

Just Shift It: Test-Time Prototype Shifting for Zero-Shot Generalization with Vision-Language Models

Supplementary Material

This document provides more details of our approach and additional experimental results, organized as follows:

- § **A** Implementation Details.
- § **B** Additional Quantitative Results with Different Random Seeds.
- § **C** Additional Ablation Studies.
- § **D** Comparison to Training-Free Methods.
- § **E** Research Impact and Limitations.

A. Implementation Details of TPS

Algorithm 2 shows more detailed pseudocode in PyTorch-like style for Test-Time Prototype Shifting over an entire dataset. We will release the models and source code to ensure reproducibility.

B. Main Results With More Random Seeds

In Sec B.1 and B.2, we run Test-Time Prototype Shifting (TPS) over 3 random seeds on both the natural distribution shifts (Table 1) and cross-dataset generalization (Table 3), respectively. The randomness comes from the image augmentation in creating a diverse minibatch for the entropy minimization objective.

B.1. Natural Distribution Shifts

From Table 7, we observe that our conclusion from Sec 4.1.1 still holds. That is, our method outperforms SoTA TPT [47] by $> 3.4\%$ on average. We also observe that augmenting the TPT-tuned class prototypes with more advanced off-the-shelf prototypes only boosts performance by a mere 0.5% on average over vanilla TPT, demonstrating TPT’s limitation in maximally leveraging these advanced prototypes.

B.2. Cross-Dataset Generalization

From Table 8, we see that our conclusion from Sec 4.1.2 remains valid. Specifically, TPS outperforms TPT [47] by $> 2\%$ on average. Similarly to Sec B.1, we observe that taking the mean of the TPT-tuned and advanced off-the-shelf prototypes increases performance by only 0.5% on average over TPT, demonstrating TPT’s inflexibility in utilizing these more robust class representations.

B.3. Context-Dependent Visual Reasoning

From Table 9, we see that our conclusion from Sec 4.2 remains valid. Specifically, TPS outperforms TPT [47] by $> 0.5\%$ on average.

C. Full Ablations

In Sec C.1, we report full ablations on TPS on the effectiveness of feature-space shift on various prototypes. These results are comparable to those reported in Sec 4.4. In Sec C.2, we include additional ablations to observe the effect of learning a class-specific shift over a universal shift for all classes. In Sec C.3, we explore variants on prototype generation using the class-agnostic CLIP ImageNet context prompt templates [42] and the class-specific descriptors generated using GPT-4 [38]. All these ablations are run over 3 random seeds.

C.1. Effect of Shift on Different Prototypes

Full comparisons between zero-shot and feature-shifted performance on all natural distribution shift and cross-domain generalization benchmark datasets over 3 random seeds are in Tables 11 and 12, respectively. We demonstrate that our conclusion from Sec 4.4.1 stills holds – that learning a small perturbation in the feature space results in performance gains of $> 4\%$ and up to 1% on average across natural distribution shift and cross-domain generalization tasks regardless of what prototypes are used.

C.2. Effect of Per-Class vs. Shared Shift

Test-time prompt tuning methods involve tuning a prompt that is shared across all classes in a dataset. Given that the tuneable prompt tokens form a portion of the text encoder input, these full prompts are then mapped to the embedding space with the encoder’s learned complex feature-space mapping. This results in non-linear perturbations from the original class prototypes. However, for our method, tuning shift parameters that are shared for all class prototypes in the feature-space means that the relative distance between class prototypes will remain constant before and after test-time shift tuning, limiting the expressive capability of the learned shift. Rather, we believe that each class prototype should be modulated by slightly different magnitudes and/or directions to provide more degrees of freedom in capturing the class-level distribution shifts in addition to the dataset-level shifts present in a domain gap.

Table 13 shows that, on average, learning a per-class shift increases performance by $> 1.2\%$ regardless of which prototypes are used. Moreover, we see that Table 14 demonstrates that, on average, learning a per-class shift increases performance by around 0.5% on average over different prototype settings. This demonstrates that learning per-class

