

FAME-Vil: Multi-Tasking Vision-Language Model for Heterogeneous Fashion Tasks

Xiao Han, Xiatian Zhu, Licheng Yu, Li Zhang, Yi-Zhe Song, Tao Xiang

xiao.han@surrey.ac.uk 🗘 https://github.com/BrandonHanx/FAME-ViL







Introduction

Background: Various real-world multi-modal Vision-and-Language (V+L) tasks, including recognition, retrieval, and image captioning, exist in the fashion domain, with applications in e-commerce to enhance product discoverability, engagement, and customer conversion rates, but they are *heterogeneous* in terms of *input/output formats* and *dataset sizes*.

Motivation: Existing methods tackle these fashion tasks independently, leading to parameter redundancy and a lack of inter-task relatedness.

Contributions: (1) A task-versatile pipeline with two novel adapters, adapting a pre-trained CLIP model for all tasks; (2) An efficient and effective multitask training strategy supporting heterogeneous task modes in one unified model; (3) SOTA performance on 5 fashion tasks with 61.5% parameter saving.

Cross-Modal Retrieval (XMR)

Long sleeve relaxed-fit silk blazer in light peach. Shawl collar. Single-button closure and patch pockets at front. Breast pocket. Slits at sleeve cuffs. Vented at back.

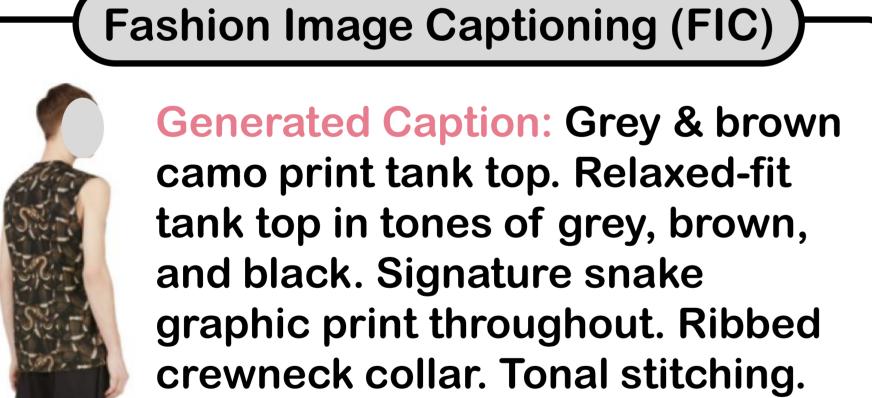


Text-Guided Image Retrieval (TGIR)



Sub-Category Recognition (SCR)

Slouchy lamb nubuck patrol hat in black. Wrinkling and light distressing throughout. Fully lined. Predicted Class: [FLAT CAPS]



Learning Objectives

 $\mathcal{L}_{ ext{XMR}} = rac{1}{2} \left[\mathcal{L}_{ ext{InfoNCE}}(extbf{T}, extbf{I}) + \mathcal{L}_{ ext{InfoNCE}}(extbf{I}, extbf{T})
ight]$ $\mathcal{L}_{\mathbf{XMR}}^{\mathbf{D}} = \frac{1}{2B} \sum_{\mathbf{S}} \left(\mathrm{KL}\left(\mathbf{s}_{b,\cdot} \parallel \tilde{\mathbf{s}}_{b,\cdot}\right) + \mathrm{KL}\left(\mathbf{s}_{\cdot,b} \parallel \tilde{\mathbf{s}}_{\cdot,b}\right) \right)$

 $\mathcal{L}_{\text{SCR}} = -\mathbb{E}_{(I,T)\sim D} \log P\left(f_{\theta}^{[f]}(I,T)\right), \mathcal{L}_{\text{SCR}}^{\text{D}} = \text{KL}\left(f_{\theta}^{[f]}(I,T) \parallel g_{\text{scr}}(I,T)\right)$

