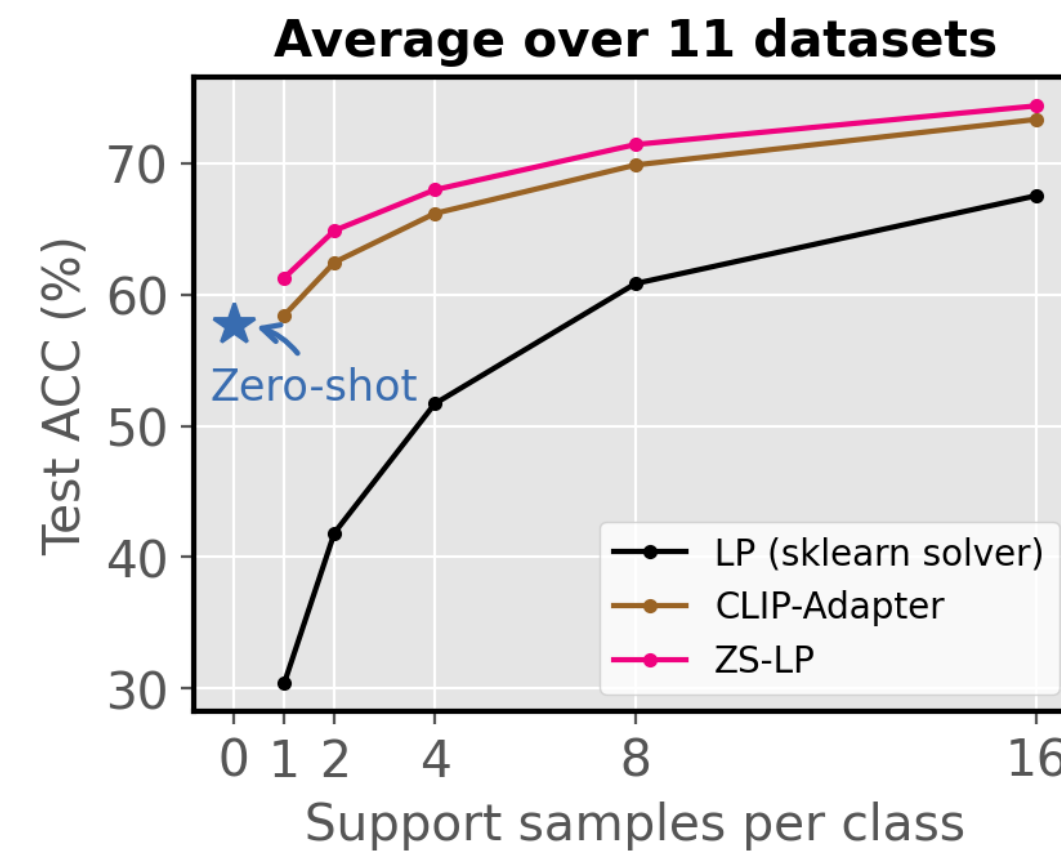


# A Closer Look at the Few-Shot Adaptation of Large Vision-Language Models

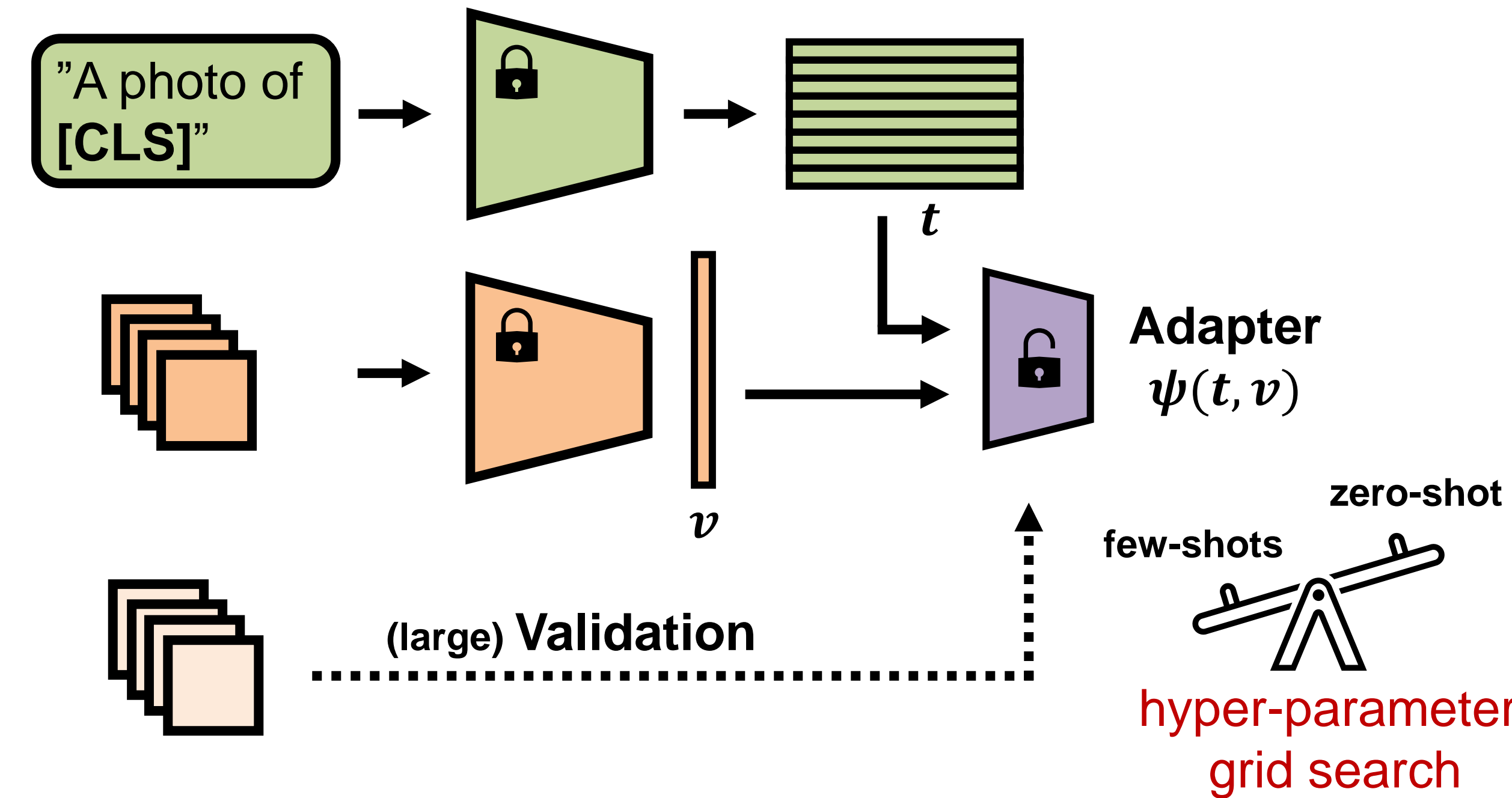
Julio Silva-Rodríguez · Sina Hajimiri · Ismail Ben Ayed · Jose Dolz · ÉTS Montreal

## VLMs Adaptation

- VLMs. present **robust zero-shot** performance
- **Few-shot Adapters** enhance the transferability **combining visual and text information**.



## Pitfalls on Existing Few-Shot Adapters



ImN	4.1	-2.1	3.1	-3.6	-2.1	1.7	-6.6	0.8	-9.9	-3	-6.3
Cal	1.9	0.4	3.6	-1.4	-1.1	1.4	-4.1	0.5	-1.4	-0.9	-2.4
Ope	1.9	0.4	3.6	-1.4	-1.1	1.4	-4.1	0.5	-1.4	-0.9	-2.4
