

1 A. Additional Technical Details

2 A.1. Summary of Notations

3 Tab. 1 summarizes the major notations.

4 A.2. Estimation of $M_w(s)$

5 This section provides the details for estimating $M_w(s)$ in §
6 4.3.1. $M_w(s)$ captures the memory overhead introduced by
7 all LoRA modules selected under strategy s , including their
8 parameters, gradients, and optimizer states. Let $M_{param}(s)$
9 be the total parameter memory of the inserted modules.
10 Gradients require the same amount of memory as parameters.
11 The optimizer state size depends on the optimizer.
12 For example, AdamW maintains two moment estimates per
13 parameter, resulting in an additional $2 \times M_{param}(s)$ mem-
14 ory cost [31]. Therefore, for AdamW, the total module-
15 related memory is $M_w(s) = 4 \cdot M_{param}(s)$, where the
16 factor 4 accounts for parameters, gradients, and two sets
17 of moment estimates. In general, $M_w(s)$ scales linearly
18 with $M_{param}(s)$ and can be computed directly once the op-
19 timizer is specified.

20 A.3. Complete Intra-Block Dependencies for $g_x(k)$

21 This section explains how to configure $g_x(k)$ for the remain-
22 ing five intra-block dependencies in § 4.3.2.

23 **Inserting a module w_k at W_k (Fig. 5a).** The forward
24 pass is $K = x_{in}(W_k + w_k)$. To compute $\nabla_{w_k} L$, gradi-
25 ents flow from $\nabla_{\text{Attn}(Q,K,V)} L$ to $\nabla_K L$ via the attention
26 operator. Since $\text{Attn}(Q, K, V) = \text{Softmax}(\frac{QK^T}{\sqrt{d}})V$, the
27 dependency between K and the output involves the inter-
28 action with Q and V . Therefore, computing $\nabla_K L$ requires
29 the stored activations Q and V . Finally, $\nabla_{w_k} L$ needs the
30 input activation x_{in} . Thus, insertion location at W_k requires
31 storing Q , V , and x_{in} .

32 **Inserting a module w_v at W_v (Fig. 5b).** The forward
33 pass is $V = x_{in}(W_v + w_v)$. To compute $\nabla_{w_v} L$, gradi-
34 ents flow to $\nabla_V L$ through the attention operator. In the
35 term $\text{Attn}(Q, K, V) = \text{Softmax}(\frac{QK^T}{\sqrt{d}})V$, the variable V
36 is multiplied by the attention weights (the result of Soft-
37 max). Consequently, the backward pass to V requires the
38 reconstruction of these attention weights, which depends on
39 the stored activations Q and K . Calculating $\nabla_{w_v} L$ further
40 requires the input x_{in} . Thus, insertion location at W_v re-
41 quires storing Q , K , and x_{in} .

42 **Inserting a module w_o at W_o (Fig. 5c).** The forward pass
43 is $x_{out}^1 = (W_o + w_o) \cdot \text{Attn}(Q, K, V) + x_{in}$. To compute
44 $\nabla_{w_o} L$, gradients flow directly from $\nabla_{x_{out}^1} L$. The local gra-
45 dient computation for the weight w_o requires the input to
46 the linear projection, which is the output of the attention
47 mechanism, $\text{Attn}(Q, K, V)$. Thus, insertion location at W_o
48 requires storing the activation vector $\text{Attn}(Q, K, V)$.

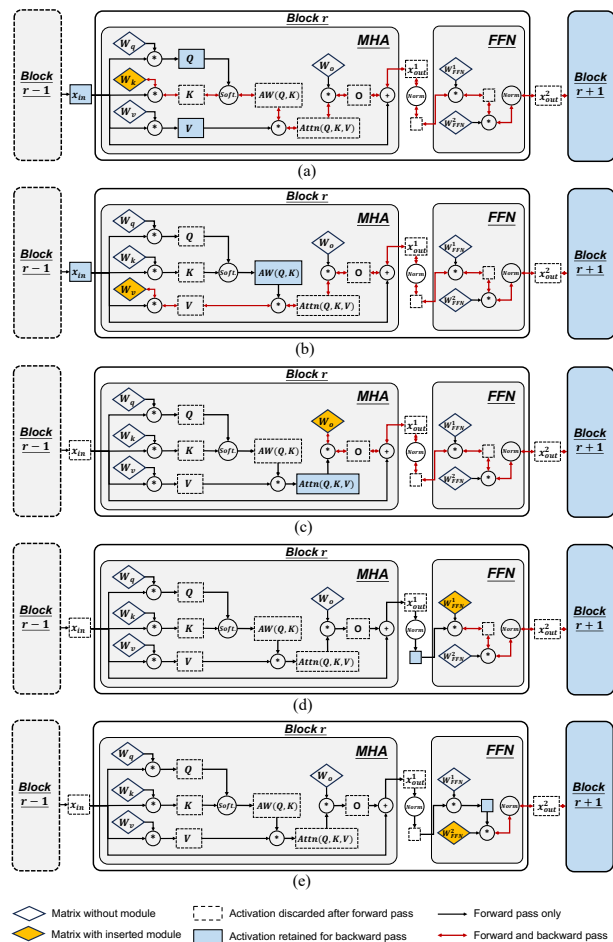


Figure 1. Intra-block activation dependency analysis. (a)-(e) demonstrate the activation preservations when inserting modules into W_k , W_v , W_o , W_{FFN}^1 , and W_{FFN}^2 matrices, respectively.

49 **Inserting a module w_{FFN}^1 at W_{FFN}^1 (Fig. 5d).** Let the
50 FFN be defined as $FFN(x) = W_{FFN}^2 \cdot \psi((W_{FFN}^1 + w_{FFN}^1) \cdot$
51 $\text{Norm}(x))$, where ψ is the activation function like ReLU. To
52 compute $\nabla_{w_{FFN}^1} L$, the backpropagation process requires the
53 input to this specific layer. The input to W_{FFN}^1 is the normal-
54 ized output of the attention block, denoted as $\text{Norm}(x_{out}^1)$.
55 Thus, insertion at W_{FFN}^1 requires storing $\text{Norm}(x_{out}^1)$.

56 **Inserting a module w_{FFN}^2 at W_{FFN}^2 (Fig. 5e).** The forward
57 pass involves $x_{out}^2 = (W_{FFN}^2 + w_{FFN}^2) \cdot v + x_{out}^1$, where $v =$
58 $\psi(W_{FFN}^1 \cdot \text{Norm}(x_{out}^1))$ is the intermediate activation of the
59 feed-forward network. To compute $\nabla_{w_{FFN}^2} L$, the gradient
60 computation requires the input to the second linear layer,
61 which is the intermediate hidden state v . Thus, insertion at
62 W_{FFN}^2 requires storing the intermediate FFN activation v .