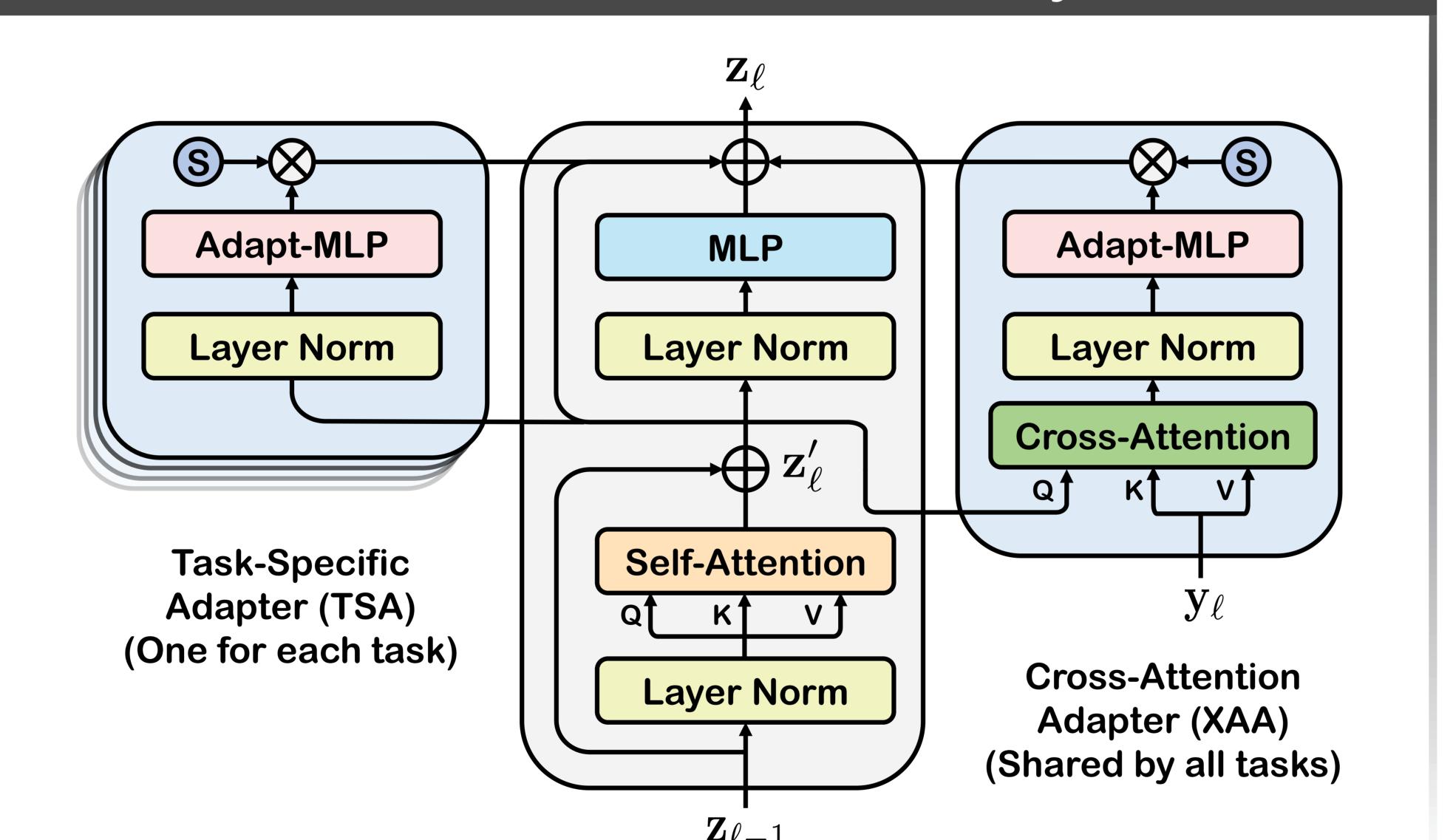
 $\mathcal{L}_{\mathbf{TGIR}} = \mathcal{L}_{\mathbf{InfoNCE}}\left((\mathbf{I}^r, \mathbf{T}), \mathbf{I}^t\right), \mathcal{L}_{\mathbf{TGIR}}^{\mathbf{D}} = \frac{1}{B} \sum_{\mathbf{KL}} \left(\mathbf{s}_{(b,b), \cdot} \parallel \tilde{\mathbf{s}}_{(b,b), \cdot}\right)$