SCa	1.9	0.4	3.6	-1.4	-1.1	1.4	-4.1	0.5	-1.4	-0.9	-2.4
Flw	-4.7	-1.6	-2.2	-6	-1.3	-5.3	-1.6	-3.8	-1.7	0.9	-2
Foo	4.1	-2.1	3.1	-3.6	-2.1	1.7	-6.6	0.8	-9.9	-3	-6.3
FGV	-4.7	-1.6	-2.2	-6	-1.3	-5.3	-1.6	-3.8	-1.7	0.9	-2
SUN	1.9	0.4	3.6	-1.4	-1.1	1.4	-4.1	0.5	-1.4	-0.9	-2.4
DTD	-3.2	-1.1	-1.4	-4.6	-2.2	-3.5	-2	-2.8	0.1	0.7	-1
EuS	-4.7	-1.6	-2.2	-6	-1.3	-5.3	-1.6	-3.8	-1.7	0.9	-1
UCF	-0.4	-0.1	1.9	-1.8	-2.0	-1.1	-2.3	-0.4	-0	0.4	0.4
ImN Cal OPe SCa Flw Foo FGV SUN DTD EuS UCF											

**! How realistic is using a validation set during few-shot adaptation?**

Otherwise, SoTA Adapters struggle to outperform a simple well-initialized Linear Probe