Algorithm 2 Test-Time Prototype Shifting Pseudocode in PyTorch-like style

```
1 # Define frozen parameters
2 image_encoder = CLIPImageEncoder()
3 prototypes = load_class_prototypes()
4
5 predictions = []
6 for img, label in data_loader:
7     # Test-Time Shifting
8     shift_params = nn.Parameter(torch.zeros(num_classes, embed_dim), requires_grad=True)
9     aug_imgs = [aug(img) for i in range(batch_size - 1)]
10    imgs = torch.stack([img] + aug_imgs, dim=0)
11    image_features = image_encoder(imgs)
12
13    text_features = prototypes + shift_params
14    text_features = F.normalize(text_features, dim=-1)
15
16    logits = (logit_scale * text_features @ image_features.T)
17
18    # Confidence selection
19    entropies = compute_batch_entropies(logits)
20    top_k_idx = torch.argsort(batch_entropy, descending=False)[:k]
21
22    loss = compute_average_entropy(logits[top_k_idx])
23    optimizer.zero_grad()
24    loss.backward()
25    optimizer.step()
26
27    # Test-Time Inference
28    new_prototypes = prototypes + shift_params
29    new_prototypes = F.normalize(new_prototypes, dim=-1)
30
31    logits = (logit_scale * new_prototypes @ image_features[0].unsqueeze(0).T)
32    pred = torch.argmax(logits)
33
34    predictions.append(pred)
35
36 return predictions
```

Method	ImageNet	ImageNet-A	ImageNet-V2	ImageNet-R	ImageNet-Sketch	Average	OOD Average
	<i>Test-Time Adaptation Baselines</i>						
TPT [47]	68.96 (± 0.03)	54.47 (± 0.26)	63.46 (± 0.07)	77.10 (± 0.04)	47.93 (± 0.03)	62.38 (± 0.05)	60.74 (± 0.06)
TPT + (templates + descriptors)*	69.51 (± 0.05)	54.94 (± 0.17)	63.86 (± 0.11)	77.57 (± 0.11)	48.38 (± 0.04)	62.85 (± 0.03)	61.19 (± 0.04)
Ours	71.43 (± 0.06)	60.78 (± 0.21)	65.00 (± 0.09)	80.06 (± 0.13)	50.97 (± 0.09)	65.65 (± 0.06)	64.20 (± 0.08)

Table 7. Acc@1 of zero-shot image classification with CLIP-ViT-B/16 backbone on ImageNet and its OOD variants over 3 random seeds. Best performances are in **bold**.

shifts allows the model to capture both dataset-level and class-level distribution shifts in a domain gap.

C.3. Prototype Variants

We explore different methods for creating class prototypes. Specifically, we experiment with different forms of aggregating the text encoded with the 80 ImageNet context

prompts from CLIP [42] and our LLM-generated descriptors. The CLIP ImageNet templates are class-agnostic and add image-level characteristics whereas the descriptors are class-specific and add class-level semantic information.

Tables 15 and 16 compare three variants of pooling these CLIP templated embeddings and descriptor embeddings to obtain a single class prototype. Similarly to the conclusion

Method	Flower102	DTD	Pets	Cars	UCF101	CalTech101	Food101	SUN397	Aircraft	EuroSAT	Average
TPT [47]	68.79 (± 1)	46.79 (± 1)	87.09 (± 1)	66.38 (± 2)	67.86 (± 1)	94.13 (± 1)	84.67 (± 1)	65.41 (± 1)	23.44 (± 3)	42.78 (± 3)	64.73 (± 1)
TPT + (templates + descriptors)*	69.67 (± 1.1)	47.56 (± 0.55)	87.88 (± 0.2)	66.91 (± 1.7)	68.35 (± 2.1)	94.17 (± 1.3)	84.89 (± 0.7)	66.23 (± 1.2)	23.55 (± 3.1)	43.12 (± 1.8)	65.23 (± 0.6)
Ours	71.47 (± 1.2)	51.00 (± 0.47)	87.45 (± 0.9)	68.99 (± 1.0)	70.98 (± 2.4)	94.90 (± 1.6)	85.15 (± 0.8)	68.85 (± 1.6)	25.82 (± 0.45)	44.61 (± 1.1)	66.92 (± 0.4)

Table 8. Acc@1 of zero-shot image classification with CLIP-ViT-B/16 backbone on cross-dataset generalization over 3 random seeds. Best performances are in **bold**.

Method	seen act., seen obj.,	unseen act., seen obj.,	seen act., unseen obj.,	unseen act., unseen obj.,	Average
TPT (reprod.)	65.81 (± 1.2)	69.15 (± 1.0)	65.69 (± 0.1)	66.87 (± 0.3)	66.88 (± 0.4)
Ours (Shift)	66.67 (± 0.68)	70.31 (± 1.67)	66.00 (± 1.38)	66.67 (± 0.40)	67.41 (± 0.98)

Table 9. Acc on the Bongard-HOI benchmark with CLIP-ResNet-50 backbone over 3 random seeds. Best performances are in **bold**.