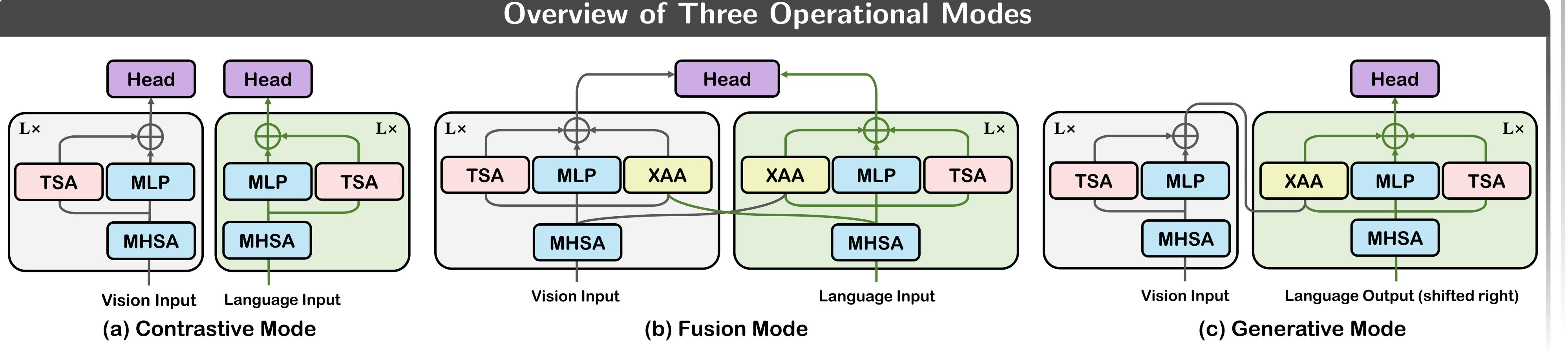
 $\mathcal{L}_{FIC} = -\mathbb{E}_{(I,T)\sim D} \sum \log P\left(T_a \middle| f_{\theta}^{[g]}(I; T_{< a})\right)$

 $\mathcal{L}_{\mathbf{FIC}}^{\mathbf{D}} = \sum \mathrm{KL}\left(f_{\theta}^{[g]}(I; T_{< a})_{a} \parallel g_{\mathrm{fic}}(I; T_{< a})_{a}\right)$

 $\mathcal{L} = \mathcal{L}_{[\mathrm{task}]} + \mathcal{L}_{[\mathrm{task}]}^{\mathrm{D}}, \quad [\mathrm{task}] \sim \{\mathrm{XMR}, \mathrm{TGIR}, \mathrm{SCR}, \mathrm{FIC}\}$

Task-Versatile Transformer Layer





Quantitative Results 1 / #Parameters Image-to-Text **Text-Guided Image Retrieval** Retrieval (Mean R@K) (Mean R@K) --- FAME-ViL (Ours) FashionViL-vit —— FashionViL (ECCV'22) KaleidoBERT (CVPR'21) --- FashionBERT(SIGIR'20) Fashion Image Text-to-Image

