

-Motivation

Traditional transfer learning

- Need to store all parameters from fine-tuning
- Requires a lot of training data for target task

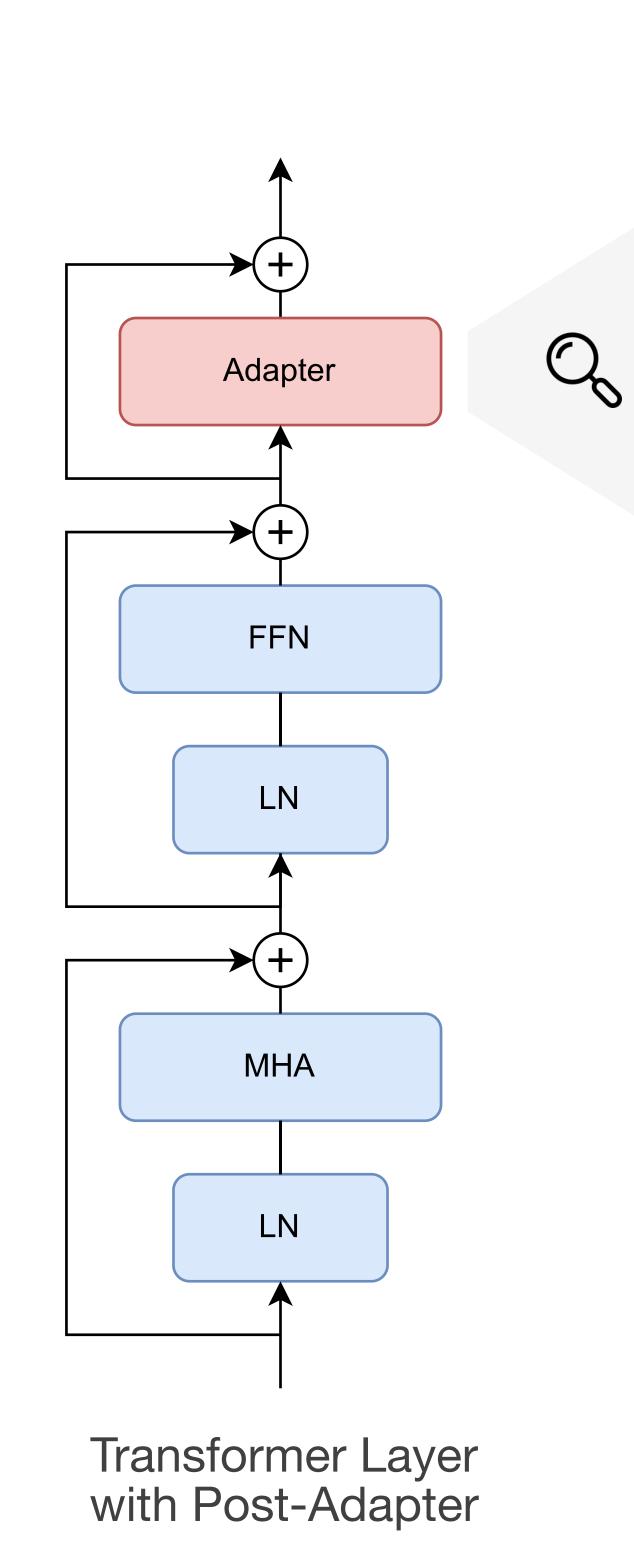
Adapter-based transfer learning

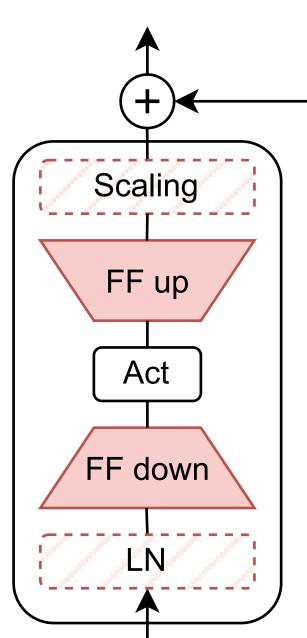
- Fine-tune only adapter modules and classifier
- Very parameter efficient: Only 0.1–0.4% of parameters required compared to fully fine-tuning a ViT-B/16
- Previous implementations introduce issues that lead to their underperformance

TL;DR

We revisit adapters and provide the first in-depth study for their use with vision transformers. We show that adapters can "strike back" and outperform more complex adaptation methods.