A.4. Pseudocode for DP-based Scheduler

Algorithm 1 illustrates the workflow of the module inser-
tion scheduler. The DP state space is $O(|\mathcal{K}| \cdot U)$. Com-

Notation	Definition
t_1, \dots, t_n	Previously learned tasks
t_{n+1}	New task to be fine-tuned
d	Batched input data of the new task
θ	Frozen pre-trained backbone parameters
\mathcal{K}	Set of projection matrices, i.e. candidate LoRA module insertion locations
W_k	Frozen projection matrix at backbone location $k \in \mathcal{K}$
w_i^k	LoRA module parameters for task t_i at location k
\mathbf{s}	Module insertion strategy (vector where s^k is the rank at location k)
M_B	Total memory budget for LoRA fine-tuning
λ_{n+1}	Predicted module importance vector for new task t_{n+1}
λ_{n+1}^k	Importance of module at location k for new task t_{n+1}
λ_i^k	Importance of module at location k for learned task t_i
$\phi(t_{n+1}, t_i)$	Task similarity between new task t_{n+1} and learned task t_i
$x(w_i^k; d)$	Output activations of LoRA module w_i^k on data d
$x(W_k; d)$	Output activations of frozen matrix W_k on data d
$x(\theta + \mathbf{w}_i; d)$	Output activations of the backbone θ merged with task modules \mathbf{w}_i
$x^*(\theta + \mathbf{w}_{n+1}; d)$	Approximated activation of the new task
$M(\mathbf{s})$	Peak fine-tuning memory usage under strategy \mathbf{s}
M_{con}	Constant memory usage (frozen backbone parameters and runtime overhead)
$M_w(\mathbf{s})$	Memory usage of LoRA modules (parameters, gradients, optimizer states)
$M_x(\mathbf{s})$	Memory usage of cached activations during backpropagation
$\mathbf{g}_x(\mathbf{s})$	Binary mask indicating activations to be retained for strategy \mathbf{s}
\mathbf{m}_x	Memory size vector of candidate activations
$W_q, W_k, W_v, W_o, W_{FFN}^1, W_{FFN}^2$	Frozen projection matrices in backbone
r	Rank size of an inserted module
Q, K, V	Query, Key, and Value matrices in attention blocks
$P[k][M_b]$	Maximum cumulative importance at location k with memory cost $\leq M_b$
k_s	Index of the nearest previously inserted module location (predecessor)
ΔM	Incremental memory cost of inserting a module given predecessor k_s
\mathcal{R}	Set of candidate LoRA ranks
$h(r)$	Rank-based importance scaling function
U	Number of memory discretization bins

Table 1. Summary of major notations.

66 putting each state $P[k][M_b]$ involves iterating over all ranks
67 r ($O(|\mathcal{R}|)$) and preceding locations k_s ($O(|\mathcal{K}|)$). Thus, the
68 overall time complexity is $O(|\mathcal{K}|^2 \cdot |\mathcal{R}| \cdot U)$.

69 B. Additional Experimental Setups

70 B.1. Details of Baselines

71 The configurations of the baselines are as follows.

72 B.1.1. Zero-FT

- 73 • Prompt Tuning [14] adds prompt vectors to the keys and
- 74 values of the attention module in each transformer block.
- 75 We adopt a deep prompt tuning approach, incorporating
- 76 16 learnable prompt vectors per block into the first 24
- 77 blocks of the backbone following [20].
- 78 • AdapterSoup [3] averages the weight space of pre-trained
- 79 domain adapters, selected via textual clustering or simi-
- 80 larity, to enhance generalization to new domains.
- 81 • AdapterFusion [22] introduces a fusion layer with an at-

tention mechanism to combine multiple pre-trained task
adapters. We implement the fusion layer with a reduction
factor of 16, balancing parameters and performance.
• LoraHub [11] applies a gradient-free evolutionary strat-
egy to learn scalar weights for combining LoRA modules.
We set the LoRA rank to $r = 8$ and randomly select 20
candidate modules for composition.

B.1.2. Non-LoRA-based PEFT

- Full tuning [12] updates all parameters of the pre-trained
backbone and classification head.
- Linear tuning [12] freezes the pre-trained backbone and
only updates the parameters of the final linear classifica-
tion head.
- BitFit [28] optimizes the bias terms within all linear lay-
ers of each transformer block. We apply this optimization
to every bias term throughout the network, following the
approach outlined in the original implementation.
- RS-Bypass [13] creates an independent bypass network

Algorithm 1 DP-based Module Insertion Scheduler

Require: Memory budget M_B , Module importance scores λ_{n+1}^k , Memory prediction function $M(s)$, Rank set $\mathcal{R} = \{0, 2^0, \dots, 2^6\}$.

```
1: Initialize DP table  $P[k][M_b] \leftarrow -\infty$  for all  $k$  and  $M_b$ 
2: Set  $P[0][M_b] \leftarrow 0$  for all  $M_b$ 
3: Set  $U \leftarrow 500$ 
4: Set  $\mathcal{R} \leftarrow [0, 1, 2, 4, 8, 16, 32, 64]$ 
5: Set the step size for memory bins  $\Delta M_b = \frac{M_B}{U}$ 
6: for  $k = 1$  to  $|\mathcal{K}|$  do
7:   for  $M_b = 0$  to  $M_B$  with step  $\Delta M_b$  do
8:     Initialize  $P[k][M_b] \leftarrow P[k-1][M_b]$  {Start with
       the value from the previous state (no insertion)}
9:     for  $r \in \mathcal{R}, r > 0$  do
10:      for  $k_s = 0$  to  $k-1$  do
11:         $\Delta M \leftarrow M(s_k) - M(s_{k_s})$ 
12:        if  $P[k_s][M_b - \Delta M] > -\infty$  then
13:           $imp \leftarrow P[k_s][M_b - \Delta M] + h(r)\lambda_{n+1}^k$ 
           {Calculate the new importance if inserting
           rank  $r$ }
14:          if  $imp > P[k][M_b]$  then
15:             $P[k][M_b] \leftarrow imp$ 
16:          end if
17:        end if
18:      end for
19:    end for
20:  end for
21: end for
22:  $s \leftarrow \max_{0 \leq k \leq |\mathcal{K}|} P[k][M_B]$ 
23: Backtrack to find optimal module insertion strategy
24: return  $s$ 
```

100 using detached Res-Tuners to avoid backpropagation
101 through the pre-trained backbone, reducing memory use-
102 age. We use a learnable fusion parameter initialized at
103 0.5, with other settings consistent with RS.
104 • Adapter [9] inserts adapters between transformer layers,
105 consisting of a down-projection, followed by a nonlinear
106 activation, and an up-projection. We set the bottleneck
107 dimension $r = 128$.
108 • AdaptFormer [2] introduces a lightweight bottleneck
109 module (AdaptMLP) alongside the frozen MLP block,
110 consisting of down-projection and up-projection layers.
111 We set the bottleneck middle dimension to $\hat{d} = 64$ and
112 the scaling factor to $s = 0.1$.
113 • FacT-TK [15] tensorizes transformer weights into a
114 3D tensor and decomposes the weight increments into
115 lightweight factors for efficient tuning. We typically use
116 a rank of $r = 8$ for the decomposition factors.
117 • RS [13] separates parameter-efficient tuners from the
118 backbone, treating them as residual modules added to
119 frozen operations, allowing flexible combinations of

adapters, prefixes, and prompts. We set the prompt/prefix
length and adapter middle dimension to 10. 120 121

- ConvPass [14] inserts trainable convolutional bypass
modules (1×1 , 3×3 , 1×1) in parallel to the atten-
tion and MLP blocks to introduce visual inductive bias. 122 123 124
We set the bottleneck dimension to $h = 8$ for the convo-
lutional layers. 125 126
- LISA [21] uses layerwise importance sampling to freeze
most layers during optimization, updating only the em-
bedding, head, and a few random intermediate layers. We 127 128
configure the strategy to update the embedding layer, the
head, and 2 random intermediate layers (E+H+2L). 129 130 131
- LoSA [20] adapts transformers by adding low-rank mixer
blocks in parallel to the frozen backbone to efficiently
process features. We set the low-rank adaptation rank to
 $r = 8$. 132 133 134 135
- HST [18] attaches a lightweight Hierarchical Side Net-
work (HSN) to the frozen backbone to extract multi-
scale features and uses a Meta-Register token to aggre-
gate global information. We use a single trainable token
for the Meta-Register and align the HSN stages with the
backbone. 136 137 138 139 140 141
- UniPT [5] supports transfer learning via a parallel net-
work with an interaction module for processing interme-
diate activations and a confidence aggregation module for
cross-layer feature integration. We adopt a reduction fac-
tor of 4 for the parallel network dimensions. 142 143 144 145 146

