few years, large pre-trained language models have made huge success in 75 NLP. Parameter-efficient transfer learning (PETL) is a research direction that reduces computational 76 cost of adapting large pre-trained models to new tasks, by avoiding updates of the entire parameters. 77 A popular line of work on PETL is to add a few trainable parameters and only tune them. For 78 example, Adapters [42, 21] are small bottleneck modules that are inserted into transformer layers, and 79 80 experiments have shown that training adapters with layer normalization layers is sufficient to achieve 81 full fine-tuning performance. In a similar trend, LoRA [22] injects trainable rank decomposition matrices into a frozen pre-trained model. Instead of inserting new parameters into pre-trained models, 82 prompt-based [26, 29] methods add trainable parameters to the input and keep the entire pre-trained 83 model unchanged during training. The inserted prompts learn to make use of the knowledge of the 84 pre-trained model to solve new tasks. Although the concept of adapter-based and prompt-based 85 approaches is different, He et al. [17] unifies the two lines of approaches (including LoRA) into 86 adapter-based methods. In addition to approaches that introduce new parameters, there are also 87 various methods [45, 52, 16] that select a sparse subset of parameters from the pre-trained model to 88 update, without adding any new parameters. One of such representative methods is BitFit [52], which 89 updates every bias term in the model. 90

PETL for CV/VL. While most of the progress in PETL is made in NLP domain, researchers have 91 also applied this technique to the CV [41, 40, 25, 53, 56, 1, 54, 20] and VL [46, 55, 56, 1, 54] domains. 92 VL-Adapter [46] benchmarks adapter-based and prompt-based methods on multiple image-text and 93 video-text tasks, and shows adapters enable us to efficiently learn fusion information of vision and 94 language. On the CV side, benefitting from the parameter efficiency of adapters and prompt-tuning, 95 some works [56, 1, 54] apply these approaches to CLIP [37] to achieve strong few-shot performance 96 in image classification tasks. Side-Tuning [53] applies an additive side network, which sums its 97 representation with the backbone network in the last layer, to solve various tasks with ResNet [18] 98 and BERT [8]. Although LST bears similarities and takes inspiration from Side-Tuning, we argue 99 that there are major differences in motivations, architecture designs, and applied tasks between the 100 two methods. LST aims to reduce the memory requirement of current PETL methods, whereas 101 Side-Tuning does not focus on memory reduction (sometimes their side network is even as big as the 102 backbone network), but instead their motivation is to ease the forgetfulness in incremental learning 103 by freezing the backbone model. Our ladder side network is more robust than their design because 104 the shortcuts fuse the intermediate information from the backbone network, and we also use layer 105 dropping and network pruning techniques to make LST more efficient and stronger. Lastly, we further 106 extend LST in VL architecture and demonstrate its usefulness on multiple VL tasks. 107

Current PETL approaches explore how to achieve competitive results using as few parameters as 108 possible. However, parameter efficiency does not necessarily mean memory efficiency. In this work, 109 we propose LST that has these two benefits simultaneously. Concurrently, Liu et al. [32] also propose 110 Y-tuning to address a similar issue; it exhausts all possible labels and feeds them into a model to 111 select the best answer from the input. However, it is intractable oftentimes to list all answers in some 112 tasks, for example, regression and open-ended generation tasks. On the other hand, LST is more 113 flexible in applying to different architectures and tasks. We show that LST can outperform Y-tuning 114 with fewer parameter updates in Table 2. 115

116 2.2 Network Pruning

In this paper, we use network pruning to extract a sub-network that contains critical information of 117 the backbone model, and use it as the initialization for our side network. Network pruning makes 118 models lighter by removing unimportant parameters or neurons. This technique aims to produce 119 pruned architectures with fast inference speed, low memory footprint, and competitive performance. 120 While PETL still uses those untrained parameters in the forward pass, network pruning entirely 121 122 discards them or sets them to zero. As a result, network pruning can generate models for faster inference speed while PETL can achieve sufficient performance by updating fewer parameters. The 123 standard procedure of network pruning is (1) learn [12, 13] or heuristically define [28] an "importance 124 measure" to identify the importance of parameters, (2) prune p% of parameters with lower importance 125 scores, (3) repeat the first and second steps until reaching the target sparsity. The rewinding procedure 126

enables pruning techniques to find a more sparse sub-network. In this paper, to keep the whole pruning process efficient, we either use weights magnitude [28] or Fisher Information [45, 30] as the importance measure, and reach the target sparsity in one shot. As a PETL method, LST makes use of the intermediate information of the backbone model as the inputs, and we empirically demonstrate those additional inputs significantly improve the performance.

132 **3 Ladder Side-Tuning (LST)**

We introduce Ladder Side-Tuning (LST), a new PETL technique that can also reduce training memory requirements by substantial amounts than previous methods. In Section 3.1, we analyze the computational cost for fine-tuning with trainable modules in backbone models. Then we explain the architectural details (Section 3.2), structural weight initialization based on network pruning (Section 3.3), and dropping side network layers for more efficiency (Section 3.4).