—Analysis: Inner Adapter Structure -





Which elements are important for the inner structure of the adapter?

We evaluate the effect of

- Biases in linear layers
- Normalization layer
- Learned, layer-wise or channel-wise scaling
- Initialization of the adapter's parameters

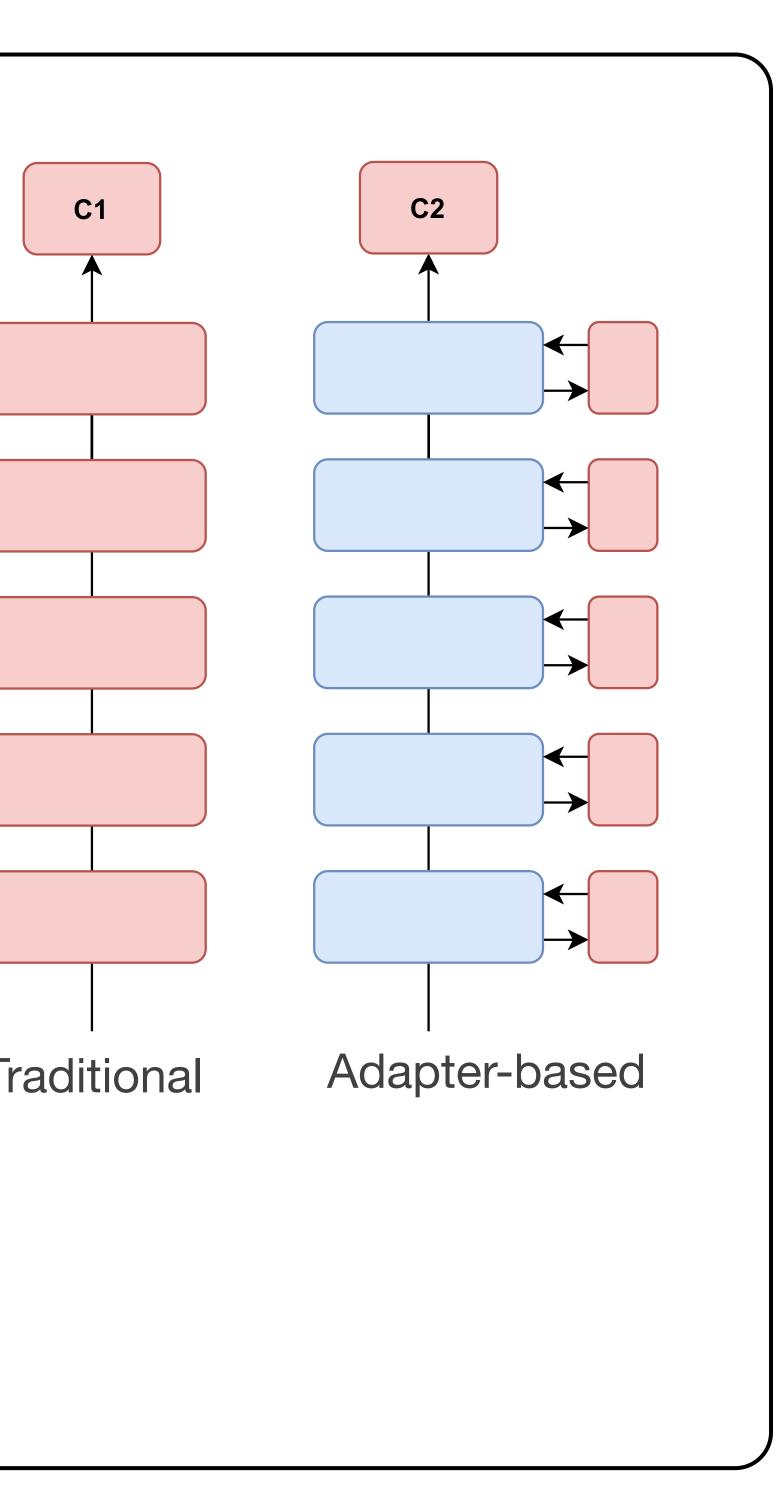
Bias	Norm	Scaling	Initialization	Accuracy (%)	$\Delta_{ extbf{base}}$
\checkmark			Houlsby	76.0	0.0
			Houlsby	75.6	-0.4
\checkmark			LoRA	75.5	-0.5
\checkmark			BERT	75.8	-0.2
\checkmark	\checkmark		Houlsby	75.9	-0.1
\checkmark	\checkmark	layer	Houlsby	75.9	-0.1
\checkmark		layer	Houlsby	76.2	+0.2
\checkmark	\checkmark	channel	Houlsby	75.8	-0.2
\checkmark		channel	Houlsby	76.5	+0.5