B.1.3. LoRA-based PEFT 147

- LoRA [10] decomposes the weight matrices into trainable
low-rank parameters to efficiently approximate weight
updates. In our approach, we apply this adaptation to
all linear projections in the transformer block (includ-
ing query, key, value, output, and MLP), using a rank of
 $r = 8$. 148 149 150 151 152 153
- NOAH [32] treats Adapter, LoRA, and VPT as modular
prompts and uses a one-shot neural architecture search to
find the optimal module design and hyperparameters for
each transformer block. We train the supernet for 500
epochs and perform the evolutionary search for 5 epochs. 154 155 156 157 158
- π -Adapter [26] uses Fisher Information Matrix embed-
dings to identify similar tasks and interpolates parameters
from pre-trained lightweight experts for task adaptation.
We use $k = 2$ auxiliary experts for interpolation and op-
timize the interpolation weights along with the target ex-
pert parameters. 159 160 161 162 163 164
- AdaLoRA [30] applies singular value decomposition to
parameterize incremental updates and dynamically allo-
cates the parameter budget by pruning singular values
based on importance scores. We use a cubic schedule to
prune the rank from an initial value to the target budget,
with orthogonal regularization $\gamma = 0.1$. 165 166 167 168 169 170
- AutoLoRA [31] uses a meta-learning framework to de-
termine layer-specific optimal ranks by learning selection 171 172

173 variables for decomposed rank-1 matrices. The selection
 174 variables start at a dimension of 8, with an alternating op-
 175 timization strategy for updating model weights and selec-
 176 tion variables.

177 B.2. Other Experimental Setups

178 **Implementation Details.** All experiments are conducted
 179 on $8 \times$ NVIDIA GeForce RTX 3090 GPUs. Unless other-
 180 wise specified, we use the default configurations of the pre-
 181 trained foundation models, including the AdamW optimizer
 182 and an input resolution of 224×224 following [20, 26, 31].
 183 When w_k^i is unavailable for a task in Eq. (3), we set $\lambda_k^i = 0$.
 184 Intuitively, if task t_i already performs well without a mod-
 185 ule at location k , then a similar new task is unlikely to re-
 186 quire one there either.

187 **Settings for Vision-Language Tasks.** We strictly follow
 188 standard evaluation protocols for all VL tasks [26]. We re-
 189 port the standard accuracy metric (ACC@0.5) on valida-
 190 tion and test splits. The training configuration includes 100
 191 epochs, a dropout rate of 0.1, a warmup ratio of 0.06, a
 192 learning rate of $1e-4$, a batch size of 4, and a label smooth-
 193 ing factor of 0.1. For LoRA-based adaptation with fixed
 194 ranks, we set a default rank of $r = 8$.

195 **Settings for Vision Tasks.** For image classification, we
 196 train for 200 epochs with a batch size of 4 and a learning
 197 rate of $1e-4$. A label smoothing ratio of 0.1 is applied.
 198 Regarding data augmentation, we follow the protocols es-
 199 tablished in prior work [26], incorporating random resize
 200 cropping, random horizontal flipping, RandAug, and ran-
 201 dom erasing. Furthermore, we utilize Mixup and CutMix
 202 with alpha values set to 0.8 and 1.0, respectively, applied
 203 with an overall probability of 0.5 per batch.

204 **Settings for Language Tasks.** Following [26], we adopt
 205 their task-specific instructions. For hyperparameter opti-
 206 mization, we perform a grid search on the development
 207 set to identify the best configuration for each task. The
 208 search space includes training epochs $\{5, 7, 10\}$, learning
 209 rates $\{3e-2, 1e-4, 5e-5\}$, and batch sizes $\{8, 16, 32\}$.

210 C. Additional Experiments

211 C.1. Impact of Task Similarity

212 **Settings.** To evaluate the effectiveness of our task similarity
 213 prediction, we analyze how the predicted similarity between
 214 a learned task and a new task correlates with the accuracy
 215 gain obtained by transferring importance from that single
 216 learned task. For each new task, we consider each learned
 217 task in turn as the sole source and plot a point with its pre-
 218 dicted similarity on the x-axis and the resulting accuracy
 219 improvement on the y-axis, as shown in Fig. 2.

220 The learned task set includes the five base tasks
 221 from § 5.2.1 and five additional tasks: video captioning

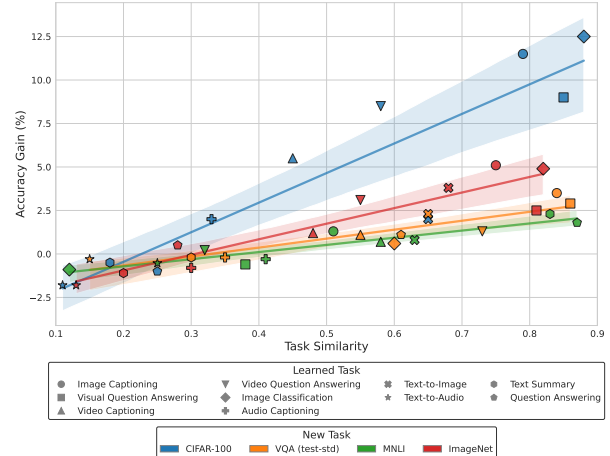


Figure 2. Impact of task similarity on fine-tuning accuracy. Each point denotes the accuracy gain when transferring module importance from a learned task (shape) to a new task (color).

(MSVD [1]), video question answering (MSVD-QA [27]),
 audio captioning (AudioCaps [16]), text-to-audio (Wav-
 Caps [19]), and text summarization (CNN/Daily Mail [7]).
 For new tasks, we use four representative benchmarks that
 cover both image and NLP modalities: CIFAR-100, VQA
 (test-std), MNLI, and ImageNet.

Results. For each new task, we fit a linear model between
 predicted similarity and accuracy improvement and report
 the coefficient of determination R^2 . The R^2 values for
 CIFAR-100, VQA, MNLI, and ImageNet are 0.82, 0.81,
 0.84, and 0.86, respectively, indicating a strong positive
 correlation between task similarity and the benefit of im-
 portance transfer. In other words, more similar learned
 tasks tend to provide larger accuracy gains when used as
 sources for module-importance transfer.

We also make a few observations from Fig. 2. (i) Similar
 learned tasks yield larger gains for relatively simple new
 tasks than for more challenging ones. CIFAR-100 (blue
 line) exhibits the steepest slope: highly similar source
 tasks can significantly boost accuracy, while dissimilar
 sources can even degrade performance. In contrast, Im-
 ageNet (red line) has a gentler slope, likely because it
 depends on diverse visual features, making it less sen-
 sitive to any single source task. (ii) Second, trans-
 ferring from the Text-to-Image task leads to a larger
 accuracy boost for ImageNet (3.8%) than for CIFAR-100
 (2.1%), even though both are image classification tasks.
 This suggests that the richer visual-text alignment
 learned in Text-to-Image pre-training aligns better with
 the complex visual concepts required by ImageNet. (iii)
 When the similarity score between a learned task and a
 new task falls below 0.35, transferring importance from
 that learned task typically reduces the new task’s
 accuracy. Motivated by this, in all experiments we re-

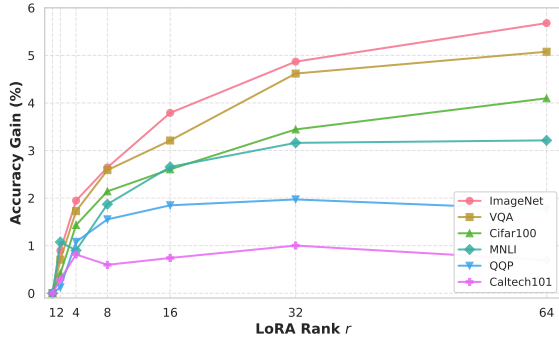


Figure 3. Accuracy gain (%) relative to $r = 1$ for varying LoRA ranks across different tasks.

