



A Closer Look at the Few-Shot Adaptation of Large Vision-Language Models





Sina Hajimiri



Ismail Ben Ayed

Jose Dolz



ÉTS Montreal

Support Set

Large Vision-Language Models

> CLIP pre-training



Zero-shot inference

$$\widehat{\boldsymbol{y}}_{c} = \frac{\exp(\boldsymbol{\nu} \cdot \boldsymbol{t}_{c}/\tau)}{\sum_{i=1}^{C} \exp(\boldsymbol{\nu} \cdot \boldsymbol{t}_{i}/\tau)}$$

 $\widehat{\mathbf{y}}_{c} = \frac{\exp(\mathbf{v} \cdot \mathbf{w}_{c}/\tau)}{\sum_{i=1}^{C} \exp(\mathbf{v} \cdot \mathbf{w}_{i}/\tau)} \longrightarrow \min_{w} \sum_{m=1}^{M} \mathcal{H}(\mathbf{y}^{(m)}, \widehat{\mathbf{y}}^{(m)})$



Vision-Language Adapters



Revisiting Linear Probing

- > The initial approximation of Linear Probing for CLIP offers limited performance.
- Zero-Shot Linear Probe (ZS-LP):
 - ✓ Pre-training head design.
 - ✓ Zero-shot weights initialization.

Method	K=1	K=2	K=4
ZS-LP	61.28	64.88	67.98
w/o DA	$57.72_{(-3.5)}\downarrow$	$61.94_{(-2.9)}\downarrow$	$65.41_{(-2.5)}\downarrow$
w/o Temp. Scaling (τ)	$58.33_{(-2.9)}\downarrow$	$59.85_{(-5.0)}\downarrow$	$59.91_{(-8.0)}\downarrow$
w/o L^2 -norm	$48.67_{(-12.6)}\downarrow$	$55.29_{(-9.6)}\downarrow$	$61.16_{(-6.8)}$
Rand. Init.	$30.42_{(-30.8)}\downarrow$	$41.86_{(-23.0)}$	$51.69_{(-16.2)}$



A well-initialized Linear Probe is all you need?

Pitfalls on Existing Few-Shot Adapters

- The strength of text prototypes varies on the difficulty of each target dataset, and the modality gap with respect to pre-training data.
- Prior works combine zero-shot and few-shot knowledge by using blending hyperparameters that control how far they go from the initial solution.

CLIP-Adapter
$$w = t$$
; $v' = v + \alpha_r f_{\varphi}(v)$



Optimum hyper-param.



SoTA Adapters require validation data to outperform a Linear Probe

How do you find the best balance per dataset?

Our Proposal: CLAP

- CLass-Adaptive linear Probing aims to train class protoypes, w_c , constraining them to remain close to text prototypes, t_c . To do so, we employ a penalty-based constraints.
- Since each category might present different zero-shot robustness and particular difficulty, we employ weighting factor per class, λ_c .

$$\min_{\boldsymbol{w}} \sum_{m=1}^{M} \mathcal{H}(\boldsymbol{y}^{(m)}, \boldsymbol{\hat{y}}^{(m)}) + \sum_{c=1}^{C} \lambda_{c} ||\boldsymbol{t}_{c} - \boldsymbol{w}_{c}||_{2}^{2}$$
Cross-entropy on Learned prototypes constrained to zero-shot

The penalty weights are treated as a learnable parameters via Augmented Lagrangian Multipliers (ALM).

$$\lambda_c^* = \frac{1}{|\beta_c^+|} \sum_{i \in \beta_c^+} \boldsymbol{y}_c^{(i)}$$