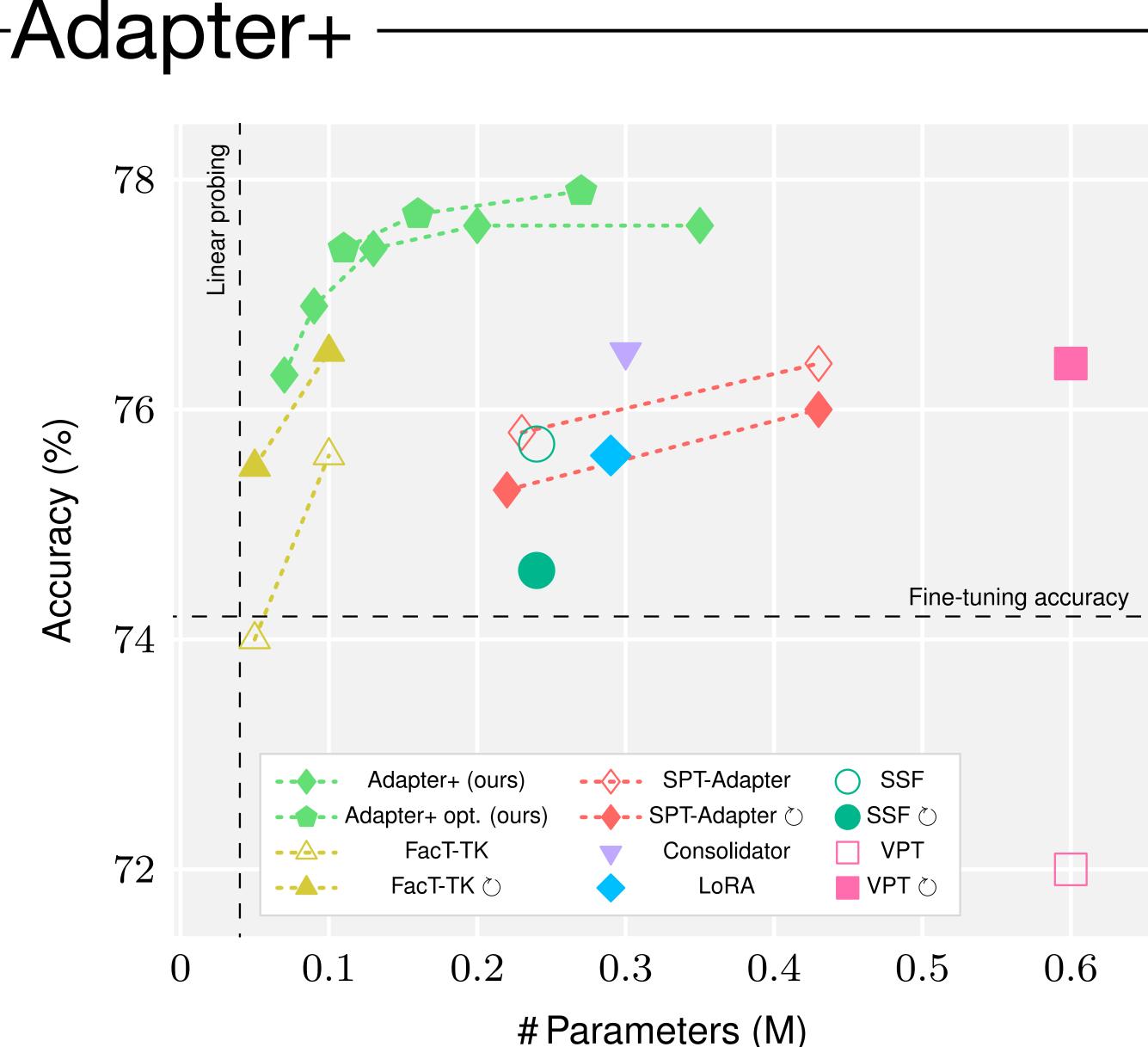
Evaluated on the VTAB *val sets*. The difference to Adapter_{base} (first row) is shown in Δ_{base} .