Method	Function $h(r)$	Cifar100	ImgNet	Caltech101	MNLI	QQP	VQA	Avg
LoRA	$r = 8$	67.9	87.6	88.0	86.3	90.1	73.4	82.2
TaskIT	$h(r) \propto 1$	77.1	83.5	86.5	80.8	85.2	68.8	80.3
	$h(r) \propto r$	73.9	88.9	91.4	86.2	83.1	74.2	83.0
	$h(r) \propto \ln r$	80.8	88.7	91.0	85.0	91.1	74.3	85.2
	$h(r) \propto \sqrt{r}$	81.2	89.1	91.9	86.1	91.5	74.9	85.8
	$h(r) \propto \sqrt{\ln r}$	83.4	89.3	92.1	86.3	89.6	75.1	86.0

Table 2. Impact of different rank scaling functions $h(r)$. We compare the accuracy (%) of TaskIT using different functions against standard LoRA with fixed rank $r = 8$.

255 strict importance transfer to learned tasks whose similar-
 256 ity exceeds 0.35, ensuring that TaskIT primarily leverages
 257 sources that yield positive transfer.

258 C.2. Impact of Normalized Rank Contribution $h(r)$

259 **Settings.** We choose $h(r)$ by studying how accuracy scales
 260 with the LoRA rank r . First, we vary the LoRA rank
 261 $r \in \{1, \dots, 64\}$ and measure the relative accuracy gain over
 262 standard LoRA with rank $r = 1$ on six representative tasks
 263 from § 5.2.1. Second, we compare five candidate forms
 264 of $h(r)$, *i.e.* Uniform (1), Linear (r), Logarithmic ($\ln r$),
 265 Square Root (\sqrt{r}), and Root-Logarithmic ($\sqrt{\ln r}$), under a
 266 fixed memory budget, and report average accuracy in Tab. 2
 267 against standard LoRA with fixed rank $r = 8$.

268 **Results.** Fig. 3 shows that accuracy grows sub-linearly with
 269 rank. Gains from increasing r are large at small ranks but
 270 become marginal from $r = 32$ to 64 despite doubling mem-
 271 ory. Harder tasks (*e.g.* ImageNet) benefit more from higher
 272 ranks, whereas easier ones (*e.g.* Caltech101) quickly satu-
 273 rate. These trends motivate a concave $h(r)$ that rewards low
 274 ranks and penalizes linear memory growth at high ranks.

275 Tab. 2 confirms this observation. Uniform scaling yields
 276 the lowest average accuracy (80.3%) due to under-utilized
 277 capacity, while Linear scaling (83.0%) over-allocates bud-
 278 get to high ranks and overfits on simpler tasks. Logarith-
 279 mic and Square Root scaling improve accuracy to 85.2%
 280 and 85.8%, respectively. Root-Logarithmic scaling ($h(r) \propto$
 281 $\sqrt{\ln r}$) achieves the best result at 86.0%, matching the em-
 282 pirically concave trend and balancing capacity with effi-

Method	Backbone	Params (10 ⁶)	Avg Mem. (GB)	Image Avg Acc.	NLP Avg Acc.	VL Avg Acc.
AdapterFusion [22]	OFA-b	0.9	1.6	58.6	76.2	67.3
	OFA-l	3.7	3.2	61.0	77.2	70.5
UniPT [5]	OFA-b	19.4	1.1	71.2	83.4	74.1
	OFA-l	49.2	2.8	73.6	84.7	76.0
AutoLoRA [31]	OFA-b	3.2	1.8	73.8	85.7	74.4
	OFA-l	8.8	4.5	76.1	86.5	78.5
TaskIT	OFA-b	1.9	0.8	73.9	84.6	75.3
	OFA-l	5.1	2.3	<u>74.9</u>	<u>86.1</u>	<u>78.1</u>

Table 3. Performance with other backbones. Only the most accurate Zero-FT, non-LoRA, and LoRA-based baselines in the cross-modal setting are included. Best results are **bolded**.

ciency. Thus, we adopt $h(r) \propto \sqrt{\ln r}$ in all experiments. 283

284 C.3. Impact of Other Backbones

Settings. To assess architectural generalization, we extend 285
 TaskIT beyond Janus-Pro and adopt OFA [25] as the vision- 286
 language backbone following [26], using both Base (OFA- 287
 b, 180M) and Large (OFA-l, 470M) variants. We compare 288
 TaskIT with three SOTA baselines from § 5.2.1: Adapter- 289
 Fusion [22] (Zero-FT), UniPT [5] (Non-LoRA), and Au- 290
 toLoRA [31] (LoRA-based). For each backbone, we report 291
 average accuracy over the cross-modal task suite in § 5.2.1, 292
 covering Image, NLP, and VL modalities, under the same 293
 memory budget. 294

Results. As shown in Tab. 3, TaskIT offers the best trade-off 295
 between memory and accuracy across all backbones. Com- 296
 pared with AutoLoRA, TaskIT achieves similar accuracy 297
 (within 0.4% on OFA-l) while reducing memory from 4.5 298
 GB to 2.3 GB, a 48.8% saving. AdapterFusion uses fewer 299
 parameters but suffers from a lower accuracy. Although 300
 UniPT is memory-efficient on OFA-b, it scales poorly on 301
 OFA-l, requiring 49.2M trainable parameters versus only 302
 5.1M for TaskIT. 303

304 C.4. Performance with Cold-start Setting

Settings. To assess the robustness of TaskIT in a resource- 305
 scarce cold-start setting, we initialize the multi-LoRA 306
 model with only a single set of modules trained on 307
 ImageNet-21k, and evaluate fine-tuning on the 19 down- 308
 stream tasks of VTAB-1K [29] using a frozen ViT-Base 309
 backbone (86M parameters). 310

Results. Tab. 4 presents the results on VTAB-1K. Even 311
 with importance transfer limited to a single learned task 312
 (ImageNet-21k), TaskIT remains highly effective. With 313
 a memory budget factor $M_B = 0.6$, our method attains 314
 79.3% average accuracy, outperforming the Zero-FT base- 315
 line AdapterFusion (59.1%) by over 20% and highlight- 316
 ing the limitations of simple upstream module fusion, par- 317
 ticularly under cold-start with few learned tasks. More- 318
 over, at $M_B = 0.6$, TaskIT attains a higher average ac- 319

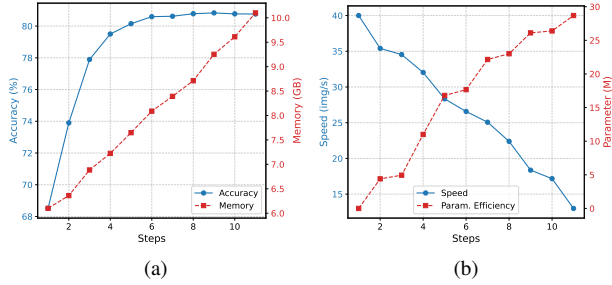


Figure 4. Impact of the memory budget M_B , showing the trade-offs between (a) memory usage (GB) and accuracy (%), and (b) training speed (img/sec) and trainable parameters (M).

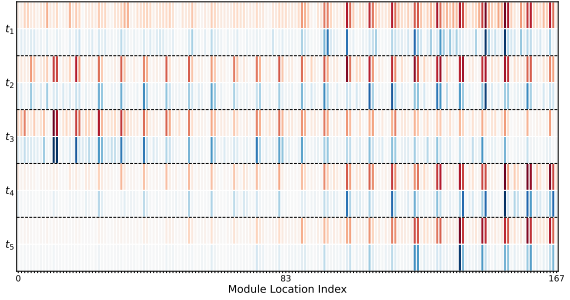


Figure 5. Predicted vs. Actual Importance of LoRA Modules. The heatmap compares the predicted importance (blue) with the actual importance (red) for LoRA modules across different location indices (x-axis) and tasks t_1 to t_5 (y-axis).

320 curacy (79.3% vs. 78.5% for UniPT) with 56% fewer
 321 trainable parameters. Compared with AutoLoRA, it nearly
 322 matches the state-of-the-art accuracy (79.5%) while using
 323 about 35% fewer parameters. This shows that our memory-
 324 aware scheduler effectively exploits the structural impor-
 325 tance of the ImageNet module to prioritize critical param-
 326 eters and maximize performance in the resource-constrained
 327 cold-start phase.