138 **3.1** Dependency on Backpropagation through Large Backbone Model

We consider a N multilayer perceptron (MLP): $f_N(f_{N-1}(...f_2(f_1(x))...))$, where the i^{th} layer $f_i(x) = \sigma_i(W_ix + b_i)$ consists of weight W_i , bias b_i , and nonlinear function σ_i . We denote the output of i^{th} layer as a_{i+1} and the pre-activation as z_{i+1} , where $a_{i+1} = \sigma_i(z_{i+1}) = \sigma_i(W_ia_i + b_i)$. In backpropagation with loss L, the gradient with respect to W_i and b_i :

$$\frac{\partial L}{\partial W_i} = \frac{\partial L}{\partial a_{i+1}} \frac{\partial a_{i+1}}{\partial z_{i+1}} \frac{\partial z_{i+1}}{\partial W_i} = \frac{\partial L}{\partial a_{i+1}} \sigma'_i a_i, \qquad \frac{\partial L}{\partial b_i} = \frac{\partial L}{\partial a_{i+1}} \sigma'_i \tag{1}$$

where σ'_i is the derivative of σ_i . $\frac{\partial L}{\partial a_{i+1}}$, the gradient with respect to a_i , can be calculated with the gradients with respect to a_{i+2} , using the chain rule:

$$\frac{\partial L}{\partial a_{i+1}} = \frac{\partial L}{\partial a_{i+2}} \frac{\partial a_{i+2}}{\partial z_{i+2}} \frac{\partial z_{i+2}}{\partial a_{i+1}} = \frac{\partial L}{\partial a_{i+2}} \sigma'_{i+1} W_{i+1}$$
(2)

As shown in Eqs. (1) and (2), during backpropagation, 1) $\{a\}$ corresponding to updated parameters 145 $\{W\}$ and $\{b\}$ and $\{o\}$ for the chain rule must be cached and thus dominate the memory footprint, 146 where we use $\{\cdot\}$ to denote a set of activations, parameters, or gradients. Existing PETL methods, 147 such as Adapters [21], LoRA [22], Prompt-tuning [26], and BitFit [52, 4] could reduce the memory 148 footprint by making |a| smaller, as they have fewer $\{W\}$ to update, but do not reduce $|\sigma'|$, where $|\cdot|$ 149 means the size of set $\{\cdot\}$. Since $|a| = |\sigma'|$ holds for most activation functions, the memory footprint 150 for backpropagation $|a| + |\sigma'|$ can be reduced by up to 50% by the PETL methods when they reduce 151 the entire memory footprint for |a|. By making the updated parameter do not require backpropagation 152 through the backbone network, our LST can achieve better memory efficiency beyond 50%, and we 153 explain it below. 154

155 3.2 Ladder Side Network for Transformers

Unlike existing transfer learning methods that insert additional parameters inside a transformer net-156 work, we propose training a *ladder* side network, a small and separate network that takes intermediate 157 activations from the backbone transformer as input and makes predictions. As illustrated in Fig. 3 158 (a), since the ladder side network parameters ϕ are not used during the forward pass of the backbone 159 transformer with parameters θ , the update of the ladder side network does not require expensive 160 backpropagation of the large backbone transformer. Note that our LST method is not limited to a 161 specific architecture. We provide a simplified overview of LST with an encoder architecture in Fig. 3 162 (a) and an illustration of LST with an encoder-decoder architecture in Fig. 3 (b). 163

Lightweight architecture. Our side network g is a lightweight version of the backbone transformer f, where all weights and hidden state dimensions of in g are $\frac{1}{r}$ times of the original weights and hidden states of f, where r is a reduction factor (e.g. r = 2, 4, 8, 16). For example, if the backbone f has a 768-dimensional hidden state, then the side network g with r = 16 has a hidden state of 48 dimensions (= 768/16). The side network g reuses frozen word embeddings ('Emb' in Fig. 3 (a))



Figure 3: Illustration of Ladder Side-Tuning (LST) with transformer described in Section 3.2. (a) shows a high-level overview of LST, and (b) shows LST with an encoder-decoder architecture.

and the language model head ('LM head' in Fig. 3 (a)) of the backbone f. Following the analysis in Section 3.1, we also examine the memory cost of LST. Recall that original memory footprint for backpropagation is $|a| + |\sigma'|$. Because we do not have to run a backward pass through the backbone network, we can only consider the gradients for the side network, whose memory footprint is $\frac{|a|+|\sigma'|}{r}$. Therefore, LST has a better memory efficiency than other PETL methods (saving up to 50%) as long as r is greater than 2 (we find 8 works well in most experiments).