The optimal structure has a channel-wise scaling but no normalization layer and uses Houlsby initialization.

Adapters Strike Back

Jan-Martin O. Steitz¹ and Stefan Roth^{1,2} ¹TU Darmstadt ²hessian.Al





Parameter-accuracy characteristics of adaptation methods on the VTAB *test sets*. O: Re-evaluations with suitable data normalization and without early-stopping based on the test set.

Configuration	# Param (м)	Natural	Specialized	Structured	A
Houlsby, $r = 8$	0.39	82.9	85.5	<u>58.9</u>	
Houlsby, $r = 4$	0.24	82.9	84.9	58.3	
Pfeiffer	0.21	82.9	<u>86.1</u>	58.4	
AdaptFormer	0.19	83.0	85.0	57.4	
Adapter+	<u>0.20</u>	83.0	86.8	59.7	

Comparison of Adapter+ with different adapter configurations. Average accuracy (in %) evaluated on the VTAB val sets.

-Analysis: Adapter Position

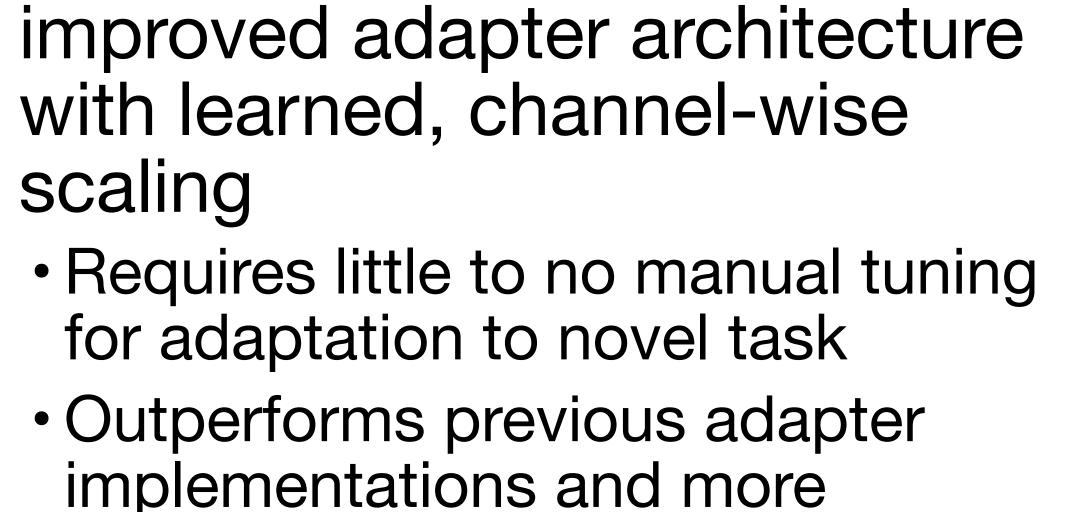
We can connect adapters in various ways to the FFN of the transformer layer. Which position is optimal?