328 C.5. Impact of Memory Budget M_B

329 **Settings.** This experiment tests how the memory budget
 330 M_B affects TaskIT by varying M_B from a no-update setting
 331 to full LoRA fine-tuning, spanning 6.4 GB to 10.3 GB in 11
 332 discrete steps. Following the cross-modal setup in § 5.2.1,
 333 we report the average fine-tuning accuracy and the number
 334 of trainable parameters at each memory level. The smallest
 335 budget corresponds to reusing existing LoRAs without inser-
 336 ting new modules, while the largest budget matches full
 337 LoRA fine-tuning.

338 **Results.** As shown in Fig. 4a, accuracy increases rapidly
 339 with M_B and saturates around step 6. Larger budgets allow
 340 TaskIT to insert more modules and assign higher ranks to
 341 important ones. Once these critical modules are fully acti-
 342 vated, allocating memory to less important modules yields
 343 diminishing returns and may even cause mild overfitting.
 344 Fig. 4b further shows that the number of trainable param-
 345 eters grows with M_B , which in turn slows training due to the
 346 additional gradient computations.

347 C.6. Effectiveness of Importance Prediction

348 **Settings.** We evaluate the effectiveness of the task impor-
 349 tance predictor by comparing the predicted importance of
 350 each module with its actual importance across locations.
 351 We consider five tasks: t_1 Text-to-Image (CUB-200 [24]),
 352 t_2 Image Captioning (Flickr8k [8]), t_3 Image Classification
 353 (Pascal VOC [6]), t_4 Question Answering (SQuAD [23]),
 354 and t_5 Text Summarization (CNN/Daily Mail [7]). Fol-
 355 lowing [17], we approximate the actual importance of each

356 module by measuring the accuracy gain from individually
 357 inserting and fine-tuning that module. Fig. 5 visualizes the
 358 predicted (blue) and actual (red) importance across loca-
 359 tions. We further quantify their alignment using Spearman’s
 360 rank correlation coefficient ρ and corresponding p -values.

361 **Results.** In Fig. 5, the predicted importance closely tracks
 362 the actual importance at each location. Spearman’s ρ con-
 363 firms strong positive correlations: $\rho = 0.8862$ for t_2 ,
 364 0.8837 for t_3 , 0.8766 for t_4 , and 0.8679 for t_5 . Even for
 365 the challenging cross-modal task t_1 , we obtain $\rho = 0.6819$.
 366 All correlations are statistically significant with $p < 10^{-24}$,
 367 indicating that TaskIT accurately predicts the relative im-
 368 portance of unknown modules.

369 We also observe task-specific patterns. For t_4 and t_5 , im-
 370 portance concentrates in deeper locations, consistent with
 371 their need for high-level semantic reasoning and genera-
 372 tion. For t_3 , importance peaks in shallow and middle loca-
 373 tions, where visual features are extracted and mapped to
 374 new classes. For t_2 , importance spreads across most loca-
 375 tions, reflecting its cross-modal nature and the need to adapt
 376 both early visual processing and later text generation stages.

377 C.7. Effectiveness of Module Memory Prediction

378 **Settings.** We evaluate the accuracy of our block-based
 379 memory predictor (§ 4.3) by comparing its predictions to
 380 ground-truth estimates from PyTorch profiling tools [4] and
 381 the linear index baseline [17], using the Janus-Pro back-
 382 bone. Fig. 6 shows the normalized memory cost profiles,
 383 comparing the linear index prediction (gray), actual esti-
 384 mated memory (green), and our TaskIT prediction (blue).

385 **Results.** As shown in Fig. 6, TaskIT’s predictions closely
 386 match the actual memory usage, achieving a Pearson corre-
 387 lation coefficient of $r = 0.94$. In contrast, the linear index
 388 prediction fails to capture the structural complexities of the
 389 Transformer architecture. We observe that memory usage
 390 within each block follows a sawtooth pattern, where mod-
 391 ules like W_k , W_v , and W_{FFN}^2 incurring higher costs due to

	Param(10^6)	Natural							Specialized				Structured						Average					
		Cifar100	Caltech101	DTD	Flower102	Pets	SVHN	Sun397	Camelyon	EuroSAT	Resisc45	Retinopathy	Clevr-Count	Clevr-Dist	DMLab	KITTI-Dist	dSpr-Loc	dSpr-Ori	sNORB-Azim	sNORB-Ele	Avg Natural	Avg Specialized	Avg Structured	Average
Full[12, 14]	85.8	68.9	87.7	64.3	97.2	86.9	87.4	38.8	79.7	95.7	84.2	73.9	56.3	58.6	41.7	65.5	57.5	46.7	25.7	29.1	75.9	83.4	47.6	68.9
Linear[12, 14]	0	64.4	85.0	63.2	97.0	86.3	36.6	51.0	78.5	87.5	68.5	74.0	34.3	30.6	33.2	55.4	12.5	20.0	9.6	19.2	69.1	77.1	26.9	57.6
BitFit[15, 28]	0.10	72.8	87.0	59.2	97.5	85.3	59.9	51.4	78.7	91.6	72.9	69.8	61.5	55.6	32.4	55.9	66.6	40.0	15.7	25.1	73.3	78.3	44.1	65.2
VPT[12, 14]	0.53	78.8	90.8	65.8	98.0	88.3	78.1	49.6	81.8	96.1	83.4	68.4	68.5	60.0	46.5	72.8	73.6	47.9	32.9	37.8	78.5	82.4	55.0	72.0
RS-Bypass. [13]	0.42	64.5	88.8	73.2	99.4	90.6	63.5	<u>57.2</u>	85.5	95.2	82.4	75.2	70.4	61.0	40.2	66.8	79.2	52.6	26.0	49.3	76.7	84.6	55.7	72.3
Adapter [9, 14]	0.16	69.2	90.1	68.0	98.8	89.9	82.8	54.3	84.0	94.9	81.9	75.5	80.9	65.3	48.6	78.3	74.8	48.5	29.9	41.6	79.0	84.1	58.5	73.9
LoRA[10]	0.29	67.1	91.4	69.4	98.8	90.4	85.3	54.0	84.9	95.3	84.4	73.6	<u>82.9</u>	69.2	49.8	78.5	75.7	47.1	31.0	44.0	79.5	84.6	59.8	74.5
AdaptFormer[2, 14]	0.16	70.8	91.2	70.5	99.1	90.9	86.6	54.8	83.0	95.8	84.4	76.3	81.9	64.3	49.3	80.3	76.3	45.7	31.7	41.1	80.6	84.9	58.8	74.7
NOAH [32]	0.36	69.6	92.7	70.2	99.1	90.4	86.1	53.7	84.4	95.4	83.9	75.8	82.8	<u>68.9</u>	49.9	81.7	81.8	48.3	32.8	44.2	80.3	84.9	61.3	75.5
FacT-TK [15]	0.06	70.6	90.6	70.8	99.1	90.7	88.6	54.1	84.8	96.2	84.5	75.7	82.6	68.2	49.8	80.7	80.8	47.4	33.2	43.0	80.6	85.3	60.7	75.6
RS [13]	0.55	75.2	92.7	71.9	99.3	91.9	86.7	58.5	86.7	95.6	85.0	74.6	80.2	63.6	50.6	80.2	85.4	55.7	31.9	42.0	82.3	85.5	61.2	76.3
Convpass [14]	0.33	72.3	91.2	72.2	99.2	90.9	91.3	54.9	84.2	96.1	85.3	75.6	82.3	67.9	<u>51.3</u>	80.0	85.9	53.1	36.4	44.4	81.7	85.3	62.7	76.6
HST [18]	0.78	76.7	94.1	74.8	99.6	91.1	<u>91.2</u>	52.3	87.1	96.3	88.6	<u>76.5</u>	85.4	63.7	52.9	81.7	87.2	56.8	35.8	52.1	82.8	87.1	64.5	78.1
LoSA [20]	0.19	82.5	92.8	76.1	<u>99.7</u>	<u>90.5</u>	82.0	55.8	86.6	97.1	87.0	76.7	81.5	62.3	48.6	82.1	94.2	61.7	47.9	45.6	82.8	86.9	65.5	78.4
AdapterFusion [22]	0.23	67.4	67.8	55.5	77.9	73.3	58.3	41.3	69.8	85.9	70.7	51.9	55.1	50.8	35.0	63.8	59.0	36.2	26.0	31.4	63.1	69.6	44.7	59.1
UniPT [5]	1.21	82.7	93.0	76.2	<u>99.7</u>	90.6	82.1	55.9	86.8	97.1	87.2	77.0	81.8	62.4	49.0	82.2	94.3	61.8	48.0	45.8	82.9	87.0	65.6	78.5
AutoLoRA [31]	0.82	83.7	<u>93.7</u>	77.7	99.8	<u>91.5</u>	82.9	56.7	87.6	97.7	<u>88.4</u>	77.6	82.2	64.2	50.6	84.0	94.9	63.2	50.2	48.5	83.7	87.8	67.2	79.5
TaskIT, $M_B = 0.4$	0.37	82.6	92.7	76.7	99.7	90.7	81.4	55.5	86.8	97.1	87.4	76.7	80.4	62.1	48.8	82.4	92.7	61.4	48.7	46.7	82.8	87.0	65.3	78.3
TaskIT, $M_B = 0.5$	0.44	83.1	93.0	77.1	99.7	91.0	82.0	56.0	87.1	97.3	87.8	77.1	81.2	63.0	49.5	83.0	93.6	62.2	49.3	47.6	83.1	87.3	66.2	78.9
TaskIT, $M_B = 0.6$	0.53	<u>83.5</u>	93.5	<u>77.5</u>	99.8	91.4	82.8	56.5	<u>87.5</u>	<u>97.6</u>	<u>88.2</u>	<u>77.5</u>	81.8	63.8	50.2	<u>83.5</u>	<u>94.5</u>	62.8	49.8	48.0	<u>83.6</u>	<u>87.7</u>	<u>66.8</u>	<u>79.3</u>