Gated ladder connections. Although Zhang et al. [53] found that late fusion to combine the 175 representations of the backbone and the side network works well with convolutional networks for 176 CV tasks, in our experiments, we find that late fusion hurts the performance of the transformer 177 architecture in NLP tasks (see Figure 8 in Section 5 for details). To address this, we use the shortcut 178 connection (called *ladder*, due to the overall shape created from the multiple shortcut connections) 179 from intermediate activations from the backbone f to the side network q and find it helpful. We learn 180 linear projections to downsample $(\times \frac{1}{x})$ the intermediate activations (including word embeddings) of 181 f to low-dimensional attention blocks in g. Then, we learn a linear projection to upsample $(\times r)$ the 182 side network output to the dimension of the original language model head. The linear projections are 183 illustrated as green trapezoids in Fig. 3 (a). The i^{th} transformer layer of the side network g combines 184 the activation of the backbone h_i^f and the activation of the previous layer of the side network h_{i-1}^g 185 with learned gating: $\mu_i * h_i^f + (1 - \mu_i) * h_{i-1}^g$, where $\mu_i = \text{sigmoid}(\frac{\alpha_i}{T})$ is a gate parameterized with a learnable zero-initialized scalar α_i and temperature T (= 0.1). 186 187

188 3.3 Structural Weight Initialization for Ladder Side Network

We find it helpful to initialize the weights of the side network ϕ from the weight of the backbone 189 network θ based on structural pruning [28], as shown in Fig. 4 (a). Concretely, given a weight matrix $W \in \mathbb{R}^{d_{out} \times d_{in}}$ of the backbone network that maps the d_{in} -dim vectors to the d_{out} -dim space, and the importance matrix of the weight $I \in \mathbb{R}^{d_{out} \times d_{in}}$, we first calculate the importance score of 190 191 192 each row $s_i = \sum_j |I_{i,j}|$, where the importance matrix I can be weight magnitude [28] (I = W) or 193 empirical Fisher Information [45] $(I = F_W = \frac{1}{N} \sum_{i=1}^{N} (\nabla_W \log p(y_i|x_i))^2; (x_i, y_i), ..., (x_N, y_N)$ are samples from data). Then, we choose the rows of W which have the top $\frac{d_{out}}{r}$ importance scores 194 195 and prune the remaining rows to obtain a new weight matrix $W^P \in \mathbb{R}^{\frac{d_{out}}{r} \times d_{in}}$. The columns of the 196 weights and the importance matrix in the next layer corresponding to the pruned feature map are 197 also pruned. By iterating this process, we obtain the set of weight matrices whose rows and columns 198 are pruned $\frac{1}{r}$ times from the backbone network and use them to initialize the side network. In our 199



Figure 4: Illustration of (a) Structural Weight Initialization (Section 3.3) and (b) Layer Dropping (Section 3.4). In our experiments, we find that initialization of side network parameters from backbone network parameters improves performance, and dropping some shortcut connections improves efficiency without hurting performance.

experiments shown in Figure 7, we find that using Fisher information as an importance score metric generally performs well, and therefore we use it in our structural weight initialization.

202 3.4 Layer Dropping in the Ladder Side Network

We explore to increase efficiency of LST even further by making side network more compact, by dropping its intermediate transformer layers, as illustrated in Fig. 4 (b). Similar to LayerDrop [11], we drop layers in the side network, and this can linearly reduce the memory and parameter requirements of LST. For instance, a side network with N layers will only have 2^{nd} , 4^{th} , 6^{th} ... layers left, after we drop half of the layers. Refer to Section 4 for more details on applying layer dropping on an encoder-decoder architecture. In Fig. 6, we show that layer dropping can greatly boost the model's efficiency without sacrificing performance.

210 4 Experiment Setup

Datasets. We evaluate LST on NLP and VL tasks. For NLP tasks, we use the GLUE [49] benchmark, which consists of seven classification and one regression task. The benchmark evaluate models on multiple diverse tasks over linguistic acceptability (CoLA [50]), sentiment analysis (SST-2 [43]), similarity and paraphrase (MRPC [9], QQP [24], STS-B [5]) and natural language inference (MNLI [51], QNLI [39], RTE [3]). For VL tasks, we experiment with visual question answering (VQA [15], GQA [23]), visual reasoning (NLVR² [44]) and image captioning (MSCOCO [6]) tasks.

Training and Evaluation Setup. For NLP tasks, we use T5 [38], a pre-trained encoder-decoder 217 218 language model as our backbone. We use T5-base in most experiments, except that we scale up LST on T5-large and T5-3B to demonstrate its memory efficiency. The training and evaluation process 219 follows the setup used by Mahabadi et al. [35]. Since there is no local test set, we split 1k samples 220 from the training set as the new validation set and use the original validation set as the test set. For 221 datasets whose samples are less than 10k (RTE, MRPC, STS-B, CoLA), we split the validation set 222 into two equal-sized subsets and treat them as a new validation and test set. We train every approach 223 with 10 epochs on large datasets and 20 epochs on small ones (RTE, MRPC, STS-B, CoLA) for 224 complete convergence. For MNLI, we use the mismatched set as the validation set and matched set as 225 the test set. We search for learning rates over $\{3 \times 10^{-4}, 1 \times 10^{-3}, 3 \times 10^{-3}\}$ for LST and LoRA[22], 226 and we use the optimal learning rates that are used by Mahabadi et al. [35] for other methods. The 227 reduction factor used in LST is set to 8 if not additionally specified. T5-base has 12 layers each in 228 encoder and decoder, while T5-large and T5-3B have 24 layers each. In our experiments, we do not 229 drop layers in T5-base if we do not specially mention it. For T5-large and T5-3B, we drop 24 layers 230 (12 layers each in encoder and decoder) and 46 layers (23 each) of the side network to make the 231 memory usage close to our baselines. The experiments on T5 take around 12 hours to train with one 232 A6000 GPU (48GB). 233

Table 1: Comparison between multiple parameter-efficient training methods on GLUE benchmark. We use T5-base if we don't additionally specify. We report accuracy for SST-2, MNLI, QNLI and RTE. For CoLA and STS-B, we use Matthew's Correlation and Pearson-Spearman Correlation as the metrics, respectively. For MRPC and QQP, we report the average of F1 score and accuracy. Each number in the table is the average result over three seeds. For the results with [†], we report the best performance out of three seeds due to the instability of the method.