Further Analysis

Subcategory

Recognition

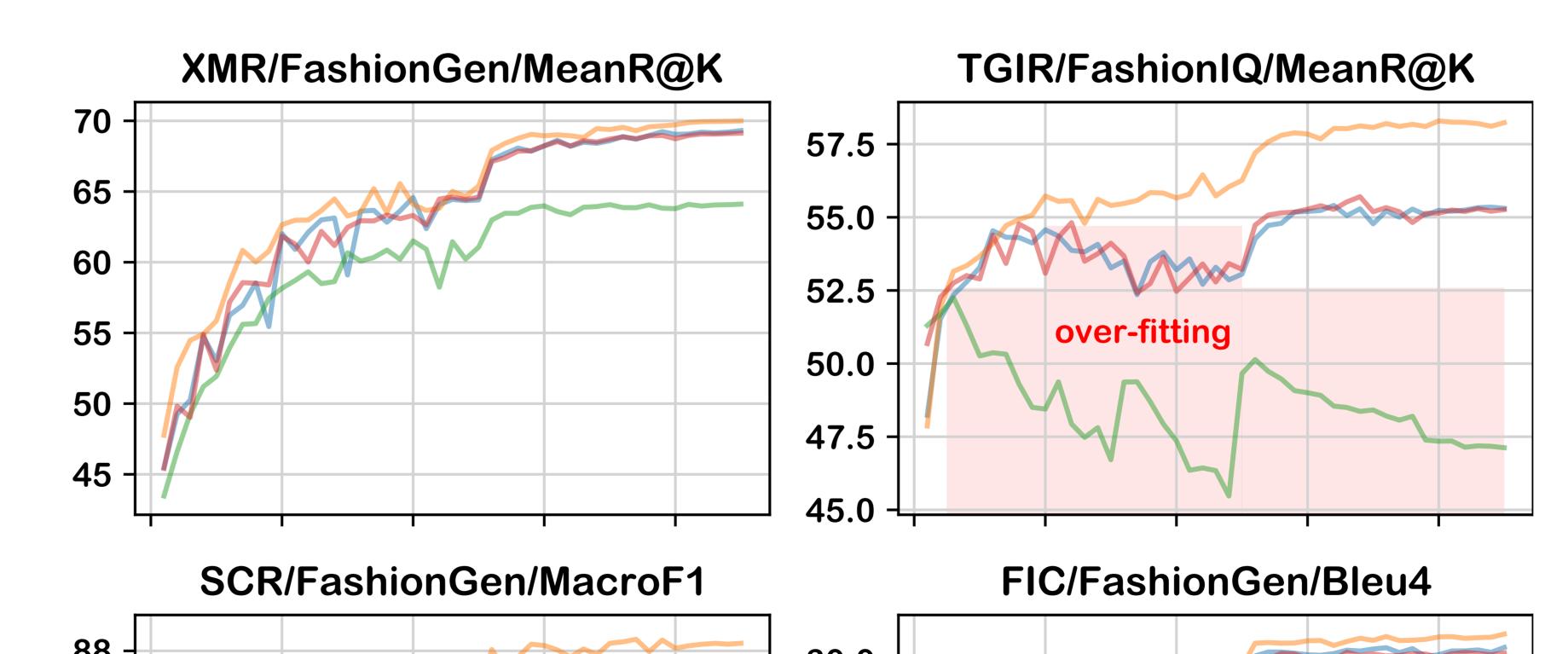
(Macro F1)