Method	ImageNet	ImageNet-R	ImageNet-Sketch	Cross-Dataset Average	Needs Support Set
SuS-X-SD-C [51]	61.65	61.69	35.88	60.49	✓
SuS-X-LC-P [51]	61.80	61.62	36.25	60.64	✓
CALIP [12]	60.57	-	-	59.34	✗
Ours (Shift + SuS-X descriptors)	63.52	63.66	37.66	61.47	✗

Table 10. Acc@1 of zero-shot image classification with CLIP-ResNet-50 backbone on ImageNet and its OOD variants. Best performances are in **bold**.

of Sec 4.4.1, we observe that in general, the gains observed using more advanced prototypes in the zero-shot setting almost directly translate to the test-time adaptation setting with shifting. In Sec 4, we present the results of our method using prototypes that are a micro average of the CLIP templates and LLM-generated descriptors.

D. Comparison to Training-Free Methods

In a similar spirit to zero-shot test-time adaptation, training-free methods perform domain adaptation without tuning any parameters. We compare our method to state-of-the-art training-free methods in Table 10. We show that our method, TPS, when using the same GPT-3-generated [2] text prompts from the official SuS-X [51] code, out-performs both CALIP [12] and SuS-X [51] without any additional image support set constructed with Stable Diffusion [44] or LAION-5B [46]. This demonstrates how a simple feature-space shift is more effective than complex training-free methods. For fair comparison, we compare TPS with SuS-X results with fixed hyperparameter settings as ours are not tuned per dataset and use the same CLIP-ResNet50 backbone.

E. Research Impact and Limitations

We propose TPS, a framework that can be used to easily and effectively improve zero-shot generalization of VLMs. Given the large-scale training of foundation VLMs, we believe it is important to understand different ways to better leverage the resulting rich multi-modal contrastive representation spaces in parameter- and runtime-constrained

settings. We propose to learn a slight perturbation to the class prototypes to maintain the overall representation quality of the pre-trained embedding space while learning a better alignment to the OOD target dataset. We hope that this framework can inspire future work to explore other tasks where learning directly in the feature space can be an efficient alternative to more complex tuning approaches.

Our work builds on the CLIP [42] representation space and uses GPT-4 [38] to generate class descriptors to create more advanced class prototypes. Thus, our model has the potential to magnify the biases of both these models. Future studies may explore how to best leverage these models' capabilities without promoting its biases.

Prompt Type	Setting	ImageNet	ImageNet-A	ImageNet-V2	ImageNet-R	ImageNet-Sketch	Average	OOD Average
Vanilla	Zero-Shot + shift	66.74 68.81 (± 0.3)	47.79 58.11 (± 1.6)	60.89 63.51 (± 1.7)	73.99 76.98 (± 0.5)	46.12 48.11 (± 0.9)	59.10 63.10 (± 0.8)	57.20 61.68 (± 0.9)
	Δ	+ 2.07	+ 10.32	+ 2.62	+ 2.99	+ 1.99	+ 4.00	+ 4.48
CoOp [60]	Zero-Shot + shift	71.51 73.76 (± 0.4)	49.71 60.43 (± 1.2)	64.20 66.84 (± 1.0)	75.21 77.39 (± 0.5)	47.99 49.08 (± 0.6)	61.72 65.50 (± 0.2)	59.28 63.44 (± 0.3)
	Δ	+ 2.25	+ 10.72	+ 2.64	+ 2.18	+ 1.09	+ 3.78	+ 4.16
CLIP templates	Zero-Shot + shift	68.35 70.39 (± 0.6)	49.95 60.47 (± 0.7)	61.97 64.66 (± 0.4)	77.59 80.70 (± 0.4)	48.21 50.38 (± 1.4)	61.21 65.32 (± 0.3)	59.43 64.05 (± 0.2)
	Δ	+ 2.04	+ 10.52	+ 2.69	+ 3.11	+ 2.17	+ 4.11	+ 4.62
Descriptors	Zero-Shot + shift	68.52 70.38 (± 0.3)	48.91 59.21 (± 0.9)	61.78 63.80 (± 0.7)	74.81 77.49 (± 1.2)	47.68 49.57 (± 0.6)	60.34 64.09 (± 0.2)	58.29 62.52 (± 0.3)
	Δ	+ 1.86	+ 10.30	+ 2.02	+ 2.68	+ 1.89	+ 3.75	+ 4.23
CLIP templates + Descriptors	Zero-Shot + shift	69.54 71.43 (± 0.6)	50.51 60.78 (± 2.1)	63.01 65.00 (± 0.9)	77.18 80.06 (± 1.3)	48.84 50.97 (± 0.9)	61.82 65.65 (± 0.6)	59.88 64.20 (± 0.8)
	Δ	+ 1.89	+ 10.27	+ 1.99	+ 2.88	+ 2.13	+ 3.83	+ 4.32

Table 11. Acc@1 for zero-shot and with feature-space shift with features initialized using different prototype generation techniques on ImageNet and its out-of-distribution variants. Results are over 3 random seeds.