Method	Update Param. per Task (%)	Memory Usage (GB)	CoLA	SST-2	MRPC	QQP	MNLI	QNLI	RTE	STS-B	Avg.
Full fine-tuning	100	17.6	62.8	93.9	91.9	89.9	86.2	92.5	74.1	90.3	85.2
Adapters		13.0	64.4	94 2	88.9	88 9	86.4	93.1	75.1	91.1	85.3
LoRA	1.71	12.6	63.3	94.3	90.1	89.0	86.3	93.2	75.5	90.9	85.3
BitFit	0.13	10.7	61.8	94.3	91.0	88.7	85.6	93.1	67.6	90.8	84.1
Prompt-tuning	0.03	11.6	0 [†]	90.3 [†]	74.6	88.5	82.5	92.5	59.5	90.1	72.2
Ladder Side-Tuning	1.74	5.5	58.1	94.1	90.4	88.8	85.6	93.3	71.9	90.7	84.1
Ladder Side-Tuning (T5-large)	1.23	12.2	65.3	95.7	91.6	89.7	88.6	94.1	79.9	92.4	87.1
Ladder Side-Tuning (T5-3B)	0.08	22.4	66.4	96.5	92.9	89.7	90.7	95.1	80.1	93.0	88.1

For VL tasks, we experiment with CLIP-T5 [46], which is a VL architecture combining CLIP [37] 234 and T5 [38]. We always freeze the CLIP and only train the T5 for new tasks. The CLIP visual 235 representation is concatenated with the text embedding, and the combined input is fed to T5 to make 236 predictions. A visual projection layer is added between CLIP and T5 to let the visual representation 237 have the same dimension as the text embedding. To avoid updating the visual projection layer by 238 the gradients from the backbone model, we do not feed combined inputs to the backbone model, 239 but only text inputs. The combined inputs are fed to the side network, so we can achieve efficient 240 training by only using the gradients from the side network. Because the backbone network only 241 uses texts as the input, the information from the backbone network via shortcut connections is 242 only summed to the text part of the side network's combined inputs. We follow the multi-tasking 243 setting for training and evaluation used in VL-Adapter [46]. We report the performance on Karpathy 244 test/test-dev/test-P/Karpathy test split for VQA/GQA/NLVR²/MSCOCO, and train models for 20 245 epochs. We search learning rates over $\{3 \times 10^{-4}, 1 \times 10^{-3}, 3 \times 10^{-3}\}$ for PETL methods, and 246 use 1×10^{-4} used by Sung et al. [46] for full fine-tuning. We set the reduction factor for the side 247 network to 4. We train CLIP-T5 for 16 hours on one A6000 GPU. 248

249 5 Experimental Results

In this section, we show experiments to justify our design of LST and demonstrate that LST performs the best among all approaches in the scenario with limited memory. As the result, LST is the most efficient tool to fine-tune large-scale pre-trained models for real-world applications.

LST outperforms other methods under similar memory usage. Figure 1 and Table 1 show the 253 results on GLUE of different approaches applying on T5-base. We drop 6 layers (3 layers each in side 254 encoder and decoder) for LST to match the parameter usage of the Adapter and LoRA. Under the 255 same parameter usage, LST can save 69% of memory cost to fully fine-tune the model, while Adapter 256 and LoRA only save 26% of that, leading to LST having a 2.7x more memory saving. Compared 257 to BitFit, LST achieves the same average performance but costs 5GB less GPU memory. LST also 258 surpasses Prompt-tuning in terms of both performance and memory cost. To further take advantage of 259 the memory efficiency of LST, we also train T5-large and T5-3B with LST. We find that with a similar 260 budget of memory usage in Adapter and LoRA, LST with T5-large can surpass the performance 261 of other methods by a large margin. The result on T5-3B also outperforms the result on T5-large, 262 demonstrating the scalability of our memory-efficiency method on large language models. 263

In Table 2, we also compare LST with a concurrent work, Y-tuning [32] on GLUE tasks (except for STS-B) with BARTlarge [27] encoder as backbone. Following the experimental setup of Liu et al. [32], for each task, we use different learning rate, and report the best accuracy out of 3 seeds. Overall, LST outperforms Y-tuning by a large margin with fewer updated parameters. Table 2: LST vs. Y-tuning.

Method	Update Param. per Task (%)	Avg. GLUE
Y-tuning	7.7	76.9
LST	2.6	82.1

Table 3: Comparison between multiple parameter-efficient training methods on VQA, GQA, NLVR², and MSCOCO. We use T5-base for all approaches. We report accuracy for VQA, GQA and NLVR while we use CIDEr to evaluate MSCOCO.