Table 4. Comparison to state-of-the-art parameter-efficient finetuning methods on VTAB-1K [29]. Following standard practice, the final ‘‘Average’’ is the average of three preceding groupwise averages. Parameters denotes the number of learnable parameters excluding the final classification layer. **Best results** are bolded, and second-best underlined.

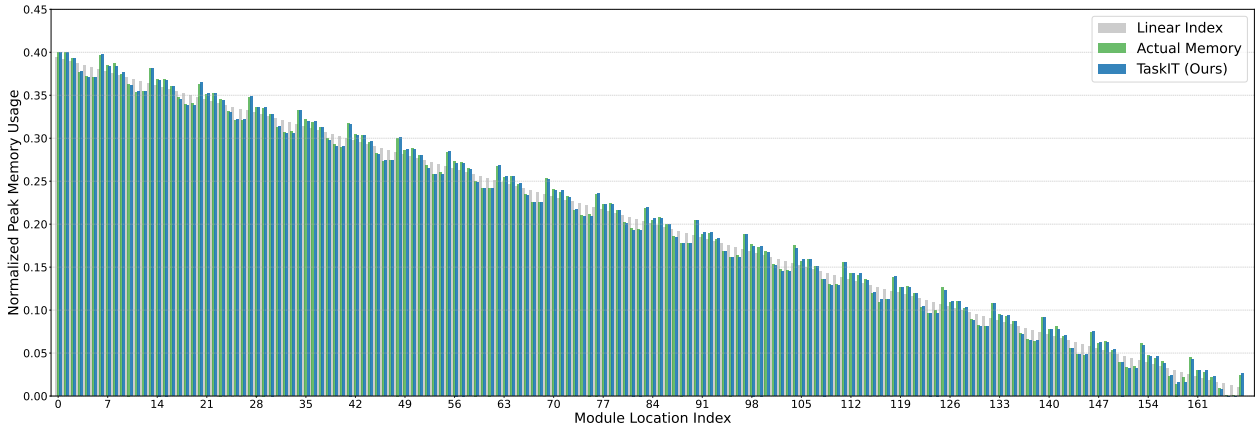


Figure 6. Comparison of normalized module memory usage predictions/estimations. TaskIT accurately captures the structural sawtooth memory pattern of Transformer blocks, aligning closely with the actual memory usage compared to the linear index baseline.

392 the retention of large activations, while W_o and W_{FFN}^1 are
393 more memory-efficient, relying on smaller activations.

394 C.8. Edge-device evaluation and cumulative over- 395 head

396 We evaluate TaskIT on a Jetson Orin Nano (8GB GPU) un-
397 der a QLoRA setting, and show the edge-device platform
398 in Fig. 7. Janus-Pro-1B is loaded in 8-bit via *BitsAndBytes*,

while fine-tuning is performed in 16-bit for 10 epochs. We 399
report accuracy, peak memory, training latency, and energy 400
consumption in Table 5, and summarize the cumulative run- 401
time and memory overhead of the major TaskIT compo- 402
nents in Table 6. The results show that TaskIT preserves the 403
best memory–accuracy trade-off under quantization, while 404
remaining practical on real edge hardware, and that its ad- 405
ditional system overhead remains modest relative to sparse 406



Figure 7. Jetson Orin Nano platform used for the edge-device evaluation.

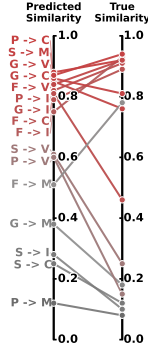


Figure 8. Validation of the task-similarity approximation by comparing the predicted similarities against the true similarities after training the new-task LoRA modules.

407 fine-tuning.

408 C.9. Task similarity approximation and sensitivity 409 analysis

410 Eq. (4) can be interpreted as a first-order Taylor approx-
411 imation of the representation update induced by adapta-
412 tion. To validate the approximation empirically, we train
413 the new-task LoRA modules w_{n+1} and compute the true
414 similarity using $x(\theta + w_i; d)$ and $x(\theta + w_{n+1}; d)$ for all
415 16 learned→new task pairs in the setting of Appendix C.1
416 in Fig. 8. The predicted similarities closely match the true
417 similarities, with Pearson correlation $r = 0.86$.

418 With batch-mean loss reduction and cosine similarity
419 dominated by direction rather than magnitude, the method
420 is not sensitive to the exact step size η . In practice, we
421 set $\eta = 10$ so that the perturbation magnitude satisfies
422 $\|\eta \nabla_x L\| / \|x\| \approx 0.1$. We further evaluate the sensitivity to
423 both η and the task-similarity threshold across four cross-
424 modal tasks, and observe stable behavior over a broad hy-
425 perparameter range in Fig. 9.