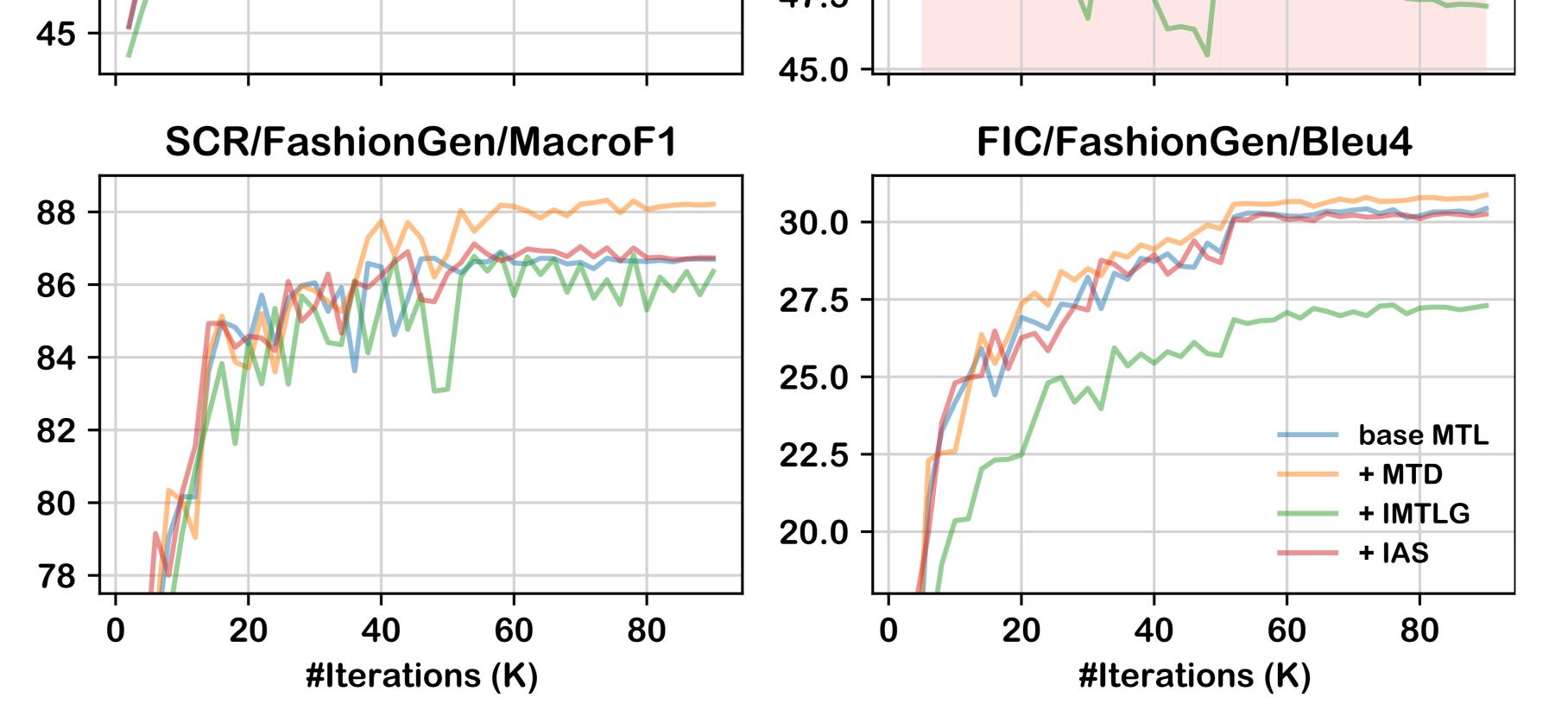
Retrieval

(Mean R@K)