Method	Update Param. (%)	Memory Usage (GB)	VQA	GQA	NLVR ²	MSCOCO	Avg.
Full fine-tuning	100	36.2	67.0	56.1	73.8	111.8	77.2
Adapters	7.98	28.4	67.2	56.4	73.1	111.9	77.2
LoRA	7.54	27.9	63.9	53.3	70.2	110.3	74.4
Ladder Side-Tuning	7.46	15.3	67.3	56.8	74.1	108.5	76.7



Figure 5: The accuracy-memory trade-off for Adapter, LoRA, and Ladder Side-Tuning over GLUE tasks. We vary the reduction factor in Ladder Side-Tuning, hidden dimension in Adapter, rank in LoRA to get the architectures with different training costs.



Figure 6: The accuracy-memory trade-off for Adapter, LoRA, BitFit, LST over GLUE tasks. we drop $N \in \{0, 6, 12, 18\}$ layers in an interleaving manner for LST while we gradually freeze the first $N \in \{0, 6, 12, 18\}$ layers in other methods (also remove inserted parameters in those layers).

LST is competitive on VL tasks. As we have mentioned beforehand, we also extend LST on a multi-modal architecture, CLIP-T5, on multiple VL tasks, and we demonstrate the outcome in

Table 3. With similar parameter usage, LST is the only method that can fit into a single 16GB GPU.

Besides the efficiency, it is as competitive as full fine-tuning and Adapter and outperforms LoRA.

LST performs the best in low-memory regime. To have a better understand of the memory 275 advantage of LST, we adjust the hyper-parameters in our method (reduction factor $\in \{32, 16, 8, 4\}$), 276 Adapter (hidden dimension $\in \{6, 12, 24, 48\}$) and LoRA (rank $\in \{4, 8, 16, 32\}$) to create multiple 277 architectures with different memory costs. Figure 5 shows the performance and memory efficiency 278 trade-off for all methods. We find the memory saving is not obvious for Adapter and LoRA, because 279 the gradients of the backbone model's intermediate outputs are still computed (see Section 3.1 for 280 281 details). Even though Adapter and LoRA still can get better memory efficiency by reducing the hidden dimension and the rank, we find that the performance drops significantly. On the other hand, 282 LST is quite robust across a wide range of side network sizes. 283

We also consider another way to compare LoRA, Adapter and BitFit to LST in different memory 284 budgets. While we drop $N \in \{0, 6, 12, 18\}$ layers in an interleaving manner to improve memory 285 efficiency, we freeze the first N layers and remove the corresponding inserted modules in other 286 approaches. With this, other methods can achieve better memory efficiency because gradients do 287 not propagate to those earlier frozen layers. We discuss the layer dropping and layer freezing with 288 details in the following. In LST, we drop $\frac{N}{2}$ layers in both side encoder and side decoder. However, 289 in other PETL approaches, we start from freezing layers in the side encoder and then turn to freeze 290 layers in the side decoder (e.g. dropping 18 layers means dropping all side encoder layers and 6 side 291 decoder layers). We display the comparison in Figure 6, showing that LST has a better performance 292 and memory trade-off and outperforms other methods in the low-memory regime. We also find that 293 layer dropping generally reduces the training cost without hurting performance. 294



Figure 7: The ablation study on different initialization strategies. Y-axis denotes the average score over GLUE tasks.



Figure 8: The comparison between network pruning, Side-Tuning, and LST on GLUE tasks.

The ablation of weight initialization on the side network. We compare the different initialization strategies for the side network and demonstrate the results in Figure 7. "Random" denotes we randomly initialize the network while we use network pruning to select initial weights for the side network based on two importance measures, "Weight Magnitude" and "Fisher Information." In general, the initialization from the pruned network helps no matter the size of the side network, showing the effectiveness of our network pruning strategy.

Comparison LST to network pruning and side-tuning. In Section 3.2, we mention that shortcut 301 connections are added to every layer of the side network. We justify this design by comparing LST to 302 two approaches: (1) network pruning, which discards all shortcut connections and the entire backbone 303 model. (2) side-tuning, which only adds one shortcut connection to merge representations right 304 before the output layer. Note that we do not drop any layer in the side network but only remove the 305 shortcuts. Figure 8 shows the comparison and LST outperforms the other two methods significantly, 306 307 suggesting the intermediate shortcut connections are useful. This also suggests that network pruning might not perform as well as PETL methods when training the same number of parameters. PETL 308 methods are stronger as they use the information from the backbone model. 309

310 6 Conclusion

We propose Ladder Side-Tuning (LST), a parameter- and memory-efficient method to adapt a large 311 backbone network. LST does not require backpropagation through the backbone network, which 312 allows for significantly lower memory requirement during training than recently proposed parameter-313 314 efficient training techniques. We demonstrate that LST allows users to adapt a larger and more 315 powerful backbone network to target tasks with a limited memory requirement, which cannot be achieved with recent parameter-efficient techniques. We also show that LST achieves more efficient 316 accuracy-memory trade-off than recent baselines and impact of weight initialization of side networks. 317 Finally, we show that the LST can be also extended beyond NLP tasks, with strong results on 318 vision-and-language tasks. We hope that LST helps users with limited computational resources tune 319 larger models in diverse domains. 320