426 C.10. Module importance prediction

427 Fig. 10 validates the predicted module importance against
428 realized performance gains after insertion. The formulation
429 in Eq. (3) is evaluated only on shared backbone projection
430 matrices W_k , making it agnostic to modality-specific en-

coders, decoders, and task heads. We additionally provide
431 calibration plots on VQA, ImageNet, and MNLI by corre-
432 lating the predicted module importance with the ablated ac-
433 curacy drop, which we treat as the true importance. 434

We also ablate Eq. (3) by removing the denomina-
435 tor $|x(W_k; d)|_1$ in Table 7, replacing new-task data with
436 learned-task data, and substituting the ℓ_1 norm with ℓ_2 . The
437 ratio-based formulation used in TaskIT remains the most ro-
438 bust among these alternatives. 439

440 C.11. Module insertion scheduler and additional 441 analyses

We compare the proposed dynamic-programming scheduler
442 with greedy search and evolutionary search in Table 8. The
443 DP scheduler achieves the best downstream accuracy while
444 keeping search overhead low. This supports the claim that a
445 more structured search is beneficial even under tight mem-
446 ory constraints. 447

The DP formulation assumes that the incremental mem-
448 ory cost ΔM depends only on the nearest predecessor k_s .
449 This is exact when k_s and the candidate location lie on the
450 same sequential backward path, including across different
451 blocks and within the same FFN. When they belong to par-
452 allel $Q/K/V$ branches of the same attention block, the ap-
453 proximation can over-count at most one branch activation,
454 and the error does not propagate across blocks. 455

For scalability with the number of learned tasks, the ideal
456 feature of the new task is obtained with one backward pass,
457 and similarities to prior tasks are then computed with n
458 forward passes. Because the learned tasks are processed se-
459 quentially, peak memory remains constant, while latency
460 scales linearly with the number of prior tasks in Fig. 11. 461

We also report explicit accuracy–memory trade-offs of
462 TaskIT at larger budgets in Table 9. Under 10.5GB and
463 14.5GB budgets, TaskIT continues to improve average ac-
464 curacy as more memory is made available, showing that the
465 framework supports flexible budget-aware adaptation. 466

467 References

- 468 [1] David Chen and William B Dolan. Collecting highly
469 parallel data for paraphrase evaluation. In *Proceedings
470 of the ACL*, pages 190–200, 2011. 4
- 471 [2] Shoufa Chen, Chongjian Ge, Zhan Tong, Jiangliu
472 Wang, Yibing Song, Jue Wang, and Ping Luo. Adapt-
473 former: Adapting vision transformers for scalable vi-
474 sual recognition. *Advances in Neural Information Pro-
475 cessing Systems*, 35:16664–16678, 2022. 3, 7
- 476 [3] Alexandra Chronopoulou, Matthew E Peters, Alexan-
477 der Fraser, and Jesse Dodge. Adaptersoup: Weight
478 averaging to improve generalization of pretrained lan-
479 guage models. *arXiv preprint arXiv:2302.07027*,
480 2023. 2

Method	CIFAR-100 (VTAB-1K)				MRPC			
	Acc. (%) \uparrow	Mem. (GB) \downarrow	Lat. (min) \downarrow	Eng. (kJ) \downarrow	Acc. (%) \uparrow	Mem. (GB) \downarrow	Lat. (min) \downarrow	Eng. (kJ) \downarrow
LoRAHub	68.2	5.2	11.5	9.6	78.5	4.7	22.3	18.9
LoRA	76.2	8.9*	112.4	94.8	87.1	8.1*	205.5	180.8
UniPT	75.1	7.7	63.5	51.3	86.8	7.3	147.4	131.6
TaskIT (Ours)	76.4	6.5	65.7	52.4	87.5	5.8	143.2	125.4

Table 5. Edge-device evaluation on Jetson Orin Nano under a QLoRA setup. * indicates that memory swapping occurred.

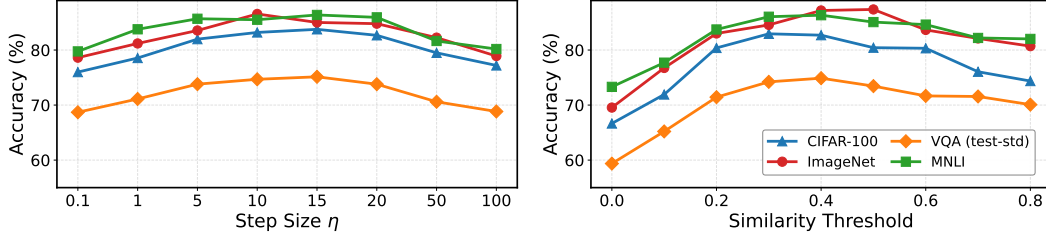


Figure 9. Sensitivity analysis with respect to the simulation step size η and the task-similarity threshold across four cross-modal tasks.

TaskIT module	Time (s)	Ratio (%)	Peak Mem. (GB)
Importance estimation	53.7	1.4	4.9
Memory profiling	92.2	2.3	4.6
DP scheduler	28.5	0.7	0.04
Sparse fine-tuning	3768.2	95.6	6.5

Table 6. Cumulative overhead of the major components of TaskIT on CIFAR-100 on Jetson. The extra overhead is modest relative to sparse fine-tuning.

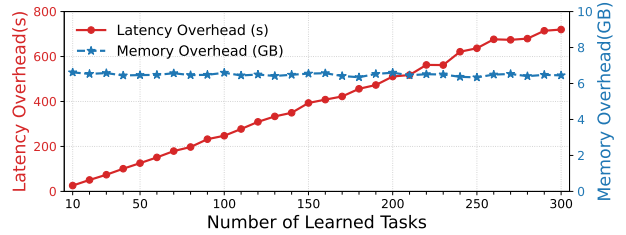


Figure 11. Scalability of task-similarity computation with respect to the number of learned tasks. Peak memory remains constant, while latency scales approximately linearly.

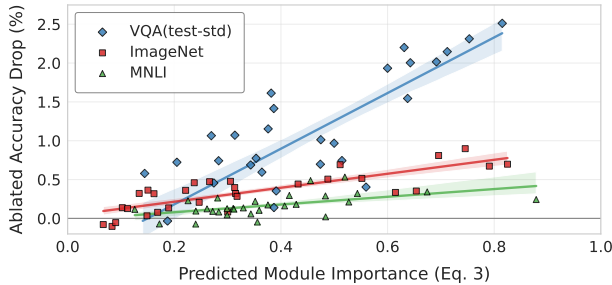


Figure 10. Calibration plots correlating predicted module importance with ablated accuracy drop on VQA, ImageNet, and MNLI.

Eq. (3) variant	Correlation \uparrow	Acc. (%) \uparrow
w/o denominator (no ratio)	0.39	80.5
w/ learned-task data	0.53	81.2
w/ ℓ_2 norm	0.68	83.3
TaskIT (ours)	0.71	83.6
Ablated accuracy drop	1.00	85.5

Table 7. Ablation of the importance metric in Eq. (3). The last row is the oracle importance measured by ablated accuracy drop after full-LoRA fine-tuning on the new task.

- [4] PyTorch Developers. Pytorch profiler. <https://docs.pytorch.org/stable/profiler.html>. 481–483
- [5] Haiwen Diao, Bo Wan, Ying Zhang, Xu Jia, Huchuan Lu, and Long Chen. Unipt: Universal parallel tuning for transfer learning with efficient parameter and memory. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 28729–28740, 2024. 3, 5, 7 484–489
- [6] Mark Everingham, Luc Van Gool, Christopher KI Williams, John Winn, and Andrew Zisserman. The pascal visual object classes (voc) challenge. *International journal of computer vision*, 88(2):303–338, 2010. 6 490–494
- [7] Karl Moritz Hermann, Tomas Kocisky, Edward Grefenstette, Lasse Espeholt, Will Kay, Mustafa Suleyman, and Phil Blunsom. Teaching machines to read and comprehend. *Advances in neural information processing systems*, 28, 2015. 4, 6 495–499
- [8] Micah Hodosh, Peter Young, and Julia Hockenmaier. Framing image description as a ranking task: Data, 501

Method	Search overhead		Sparse LoRA tuning		
	Time (s) ↓	Memory (MB) ↓	Acc. (%) ↑	Param (M) ↓	Memory (GB) ↓
Greedy search	1.2	8.2	77.5	9.1	7.9
Evolutionary search	628.5	26.7	80.1	15.4	8.0
DP search (ours)	13.2	42.3	80.4	16.8	7.8

Table 8. Comparison of the scheduler in TaskIT with greedy and evolutionary search.