Quantitative Evaluation

Validation-free comparison

Method	K=1	K=2	K=4	K=8	K = 16
Prompt-learning methods					
CoOp _{IJCV'22} [46]	59.56	61.78	66.47	69.85	73.33
ProGrad _{ICCV'23} [13]	62.61	64.90	68.45	71.41	74.28
PLOT _{ICLR'23} [6]	62.59	65.23	68.60	71.23	73.94
Efficient transfer learning - <i>a.k.a Adapters</i>					
Zero-Shot _{ICML'21} [30]	57.71	57.71	57.71	57.71	57.71
Rand. Init LP _{ICML'21} [30]	30.42	41.86	51.69	60.84	67.54
CLIP-Adapter _{IJCV'23} [11]	58.43	62.46	66.18	69.87	73.35
TIP-Adapter _{ECCV'22} [42]	58.86	60.33	61.49	63.15	64.61
TIP-Adapter(f) ECCV'22[42]	60.29	62.26	65.32	68.35	71.40
CrossModal-LP _{CVPR'23} [24]	62.24	64.48	66.67	70.36	73.65
TaskRes(e) _{CVPR'23} [40]	61.44	65.26	68.35	71.66	74.42
ZS-LP	61.28	64.88	67.98	71.43	74.37
CLAP	62.79	66.07	69.13	72.08	74.57

Using a few-shot validation set

Method	K=1	K=2	K=4	K = 8	K = 16
Protocol in [24]: K-shots for train + $min(K, 4)$ for validation					
TIP-Adapter [42] CrossModal LP [24] CrossModal Adapter [24] CrossModal PartialFT [24]	63.3 64.1 64.4 64.7	65.9 67.0 67.6 67.2	69.0 70.3 70.8 70.5	72.2 73.0 73.4 73.6	75.1 76.0 75.9 77.1
Ours: using $K + min(K, 4)$ shots for training					
ZS-LP CLAP	64.9 66.1	68.0 69.1	71.4 72.1	73.1 73.5	75.0 75.1

Take-Home Messages

- Linear Probing (if properly designed) is a strong baseline for few-shot CLIP Adaptation.
- > Few-shot adapters should include model selection strategies based on support data.
- > CLAP is largely competitive and does not require ad-hoc adjustments per dataset.

Kine Market A Closer Look Julio Silva-F	at the Few-Shot Adaptation of Large Vision-L codriguez - Sina Hajimiri · Ismail Ben Ayed · Jose Dolz	anguage Models	
CLIP Adaptation	Our Few-Shot Adapter: CLAP > We propose a rowal water requiring hyper-parameter tuning. > We introduce Class-Adatow linear Probe (CLAP) a linear	SoTA Adapters Comparisons Average over 11 datasets Mond K-1 K-1 K-1 K-1 K-16 //mays-bounds office (K-16) //mays-bounds of (K-16) //mays-bounds of (K-16) //mays-bounds of (K-16)	Know more
transferability combining visual and text information. Pitfalls on Existing Few-Shot Adapters (CLS): $A = \{ b, c \} \ b \in \mathbb{C} \ $	classifier with prototypes constrained to remain close to the initial robust zero-shot prototypes. $\min_{m=1}^{M} \frac{\mathcal{H}(\mathbf{y}^{(m)}, \mathbf{y}^{(m)})}{\mathcal{H}(\mathbf{y}^{(m)}, \mathbf{y}^{(m)})} + \sum_{i=1}^{C} \lambda_{ii} _{\mathcal{L}_{i}} - \mathbf{w}_{ii} _{2}^{2}$ Cross-entropy Learned prototypes and	Performance Color 4000 GAD	
$ \begin{array}{c} & & & \\ & & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & & \\ & & $	 For each class, A, is fixed using zero-shot performance on support samples. Thus, better performance - larger A_e. 	Generalization on Imagenet shifts Model Search (Search) Model (Search) Tech (Search) Search (Search) Search (Search) T	
How realistic is using a validation set during few-shot adaptation?	$ \begin{array}{c} c \\ (cl.s) \\ \hline \\ cl.s) \\ cl.s) \\ \hline \\ cl.s) \\ cl.s) \\ \hline \\ cl.s) \\ c$	Conclusions > Fewstor adapters should include model selection strategies based on support data.	
A men and a set of the	$ \begin{array}{c} & & \\ & & $	CLAP is largely competitive (especially for domain generalization) and does not require ad-hoc adjustments per dataset.	JUSITO/CLAP