321 References

- [1] Anonymous. How to adapt your large-scale vision-and-language model. In Submitted to The Tenth
 International Conference on Learning Representations, 2022. URL https://openreview.net/forum?
 id=EhwEUb2ynIa. under review.
- [2] Hangbo Bao, Li Dong, and Furu Wei. Beit: BERT pre-training of image transformers. CoRR,
 abs/2106.08254, 2021. URL https://arxiv.org/abs/2106.08254.
- [3] Luisa Bentivogli, Ido Dagan, Hoa Trang Dang, Danilo Giampiccolo, and Bernardo Magnini. The fifth pascal recognizing textual entailment challenge. In *In Proc Text Analysis Conference (TAC'09*, 2009.
- [4] Han Cai, Chuang Gan, Ligeng Zhu, and Song Han. Tinytl: Reduce memory, not parameters for efficient on-device learning. In H. Larochelle, M. Ranzato, R. Hadsell, M.F. Balcan, and

- H. Lin, editors, Advances in Neural Information Processing Systems, volume 33, pages 11285-331 11297. Curran Associates, Inc., 2020. URL https://proceedings.neurips.cc/paper/2020/file/ 332 81f7acabd411274fcf65ce2070ed568a-Paper.pdf. 333
- 334 [5] Daniel Cer, Mona Diab, Eneko Agirre, Inigo Lopez-Gazpio, and Lucia Specia. Semeval-2017 task 1: Semantic textual similarity-multilingual and cross-lingual focused evaluation. In SemEval, 2017.

335

- [6] Xinlei Chen, Hao Fang, Tsung-Yi Lin, Ramakrishna Vedantam, Saurabh Gupta, Piotr Dollár, and 336 C. Lawrence Zitnick. Microsoft COCO captions: Data collection and evaluation server. CoRR, 337 abs/1504.00325, 2015. URL http://arxiv.org/abs/1504.00325. 338
- [7] Jaemin Cho, Jie Lei, Hao Tan, and Mohit Bansal. Unifying vision-and-language tasks via text generation. 339 In Proceedings of the 38th International Conference on Machine Learning, volume 139, pages 1931–1942. 340 PMLR, 18-24 Jul 2021. URL https://proceedings.mlr.press/v139/cho21a.html. 341
- [8] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. Bert: Pre-training of deep bidirec-342 tional transformers for language understanding. In NAACL, 2019. 343
- [9] Bill Dolan and Chris Brockett. Automatically constructing a corpus of sentential paraphrases. In 344 Third International Workshop on Paraphrasing (IWP2005). Asia Federation of Natural Language 345 Processing, January 2005. URL https://www.microsoft.com/en-us/research/publication/ 346 automatically-constructing-a-corpus-of-sentential-paraphrases/. 347
- [10] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas 348 Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, 349 and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. In 350 International Conference on Learning Representations, 2021. URL https://openreview.net/forum? 351 id=YicbFdNTTy. 352
- [11] Angela Fan, Edouard Grave, and Armand Joulin. Reducing transformer depth on demand with structured 353 354 dropout. In International Conference on Learning Representations, 2020. URL https://openreview. net/forum?id=Syl02yStDr. 355
- 356 [12] Jonathan Frankle and Michael Carbin. The lottery ticket hypothesis: Finding sparse, trainable neural networks. ICLR, 2019. 357
- [13] Jonathan Frankle, Gintare Karolina Dziugaite, Daniel M. Roy, and Michael Carbin. Linear mode connec-358 tivity and the lottery ticket hypothesis. ICML, 2020. 359
- [14] Peng Gao, Shijie Geng, Renrui Zhang, Teli Ma, Rongyao Fang, Yongfeng Zhang, Hongsheng Li, and 360 Yu Jiao Qiao. Clip-adapter: Better vision-language models with feature adapters. ArXiv, abs/2110.04544, 361 2021. 362
- [15] Yash Goyal, Tejas Khot, Douglas Summers-Stay, Dhruv Batra, and Devi Parikh. Making the v in vqa 363 matter: Elevating the role of image understanding in visual question answering. 2017 IEEE Conference on 364 Computer Vision and Pattern Recognition (CVPR), pages 6325-6334, 2017. 365
- [16] Demi Guo, Alexander Rush, and Yoon Kim. Parameter-efficient transfer learning with diff pruning. 366 In Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 367 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 368 4884–4896, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021. 369 acl-long.378. URL https://aclanthology.org/2021.acl-long.378. 370
- [17] Junxian He, Chunting Zhou, Xuezhe Ma, Taylor Berg-Kirkpatrick, and Graham Neubig. Towards a unified 371 view of parameter-efficient transfer learning. In International Conference on Learning Representations, 372 2022. URL https://openreview.net/forum?id=ORDcd5Axok. 373
- [18] Kaiming He, X. Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. 2016 374 IEEE Conference on Computer Vision and Pattern Recognition (CVPR), pages 770–778, 2016. 375
- [19] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross B. Girshick. Masked 376 autoencoders are scalable vision learners. CoRR, abs/2111.06377, 2021. URL https://arxiv.org/ 377 abs/2111.06377. 378
- 379 [20] Xuehai He, Chengkun Li, Pengchuan Zhang, Jianwei Yang, and Xin Wang. Parameter-efficient fine-tuning for vision transformers. ArXiv, abs/2203.16329, 2022. 380