Method	Nat. Avg Acc. ↑	Spec. Avg Acc. ↑	Struc. Avg Acc. ↑	Memory (GB) ↓
TaskIT ($M_B = 10.5$ GB)	84.1	87.6	68.5	10.3
TaskIT ($M_B = 14.5$ GB)	84.7	87.7	70.2	14.4

Table 9. Accuracy–memory trade-offs of TaskIT under larger memory budgets.

- 502 models and evaluation metrics. *Journal of Artificial*
503 *Intelligence Research*, 47:853–899, 2013. 6
- 504 [9] Neil Houlsby, Andrei Giurgiu, Stanislaw Jastrzebski, Bruna Morrone, Quentin De Laroussilhe, Andrea Gesmundo, Mona Attariyan, and Sylvain Gelly. Parameter-efficient transfer learning for nlp. In *Proceedings of the ICML*, pages 2790–2799. PMLR, 2019. 3, 7
- 505
506
507
508
509
- 510 [10] Edward J Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-Zhu, Yanzhi Li, Shean Wang, Lu Wang, Weizhu Chen, et al. Lora: Low-rank adaptation of large language models. *ICLR*, 1(2):3, 2022. 3, 7
- 511
512
513
- 514 [11] Chengsong Huang, Qian Liu, Bill Yuchen Lin, Tianyu Pang, Chao Du, and Min Lin. Lorahub: Efficient cross-task generalization via dynamic lora composition. In *Proceedings of the COLM*, 2024. 2
- 515
516
517
- 518 [12] Menglin Jia, Luming Tang, Bor-Chun Chen, Claire Cardie, Serge Belongie, Bharath Hariharan, and Ser-Nam Lim. Visual prompt tuning. In *European conference on computer vision*, pages 709–727. Springer, 2022. 2, 7
- 519
520
521
522
- 523 [13] Zeyinzi Jiang, Chaojie Mao, Ziyuan Huang, Ao Ma, Yiliang Lv, Yujun Shen, Deli Zhao, and Jingren Zhou. Res-tuning: A flexible and efficient tuning paradigm via unbinding tuner from backbone. *Advances in Neural Information Processing Systems*, 36:42689–42716, 2023. 2, 3, 7
- 524
525
526
527
528
- 529 [14] Shibo Jie and Zhi-Hong Deng. Convolutional bypasses are better vision transformer adapters. *arXiv preprint arXiv:2207.07039*, 2022. 2, 3, 7
- 530
531
- 532 [15] Shibo Jie and Zhi-Hong Deng. Fact: Factor-tuning for lightweight adaptation on vision transformer. In *Proceedings of the AAAI conference on artificial intelligence*, pages 1060–1068, 2023. 3, 7
- 533
534
535
- 536 [16] Chris Dongjoo Kim, Byeongchang Kim, Hyunmin Lee, and Gunhee Kim. Audiocaps: Generating captions for audios in the wild. In *Proceedings of the 2019*
537
538
- Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 119–132, 2019. 4
- 539
540
541
542
- [17] Ji Lin, Ligeng Zhu, Wei-Ming Chen, Wei-Chen Wang, Chuang Gan, and Song Han. On-device training under 256kb memory. *Advances in Neural Information Processing Systems*, 35:22941–22954, 2022. 6
- 543
544
545
546
- [18] Weifeng Lin, Ziheng Wu, Wentao Yang, Mingxin Huang, Jun Huang, and Lianwen Jin. Hierarchical side-tuning for vision transformers. *arXiv preprint arXiv:2310.05393*, 2023. 3, 7
- 547
548
549
550
- [19] Xinhao Mei, Chutong Meng, Haohe Liu, Qiuqiang Kong, Tom Ko, Chengqi Zhao, Mark D Plumbley, Yuexian Zou, and Wenwu Wang. Wavcaps: A chatgpt-assisted weakly-labelled audio captioning dataset for audio-language multimodal research. *arXiv preprint arXiv:2303.17395*, 2023. 4
- 551
552
553
554
555
556
- [20] Otniel-Bogdan Mercea, Alexey Gritsenko, Cordelia Schmid, and Anurag Arnab. Time-memory-and parameter-efficient visual adaptation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 5536–5545, 2024. 2, 3, 4, 7
- 557
558
559
560
561
562
- [21] Rui Pan, Xiang Liu, Shizhe Diao, Renjie Pi, Jipeng Zhang, Chi Han, and Tong Zhang. Lisa: Layerwise importance sampling for memory-efficient large language model fine-tuning. *Advances in Neural Information Processing Systems*, 37:57018–57049, 2024. 3
- 563
564
565
566
567
- [22] Jonas Pfeiffer, Aishwarya Kamath, Andreas Rücklé, Kyunghyun Cho, and Iryna Gurevych. Adapterfusion: Non-destructive task composition for transfer learning. 2021. 2, 5, 7
- 568
569
570
571
- [23] Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. Squad: 100,000+ questions for machine comprehension of text. *arXiv preprint arXiv:1606.05250*, 2016. 6
- 572
573
574
575

- 576 [24] Catherine Wah, Steve Branson, Peter Welinder, Pietro
577 Perona, and Serge Belongie. The caltech-ucsd birds-
578 200-2011 dataset. 2011. 6
- 579 [25] Peng Wang, An Yang, Rui Men, Junyang Lin, Shuai
580 Bai, Zhikang Li, Jianxin Ma, Chang Zhou, Jin-
581 gren Zhou, and Hongxia Yang. Ofa: Unifying ar-
582 chitectures, tasks, and modalities through a simple
583 sequence-to-sequence learning framework. In *In-
584 ternational conference on machine learning*, pages
585 23318–23340. PMLR, 2022. 5
- 586 [26] Chengyue Wu, Teng Wang, Yixiao Ge, Zeyu Lu,
587 Ruisong Zhou, Ying Shan, and Ping Luo. π -tuning:
588 Transferring multimodal foundation models with op-
589 timal multi-task interpolation. In *Proceedings of the
590 ICML, 2023*. 3, 4, 5
- 591 [27] Dejing Xu, Zhou Zhao, Jun Xiao, Fei Wu, Hanwang
592 Zhang, Xiangnan He, and Yueting Zhuang. Video
593 question answering via gradually refined attention
594 over appearance and motion. In *Proceedings of the
595 25th ACM international conference on Multimedia*,
596 pages 1645–1653, 2017. 4
- 597 [28] Elad Ben Zaken, Shauli Ravfogel, and Yoav Gold-
598 berg. Bitfit: Simple parameter-efficient fine-tuning for
599 transformer-based masked language-models. In *Pro-
600 ceedings of the ACL, 2022*. 2, 7
- 601 [29] Xiaohua Zhai, Joan Puigcerver, Alexander
602 Kolesnikov, Pierre Ruysen, Carlos Riquelme,
603 Mario Lucic, Josip Djolonga, Andre Susano Pinto,
604 Maxim Neumann, Alexey Dosovitskiy, et al. A
605 large-scale study of representation learning with the
606 visual task adaptation benchmark. *arXiv preprint
607 arXiv:1910.04867*, 2019. 5, 7
- 608 [30] Qingru Zhang, Minshuo Chen, Alexander Bukharin,
609 Nikos Karampatziakis, Pengcheng He, Yu Cheng,
610 Weizhu Chen, and Tuo Zhao. Adalora: Adaptive bud-
611 get allocation for parameter-efficient fine-tuning. In
612 *Proceedings of the ICLR, 2023*. 3
- 613 [31] Ruiyi Zhang, Rushi Qiang, Sai Ashish Somayajula,
614 and Pengtao Xie. Autolora: Automatically tuning ma-
615 trix ranks in low-rank adaptation based on meta learn-
616 ing. In *Proceedings of the ACL, 2024*. 1, 3, 4, 5, 7
- 617 [32] Yuanhan Zhang, Kaiyang Zhou, and Ziwei Liu. Neu-
618 ral prompt search. *IEEE Transactions on Pattern
619 Analysis and Machine Intelligence*, 2024. 3, 7