Position	Natural	Specialized	Structured	Average
Pre	82.4	86.2	57.5	75.3
Intermediate	83.0	85.0	57.2	75.1
Parallel	83.0	86.2	57.7	75.6
Post	83.0	85.7	59.1	76.0

Accuracy (in %) evaluated on the VTAB val sets.

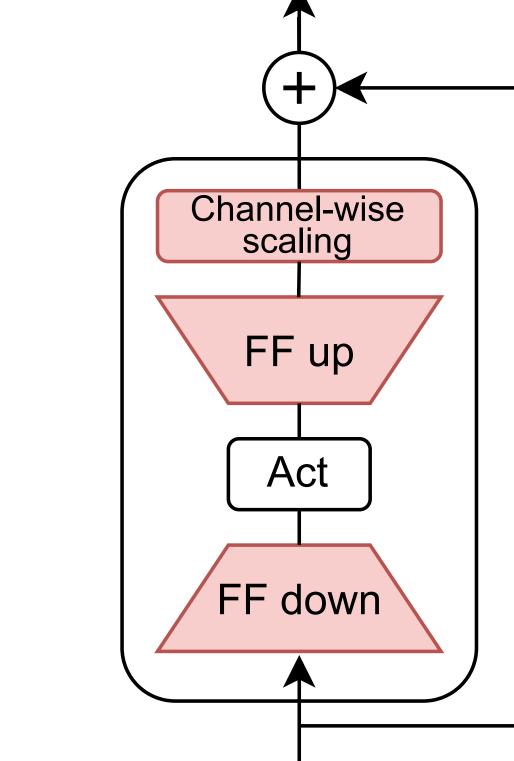
The Post-Adapter is the best configuration in the vision transformer layer.