- [21] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Ges mundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *International Conference on Machine Learning*, pages 2790–2799. PMLR, 2019.
- [22] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yuanzhi Li, Shean Wang, and Weizhu
 Chen. Lora: Low-rank adaptation of large language models. *CoRR*, abs/2106.09685, 2021. URL
 https://arxiv.org/abs/2106.09685.
- [23] Drew A Hudson and Christopher D Manning. Gqa: A new dataset for real-world visual reasoning and
 compositional question answering. *Conference on Computer Vision and Pattern Recognition (CVPR)*,
 2019.
- [24] Shankar Iyer, Nikhil Dandekar, and Kornel Csernai. First quora dataset release: Question pairs. 2017.
 URL https://data.quora.com/First-Quora-Dataset-Release-Question-Pairs.
- [25] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge J. Belongie, Bharath Hariharan, and
 Ser Nam Lim. Visual prompt tuning. *ArXiv*, abs/2203.12119, 2022.
- Brian Lester, Rami Al-Rfou, and Noah Constant. The power of scale for parameter-efficient prompt tuning.
 In *EMNLP*, 2021.
- [27] Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy,
 Veselin Stoyanov, and Luke Zettlemoyer. Bart: Denoising sequence-to-sequence pre-training for natural
 language generation, translation, and comprehension. In *ACL*, 2020.
- Hao Li, Asim Kadav, Igor Durdanovic, Hanan Samet, and Hans Peter Graf. Pruning filters for efficient
 convnets. *ICLR*, 2017.
- [29] Xiang Lisa Li and Percy Liang. Prefix-tuning: Optimizing continuous prompts for generation. In
 Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers), pages 4582–
 4597, Online, August 2021. Association for Computational Linguistics. doi: 10.18653/v1/2021.acl-long.
 353. URL https://aclanthology.org/2021.acl-long.353.
- [30] Liyang Liu, Shilong Zhang, Zhanghui Kuang, Aojun Zhou, Jingliang Xue, Xinjiang Wang, Yimin Chen,
 Wenming Yang, Qingmin Liao, and Wayne Zhang. Group fisher pruning for practical network compression.
 In *ICML*, 2021.
- [31] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis,
 Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized BERT pretraining approach.
 CoRR, abs/1907.11692, 2019. URL http://arxiv.org/abs/1907.11692.
- 412 [32] Yitao Liu, Chen An, and Xipeng Qiu. Y-tuning: An efficient tuning paradigm for large-scale pre-trained 413 models via label representation learning. *ArXiv*, abs/2202.09817, 2022.
- [33] Jiasen Lu, Dhruv Batra, Devi Parikh, and Stefan Lee. Vilbert: Pretraining task-agnostic visiolinguistic
 representations for vision-and-language tasks. In *NeurIPS*, 2019.
- [34] Rabeeh Mahabadi, Sebastian Ruder, Mostafa Dehghani, and James Henderson. Parameter-efficient multi-task fine-tuning for transformers via shared hypernetworks. In *Annual Meeting of the Association for Computational Linguistics*, 2021.
- [35] Rabeeh Karimi Mahabadi, James Henderson, and Sebastian Ruder. Compacter: Efficient low-rank
 hypercomplex adapter layers. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021.
 URL https://openreview.net/forum?id=bqGK5PyI6-N.
- Yuning Mao, Lambert Mathias, Rui Hou, Amjad Almahairi, Hao Ma, Jiawei Han, Wen-tau Yih, and
 Madian Khabsa. Unipelt: A unified framework for parameter-efficient language model tuning. *CoRR*,
 abs/2110.07577, 2021. URL https://arxiv.org/abs/2110.07577.
- [37] Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya Ramesh, Gabriel Goh, Sandhini Agarwal, Girish
 Sastry, Amanda Askell, Pamela Mishkin, Jack Clark, Gretchen Krueger, and Ilya Sutskever. Learning
 transferable visual models from natural language supervision. In *ICML*, 2021.
- [38] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou,
 Wei Li, and Peter J. Liu. Exploring the limits of transfer learning with a unified text-to-text transformer.
 Journal of Machine Learning Research, 21(140):1–67, 2020. URL http://jmlr.org/papers/v21/
- 431 20-074.html.