Prompt Type	Setting	Flower102	DTD	Pets	Cars	UCF101	CalTech101	Food101	SUN397	Aircraft	EuroSAT	Average
Vanilla	Zero-Shot + shift	67.28 67.75 (± 1.0)	44.44 45.69 (± 1.0)	87.98 87.57 (± 1.0)	65.24 67.60 (± 2.3)	65.08 66.79 (± 2.1)	92.98 93.79 (± 0.8)	83.80 84.62 (± 0.3)	62.55 64.58 (± 0.3)	23.70 24.75 (± 3.9)	41.42 41.35 (± 0.3)	63.45 64.45 (± 0.4)
	Δ	+ 0.47	+ 1.25	- 0.41	+ 2.36	+ 1.71	+ 0.81	+ 0.82	+ 2.03	+ 1.05	- 0.07	+ 1.00
CLIP templates	Zero-Shot + shift	65.57 66.41 (± 0.5)	44.86 45.61 (± 1.9)	88.25 87.99 (± 1.0)	66.19 68.66 (± 3.1)	67.46 68.02 (± 1.1)	93.67 93.85 (± 1.4)	83.77 84.54 (± 0.8)	65.78 67.19 (± 0.5)	23.64 24.66 (± 1.3)	47.74 48.28 (± 2.0)	64.69 65.52 (± 0.5)
	Δ	+ 0.84	+ 0.75	- 0.26	+ 2.47	+ 0.56	+ 0.18	+ 0.77	+ 1.41	+ 1.02	+ 0.54	+ 0.83
Descriptors	Zero-Shot + shift	71.13 71.69 (± 1.5)	52.72 53.80 (± 2.1)	86.75 87.82 (± 1.9)	65.15 67.00 (± 1.4)	70.53 71.18 (± 1.5)	94.08 94.56 (± 0.8)	84.12 84.78 (± 0.5)	67.10 68.25 (± 1.8)	25.26 26.27 (± 0.9)	43.31 42.11 (± 1.8)	66.02 66.75 (± 0.6)
	Δ	+ 0.56	+ 1.08	+ 1.07	+ 1.85	+ 0.65	+ 0.48	+ 0.66	+ 1.15	+ 1.01	- 1.20	+ 0.73
CLIP templates + Descriptors	Zero-Shot + shift	70.52 71.47 (± 1.2)	49.94 51.00 (± 4.7)	87.22 87.45 (± 0.9)	66.48 68.99 (± 1.0)	70.24 70.98 (± 2.4)	94.12 94.90 (± 1.6)	84.47 85.15 (± 0.8)	67.55 68.85 (± 1.6)	24.69 25.82 (± 4.5)	44.14 44.61 (± 1.1)	65.94 66.92 (± 0.4)
	Δ	+ 0.95	+ 1.06	+ 0.23	+ 2.51	+ 0.74	+ 0.78	+ 0.68	+ 1.30	+ 1.13	+ 0.47	+ 0.98

Table 12. Acc@1 for zero-shot and with feature-space shift with features initialized using different prototype generation techniques on cross-domain generalization datasets. Results are over 3 random seeds.

Method	ImageNet	ImageNet-A	ImageNet-V2	ImageNet-R	ImageNet-Sketch	Average	OOD Average
Shared	71.23 (± 0.2)	56.57 (± 1.9)	64.98 (± 0.3)	79.31 (± 0.3)	50.80 (± 0.6)	64.58 (± 0.4)	62.92 (± 0.6)
Per class	71.43 (± 0.6)	60.78 (± 2.1)	65.00 (± 0.9)	80.06 (± 1.3)	50.97 (± 0.9)	65.65 (± 0.6)	64.20 (± 0.8)

Table 13. Acc@1 for learning a shared vs. per-class shift on top of different prototypes over 3 random seeds. Best performances are in **bold**.

Method	Flower102	DTD	Pets	Cars	UCF101	CalTech101	Food101	SUN397	Aircraft	EuroSAT	Average
Shared	71.36 (± 1.2)	50.49 (± 1.2)	87.46 (± 1.2)	67.33 (± 0.6)	70.77 (± 1.2)	94.35 (± 0.6)	84.82 (± 0.1)	68.12 (± 0.4)	25.27 (± 0.2)	44.67 (± 0.6)	66.47 (± 0.3)
Per-class	71.47 (± 1.2)	51.00 (± 4.7)	87.45 (± 0.9)	68.99 (± 1.0)	70.98 (± 2.4)	94.90 (± 1.6)	85.15 (± 0.8)	68.85 (± 1.6)	25.82 (± 4.5)	44.61 (± 1.1)	66.92 (± 0.4)

Table 14. Acc@1 for learning a shared vs. per-class shift on top of different prototypes over 3 random seeds. Best performances are in **bold**.