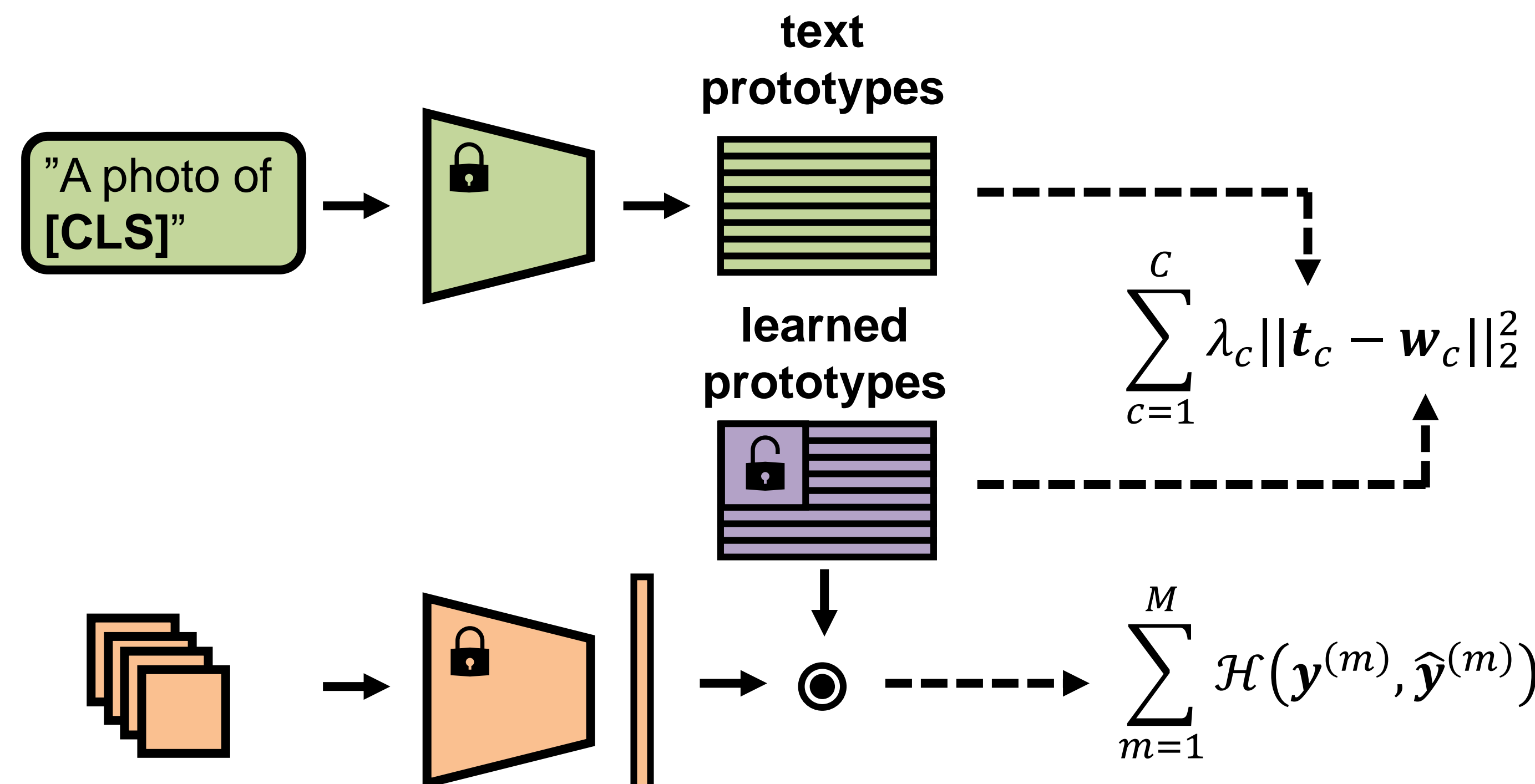
CLIP-Adapter vs. (ZS) Linear Probe

## Our Few-Shot Adapter: CLAP

- We propose a novel and simple approach that meets challenges of real-world scenarios: **not requiring hyper-parameter tuning**.
- We introduce **CLass-Adaptive linear Probe (CLAP)**, a linear classifier with **prototypes constrained to remain close to the initial, robust zero-shot prototypes**.

$$\min_w \underbrace{\sum_{m=1}^M \mathcal{H}(y^{(m)}, \hat{y}^{(m)})}_{\text{Cross-entropy on few shots}} + \underbrace{\sum_{c=1}^C \lambda_c ||t_c - w_c||_2^2}_{\text{Learned prototypes constrained to zero-shot}}$$

- For each class,  $\lambda_c$  is fixed using zero-shot performance on support samples. Thus, better performance  $\rightarrow$  larger  $\lambda_c$ .



## SoTA Adapters Comparisons

### Validation-free comparison

Method	K=1	K=2	K=4	K=8	K=16
<i>Prompt-learning methods</i>					
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
<i>Efficient transfer learning - 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	<b>62.79</b>	<b>66.07</b>	<b>69.13</b>	<b>72.08</b>	<b>74.57</b>

### 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]	63.3	65.9	69.0	72.2	75.1
CrossModal LP [24]	64.1	67.0	70.3	73.0	76.0
CrossModal Adapter [24]	64.4	67.6	70.8	73.4	75.9
CrossModal PartialFT [24]	64.7	67.2	70.5	<b>73.6</b>	<b>77.1</b>
Ours: using K + min(K, 4) shots for training					
ZS-LP	64.9	68.0	71.4	73.1	75.0
CLAP	<b>66.1</b>	<b>69.1</b>	<b>72.1</b>	<b>73.5</b>	<b>75.1</b>

## Conclusions

- 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**.

Know more!



jusiro/CLAP