- [39] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for
 machine comprehension of text. In *EMNLP*, 2016.
- 434 [40] Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. Learning multiple visual domains with 435 residual adapters. In *NIPS*, 2017.
- [41] Sylvestre-Alvise Rebuffi, Hakan Bilen, and Andrea Vedaldi. Efficient parametrization of multi-domain
 deep neural networks. 2018 IEEE/CVF Conference on Computer Vision and Pattern Recognition, pages
 8119–8127, 2018.
- [42] Andrei A. Rusu, Neil C. Rabinowitz, Guillaume Desjardins, Hubert Soyer, James Kirkpatrick, Koray
 Kavukcuoglu, Razvan Pascanu, and Raia Hadsell. Progressive neural networks. *ArXiv*, abs/1606.04671,
 2016.
- [43] Richard Socher, Alex Perelygin, Jean Wu, Jason Chuang, Christopher D. Manning, Andrew Ng, and Christopher Potts. Recursive deep models for semantic compositionality over a sentiment treebank. In *Proceedings of the 2013 Conference on Empirical Methods in Natural Language Processing*, pages 1631–1642, Seattle, Washington, USA, October 2013. Association for Computational Linguistics. URL https://aclanthology.org/D13-1170.
- [44] Alane Suhr, Stephanie Zhou, Ally Zhang, Iris Zhang, Huajun Bai, and Yoav Artzi. A corpus for reasoning about natural language grounded in photographs. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6418–6428, Florence, Italy, July 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1644. URL https://aclanthology.org/P19-1644.
- [45] Yi-Lin Sung, Varun Nair, and Colin Raffel. Training neural networks with fixed sparse masks. In
 Advances in Neural Information Processing Systems, 2021. URL https://openreview.net/forum?
 id=Uwh-v1HSw-x.
- 454 [46] Yi-Lin Sung, Jaemin Cho, and Mohit Bansal. Vl-adapter: Parameter-efficient transfer learning for vision-455 and-language tasks. *CVPR*, 2022.
- [47] Hao Tan and Mohit Bansal. Lxmert: Learning cross-modality encoder representations from transformers.
 In *EMNLP*, 2019.
- [48] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Ł ukasz
 Kaiser, and Illia Polosukhin. Attention is all you need. In I. Guyon, U. Von Luxburg, S. Bengio, H. Wallach,
 R. Fergus, S. Vishwanathan, and R. Garnett, editors, *Advances in Neural Information Processing Systems*,
 volume 30. Curran Associates, Inc., 2017. URL https://proceedings.neurips.cc/paper/2017/
 file/3f5ee243547dee91fbd053c1c4a845aa-Paper.pdf.
- [49] Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. GLUE: A
 multi-task benchmark and analysis platform for natural language understanding. In *International Confer- ence on Learning Representations*, 2019. URL https://openreview.net/forum?id=rJ4km2R5t7.
- [50] Alex Warstadt, Amanpreet Singh, and Samuel R Bowman. Neural network acceptability judgments. *arXiv preprint arXiv:1805.12471*, 2018.
- 468 [51] Adina Williams, Nikita Nangia, and Samuel R Bowman. A broad-coverage challenge corpus for sentence
 469 understanding through inference. In *NAACL*, 2018.
- [52] Elad Ben Zaken, Shauli Ravfogel, and Yoav Goldberg. Bitfit: Simple parameter-efficient fine-tuning for
 transformer-based masked language-models. *CoRR*, abs/2106.10199, 2021. URL https://arxiv.org/
 abs/2106.10199.
- [53] Jeffrey O. Zhang, Alexander Sax, Amir Roshan Zamir, Leonidas J. Guibas, and Jitendra Malik. Side-tuning:
 Network adaptation via additive side networks. *ECCV*, 2020.
- [54] Renrui Zhang, Rongyao Fang, Peng Gao, Wei Zhang, Kunchang Li, Jifeng Dai, Yu Qiao, and Hongsheng
 Li. Tip-adapter: Training-free clip-adapter for better vision-language modeling. *ArXiv*, abs/2111.03930,
 2021.
- [55] Zhengkun Zhang, Wenya Guo, Xiaojun Meng, Yasheng Wang, Yadao Wang, Xin Jiang, Qun Liu, and
 Zhenglu Yang. Hyperpelt: Unified parameter-efficient language model tuning for both language and
 vision-and-language tasks. *ArXiv*, abs/2203.03878, 2022.
- [56] Kaiyang Zhou, Jingkang Yang, Chen Change Loy, and Ziwei Liu. Learning to prompt for vision-language
 models. *ArXiv*, abs/2109.01134, 2021.

483 Checklist

484	1. For all authors
485 486	(a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes]
487	(b) Did you describe the limitations of your work? [Yes] In appendix.
488 489	(c) Did you discuss any potential negative societal impacts of your work? [N/A] We believe we do not have negative societal impacts due to it is a technical paper.
490 491	(d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
492	2. If you are including theoretical results
493 494	(a) Did you state the full set of assumptions of all theoretical results? [N/A] We do not have any theoretical results.
495	(b) Did you include complete proofs of all theoretical results? [N/A]
496	3. If you ran experiments
497	(a) Did you include the code, data, and instructions needed to reproduce the main experi-
498	mental results (either in the supplemental material or as a URL)? [Yes] See supplemen-
499	tary materials.
500 501	(b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] See Section 4.
502 503	(c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [Yes] See appendix.
504 505	(d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes]
506	4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets
507	(a) If your work uses existing assets, did you cite the creators? [Yes]
508	(b) Did you mention the license of the assets? [Yes]
509 510	(c) Did you include any new assets either in the supplemental material or as a URL? [Yes] Our codes.
511	(d) Did you discuss whether and how consent was obtained from people whose data you're
512	using/curating? [N/A]
513 514	(e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [N/A]
515	5. If you used crowdsourcing or conducted research with human subjects
516 517	(a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
518 519	(b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
520 521	(c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]