verage <u>75.8</u> 75.4 <u>75.8</u> 75.2 **76.5**

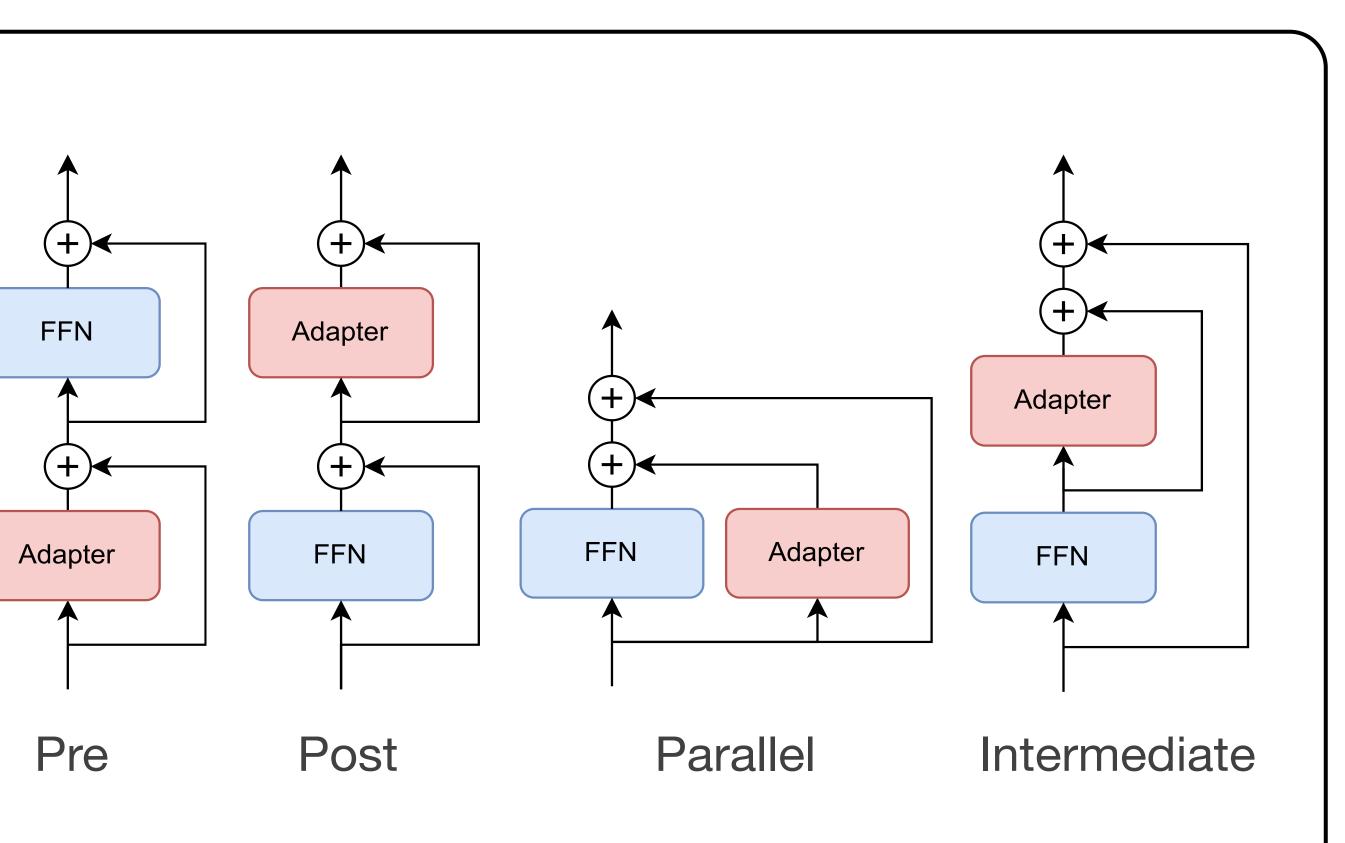


Adapter+ "strikes back": An

- implementations and more complex adaptation methods
- Reaches state-of-the-art average accuracy on the VTAB and FGVC benchmarks
- Excellent parameter-accuracy trade-off compared to other work



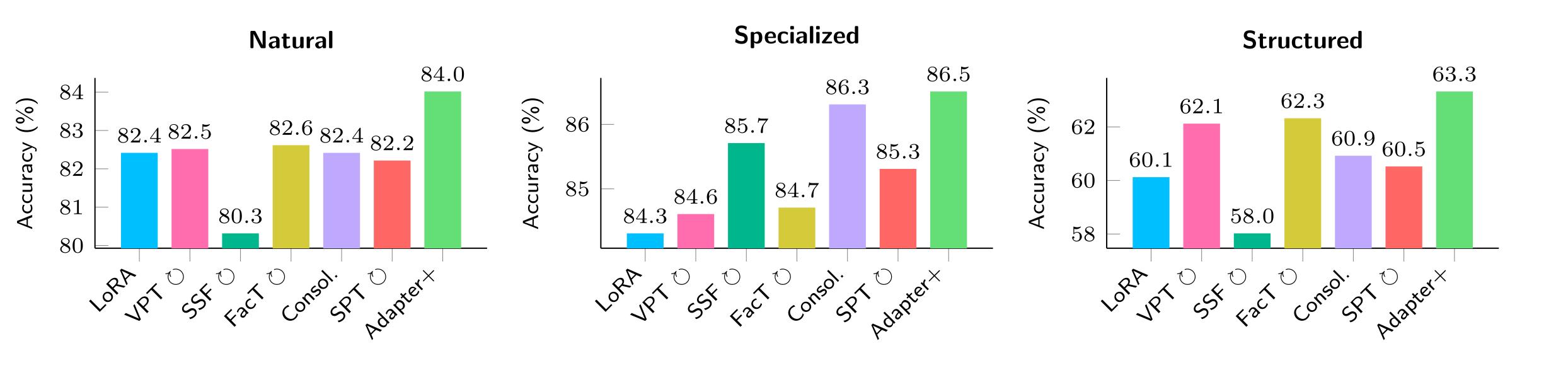
Inner structure of Adapter+ with channel-wise scaling