Prompt Type(s)	Pooling Method	ImageNet	ImageNet-A	ImageNet-V2	ImageNet-R	ImageNet-Sketch	Average	OOD Average
Vanilla prompt	N/A	66.74	47.79	60.89	73.99	46.12	59.10	57.20
CLIP templates + Descriptors	Macro	68.73	50.32	62.31	77.67	48.56	61.52	59.72
CLIP templates + Descriptors	Micro	69.54	50.51	63.01	77.18	48.84	61.82	59.88
CLIP templates \times Descriptors	Macro	69.03	50.73	62.22	76.91	49.07	61.59	59.73
Vanilla prompt	N/A	68.81 (± 0.3)	58.11 (± 1.6)	63.51 (± 1.7)	76.98 (± 0.5)	48.11 (± 0.9)	63.10 (± 0.8)	61.68 (± 0.9)
CLIP templates + Descriptors	Macro	70.75 (± 0.8)	60.86 (± 0.9)	64.95 (± 1.1)	80.84 (± 0.3)	50.70 (± 1.1)	65.62 (± 0.2)	64.34 (± 0.2)
CLIP templates + Descriptors	Micro	71.43 (± 0.6)	60.78 (± 2.1)	65.00 (± 0.9)	80.06 (± 1.3)	50.97 (± 0.9)	65.65 (± 0.6)	64.20 (± 0.8)
CLIP templates \times Descriptors	Macro	70.82 (± 0.2)	60.42 (± 0.6)	64.50 (± 0.5)	79.53 (± 0.9)	51.13 (± 0.2)	65.28 (± 0.1)	63.89 (± 0.2)

Table 15. Acc@1 for different variants of prototype generation, i.e. ways of combining templates and descriptors, on natural distribution shifts, over 3 random seeds. Best performances for each setting are in **bold**.

Prompt Type(s)	Pooling Method	Flower102	DTD	Pets	Cars	UCF101	CalTech101	Food101	SUN397	Aircraft	EuroSAT	Average
							<i>Zero-Shot</i>					
Vanilla prompt	N/A	67.28	44.44	87.98	65.24	65.08	92.98	83.80	62.55	23.70	41.42	63.45
CLIP templates + Descriptors	Macro	66.91	45.86	88.33	66.46	68.12	93.83	83.97	66.34	24.03	46.62	65.05
CLIP templates + Descriptors	Micro	70.52	49.94	87.22	66.48	70.24	94.12	84.47	67.55	24.69	44.14	65.94
CLIP templates × Descriptors	Macro	72.03	50.83	86.21	66.12	70.90	94.16	83.73	67.98	25.53	47.19	66.47
							<i>With Shift</i>					
Vanilla prompt	N/A	67.75 (±.10)	45.69 (±.10)	87.57 (±.10)	67.60 (±.23)	66.79 (±.21)	93.79 (±.08)	84.62 (±.03)	64.58 (±.03)	24.75 (±.39)	41.35 (±.03)	64.45 (±.04)
CLIP templates + Descriptors	Macro	67.52 (±.27)	46.43 (±.28)	88.00 (±.13)	69.04 (±.16)	68.67 (±.18)	94.16 (±.18)	84.77 (±.04)	67.70 (±.08)	24.79 (±.30)	47.09 (±.19)	65.82 (±.06)
CLIP templates + Descriptors	Micro	71.47 (±.12)	51.00 (±.47)	87.45 (±.09)	68.99 (±.10)	70.98 (±.24)	94.90 (±.16)	85.15 (±.08)	68.85 (±.16)	25.82 (±.45)	44.61 (±.11)	66.92 (±.04)
CLIP templates × Descriptors	Macro	72.53 (±.12)	52.56 (±.09)	86.15 (±.05)	68.89 (±.07)	71.44 (±.20)	94.43 (±.06)	84.44 (±.08)	69.04 (±.02)	26.51 (±.26)	45.65 (±.15)	67.16 (±.03)

Table 16. Acc@1 for different variants of knowledge injection, i.e. ways of combining templates and descriptors, over 3 random seeds on cross-dataset generalization tasks. Best performances in each setting are in **bold**.