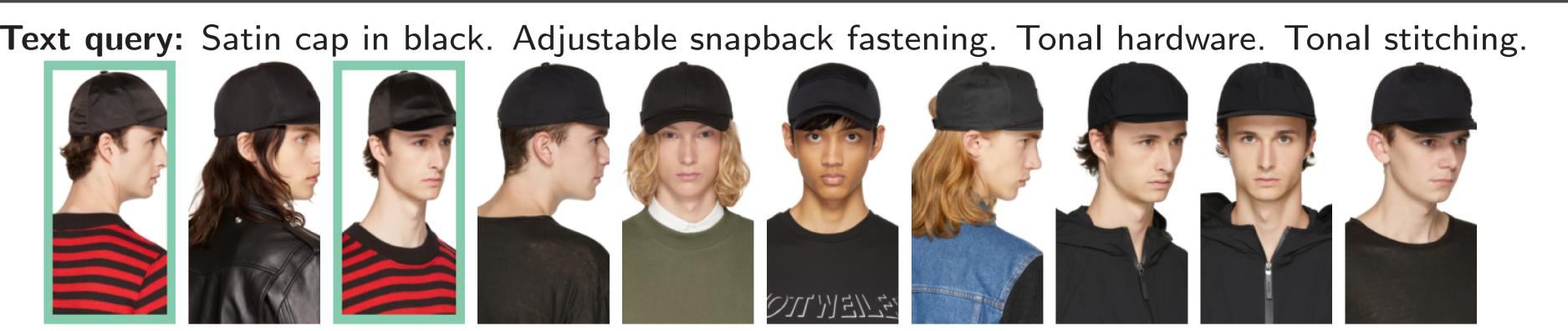
Captioning

Methods	#Params (%)	\mathcal{T}_1 : XMR		\mathcal{T}_2 : TGIR		\mathcal{T}_3 : SCR		\mathcal{T}_4 : FIC			$\bar{\Lambda}$
		μ	Δ	μ	Δ	μ	Δ	μ	Δ	μ	Δ
STL	0.0	66.30	0.0	51.87	0.0	90.34	0.0	-	-	52.13	0.0
STL + TSA	+1.35	69.99	+5.56	52.59	+1.39	90.10	-0.27	-	-	53.25	+1.67
STL + XAA	+14.70	66.30	0.0	53.83	+3.78	89.89	-0.50	63.70	0.0	68.43	+0.82
STL + TSA + XAA (FAME-ViL(ST))	+15.96	69.99	+5.56	55.47	+6.94	90.27	-0.07	63.67	-0.05	69.85	+3.10
MTL	-70.43	57.65	-13.05	49.57	-4.43	85.95	-4.86	-	-	48.29	-5.59
MTL + TSA	-70.11	67.97	+2.52	52.04	+0.33	90.32	-0.02	-	-	52.58	+0.71
MTL + XAA	-67.65	65.87	-0.65	52.59	+1.39	90.93	+0.65	60.99	-4.25	67.60	-0.72
MTL + TSA + XAA (base MTL)	-67.33	69.31	+4.54	55.41	+6.82	90.84	+0.55	65.17	+2.31	70.18	+3.56
base MTL + MTD (FAME-ViL)	-67.33	70.00	+5.56	58.29	+12.38	91.44	+1.22	65.50	+2.83	71.31	+5.50
base MTL + MTD + Uniform	-67.33	67.70	+2.11	57.31	+10.49	91.36	+1.13	65.12	+2.23	70.37	+3.99
base MTL + MTD + Round-robin	-67.33	67.79	+2.25	57.47	+10.80	91.35	+1.12	64.87	+1.84	70.37	+4.00
base MTL + IAS [32]	-67.33	69.13	+4.27	55.26	+6.54	90.51	+0.19	63.67	-0.05	69.64	+2.74
base $MTL + MTD + IAS [32]$	-67.33	70.11	+5.75	57.97	+11.76	90.88	+0.60	65.66	+3.08	71.16	+5.30
base MTL + IMTLG [46]	-67.33	64.11	-3.30	47.12	-9.16	90.21	-0.14	55.61	-12.70	64.26	-6.33
base MTL + MTD + IMTLG [46]	-67.33	67.14	+1.27	57.22	+10.31	90.09	-0.28	58.14	-9.56	68.15	+0.44
FAME-ViL (bottleneck dim. = 128)	-65.14	70.73	+6.68	58.03	+11.88	91.54	+1.33	66.20	+3.92	71.63	+5.95
FAME-ViL (bottleneck dim. = 256)	-62.67	71.77	+8.25	58.45	+12.69	91.10	+0.84	66.81	+4.88	72.03	+6.67
FAME-ViL (bottleneck dim. = 512)	-57.73	72.32	+9.08	58.51	+12.80	90.96	+0.69	66.92	+5.05	72.18	+6.91





Qualitative Results





Pleats at front. Two-pocket styling. Unlined.



Modifying text: is a green t-shirt with a light material, is more colorful.







Ground Truth Captions



White logo tank top. Relaxed-fit tank White logo tank top. Racertop in white. Ribbed scoopneck collar and armscyes. Logo print at black. Tonal logo embroidered at back hem. Tonal stitching.



Black python print shirt. Short sleeve shirt in tones of grey and black. Detailed python scale print throughout with ombre effect at bottom portions. Spread collar. Button closure at front. Tonal stitching. Single-button barrel cuffs with buttoned sleeve placket.

Generated Captions

back tank top in white. printed at front in black. Curved hem. Tonal stitching.