-Results

					Na	tural					Sp	ecial	ized					St	tructu	ired				
	# Param (M)	Cifar100	Caltech101	DTD	Flower102	Pets	SVHN	Sun397	Average	Camelyon	EuroSAT	Resisc45	Retinopathy	Average	Clevr-Count	Clevr-Dist.	DMLab	KITTI-Dist.	dSpr-Loc.	dSpr-Ori.	sNORB-Azi.	sNORB-Ele.	Average	Global Average
Full Linear	85.8 0.04							39.6 56.9	78.6 73.9			<u>87.5</u> 73.7	74.0 73.7	<u>86.3</u> 79.5						51.9 26.2			57.8 29.6	74.2 61.0
LoRA VPT-Deep \checkmark NOAH $\checkmark \dagger \bullet$ SSF \checkmark FacT-TK ₈ \bigstar FacT-TK _{≤ 32} \bigstar Consolidator ² SPT-Adapter \circlearrowright SPT-Adapter \circlearrowright	0.29 0.60 0.43 0.24 0.05 0.10 0.30 0.22 0.43	83.0 69.6 61.9 74.9 74.6 74.2 74.7	93.0 92.7 92.3 92.7 93.7 90.9 <u>94.1</u>	71.2 70.2 73.4 73.7 <u>73.6</u> 73.9 73.0	99.0 99.1 99.4 99.1 <u>99.3</u> 99.4 99.1	91.3 90.4 92.0 91.3 90.6 <u>91.6</u> 91.2	84.1 86.1 <u>90.8</u> 85.5 88.7 91.5 84.5	52.0	82.4 82.5 80.2 80.3 82.1 82.6 82.4 82.0 82.2	84.9 84.4 86.5 86.8 87.6 86.9 85.7	96.6 95.4 95.8 94.9 95.4 95.7 94.9	82.5 83.9 <u>87.5</u> 84.1 85.5 86.6 85.7	71.9 74.5 <u>75.8</u> 72.8 70.9 70.4 75.9 70.2 72.4	84.3 84.6 84.9 85.7 84.2 84.7 <u>86.3</u> 84.1 85.3	77.5 82.8 77.4 81.9 84.3 81.2 81.3	58.7 68.9 57.6 64.1 62.6 <u>68.2</u> 63.2	49.7 49.9 53.4 49.2 51.9 51.6 49.1	79.6 81.7 77.0 77.2 79.2 83.5 80.7	86.2 81.8 78.2 83.8 85.5 79.8 83.5	51.7 56.1 48.3 54.3 53.1 52.0 52.3 52.0 51.4	37.9 32.8 30.3 28.2 36.4 31.9 26.4	50.7 44.2 36.1 44.7 46.6 38.5 41.5	60.1 62.1 61.3 58.0 60.3 62.3 60.9 59.7 60.5	75.6 76.4 75.5 74.6 75.5 76.5 76.5 76.5
Adapter+, $r = 1$ Adapter+, $r = 2$ Adapter+, $r = 4$ Adapter+, $r = 8$ Adapter+, $r = 16$	0.07 0.09 0.13 0.20 0.35	85.4 <u>84.8</u> 84.6	93.0 93.8 94.2	72.7 72.7 72.3	99.2 99.2 <u>99.3</u>	90.6 90.6 90.7	85.3 86.5 87.6		83.6 83.6	87.9 87.5 <u>87.7</u>	96.8 <u>96.9</u> 97.0	85.5 85.9 86.7	71.5	85.5 85.4 85.4 85.9 86.2	83.2 <u>83.4</u> 83.2	61.0 61.6 60.9	51.6 53.6 53.8	80.1 81.4 80.3	86.1 87.3 <u>88.1</u>		30.7 34.4 35.7	46.5	60.1 61.9 <u>63.1</u> <u>63.1</u> 63.3	76.3 76.9 77.4 77.6 77.6
Adapter+, $r \in [14]$ # Adapter+, $r \in [18]$ # Adapter+, $r \in [132]$ #	0.11 0.16 0.27	85.4	93.8	72.7	99.1	90.7	87.6		83.7 <u>83.9</u> 84.0	<u>87.7</u>	96.8	86.7	72.3	85.4 85.9 86.5	83.4	60.9	53.8	80.3	<u>88.1</u>		35.7	<u>48.1</u> 47.7 47.7		77.4 <u>77.7</u> 77. 9

etailed accuracy results (in %) on the VTAB test sets. 2: Per-task hyperparameter optimization. (): re-evaluation with suitable data normalization and without early-stopping based on the test set.



-Analysis: Regularizationand Data Normalization

We investigate the influence of training regularization for transfer learning with adapters. We use Stochastic Depth and additional Dropout inside the adapter.

		Stoch
ViT	Stochastic Depth None	
	Average accuracy (in %)	evaluated

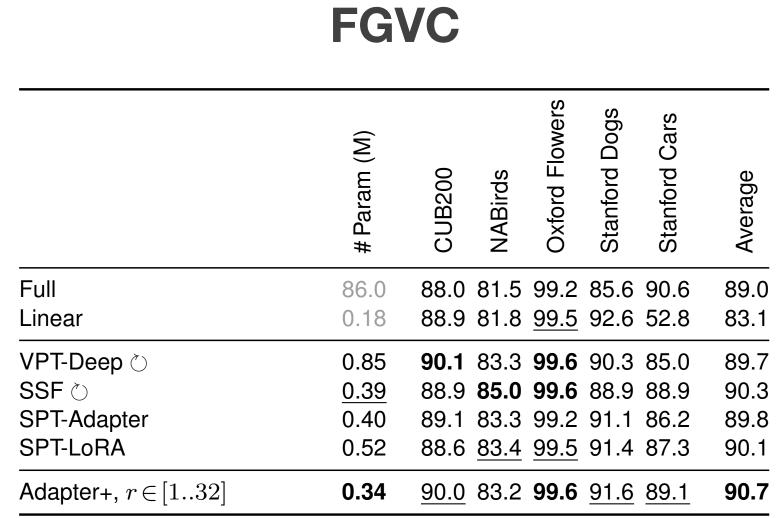
This work has been funded by the LOEWE initiative (Hesse, Germany) within the emergenCITY center.





We evaluate Adapter+ on the VTAB and FGVC benchmarks Adapter+ sets a new state-of-the-art for parameter-efficient adaptation Best on VTAB even without per-task hyperparameter optimization

VTAB



data normalization and without early ping based on the *test se*

Adapter										
hastic Depth	Dropout	None								
76.0	75.4	75.3								
74.5	74.3	73.7								

on the VTAB val sets.

Using a suitable data normalization is essential so as not to lose adaptation capabilities to compensate for the shift in normalization.

	ImageNet n	orm	In	Inception norm							
	Natural Specialized Structured	Average	Natural	Specialized	Structured	Average	Δ Average				
VPT	79.2 83.0 53.8	3 72.0	82.2	86.2	57.9	75.4	3.4				
LoRA	78.4 84.1 53.2	2 71.9	82.0	85.8	56.4	74.7	2.8				
FacT-TK	78.0 83.3 <u>5</u> 6.1	72.4	81.6	85.6	58.1	75.1	<u>2.7</u>				
Adapter+	80.5 85.0 56.0	73.9	83.0	86.8	59.7	76.5	2.6				

Average accuracy (in %) evaluated on the VTAB val sets.







