

Masked Vision-Language Transformer in Fashion

Ge-Peng Ji^{†1}, Mingcheng Zhuge^{†1}, Dehong Gao¹, Deng-Ping Fan^{*2}, Christos Sakaridis² and Luc Van Gool²

¹International Core Business Unit, Alibaba Group, Hangzhou 310051, China.

²Computer Vision Lab, ETH Zürich, Zürich 8092, Switzerland.

Abstract

We present a masked vision-language transformer (MVL^T) for fashion-specific multi-modal representation. Technically, we simply utilize vision transformer architecture for replacing the BERT in the pre-training model, making MVL^T the first end-to-end framework for the fashion domain. Besides, we designed masked image reconstruction (MIR) for a fine-grained understanding of fashion. MVL^T is an extensible and convenient architecture that admits raw multi-modal inputs without extra pre-processing models (*e.g.*, ResNet), *implicitly* modeling the vision-language alignments. More importantly, MVL^T can easily generalize to various matching and generative tasks. Experimental results show obvious improvements in retrieval (rank@5: **17%**) and recognition (accuracy: **3%**) tasks over the Fashion-Gen 2018 winner Kaleido-BERT. Code is made available at <https://github.com/GewelsJI/MVLt>.

Keywords: Vision-language, masked image reconstruction, transformer, fashion, e-commercial.

1 Introduction

The emergence of transformer is drawing enormous attention from the academic community, facilitating the advancement of computer vision (CV) [3, 4] and natural language processing (NLP) [5, 6]. Benefiting from the robustness of transformers, researchers also contribute to the vision-language (VL) field [7–11] with zeal. To better utilize the pre-trained models in CV and NLP, existing general VL models are mainly based on the BERT model [12] or adopt the well-pretrained vision extractors [13, 14] or both. However, general VL methods [15–17] still struggle when applied to the fashion domain in e-commerce because they suffer from the two main issues:

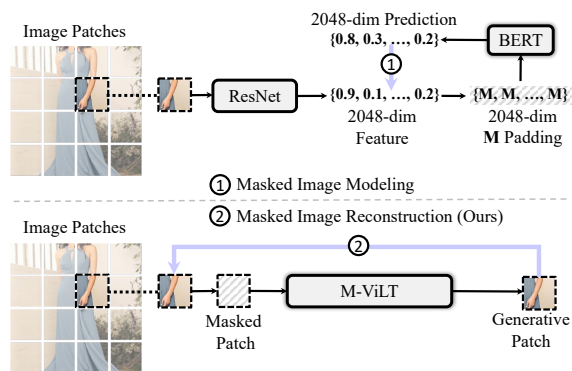


Fig. 1 Different visual reconstruction tasks for VL pre-training [1, 2] utilize masked image modeling (top) with the random masking strategy (*i.e.*, to use M padding to replace raw vectors), which reconstructs pre-extracted visual semantics (*i.e.*, probabilities) at the feature-level. We introduce a generative task named masked image reconstruction (bottom), which directly reconstructs image patches at the pixel level.

[†] Contributed equally. ^{*} Corresponding author. Work was done while Ge-Peng Ji was an research intern in Alibaba Group.

a) Insufficient Granularity. Unlike the general objects with complex backgrounds, only focusing on coarse-grained semantics is insufficient for a fashion product [18–20], as it would lead the network to generate sub-optimal results. Contrarily, the fashion-oriented framework requires more fine-grained representations, such as a suit with different materials (*e.g.*, wool, linen, and cotton) or collars (*e.g.*, band, camp, and windsor). **b) Bad Transferability.** The pre-extracted visual features are not discriminative for fashion-oriented tasks, restricting the cross-modal representations.

To address the above issues, we present a novel VL framework, termed masked vision-language transformer (MVLT). Specifically, we introduce a generative task, masked image reconstruction (MIR), for the fashion-based VL framework. Compared to previous pre-training tasks, such as masked image modeling (regression task) or masked image classification (classification task), MIR enables the network to learn more fine-grained representations via pixel-level visual knowledge (see Fig. 1). Further, inspired by pyramid vision transformer (PVT) [21], we utilize a pyramid architecture for our VL transformer. Then, we introduce the MIR task. These two improvements significantly enhance the ability to adapt to fashion-specific understanding and generative tasks, and can conduct in an end-to-end manner. To this end, MVLT can directly process the raw multi-modal inputs in dense formats (*i.e.*, linguistic tokens and visual patches) without extra (*e.g.*, ResNet) pre-processing models [22, 23]. Our main contributions are summarized as follows:

- We introduce a novel masked image reconstruction (**MIR**) task, which is the first real pixel-level generative strategy utilized in VL pre-training.
- Based on the MIR task, we present an end-to-end VL framework, called **MVLT**, for the fashion domain, greatly promoting the transferability to the downstream tasks and large-scale web applications.
- Extensive experiments show that MVLT significantly outperforms the state-of-the-art models on matching and generative tasks.

2 Background

In recent years, BERT-based pre-training models have been widely investigated in VL tasks. Many previous attempts, such as LXMERT [24], VL-BERT [25], and FashionBERT [1], were successful in a wide range of downstream applications. Experiments and discussions show that BERT is a powerful method for learning multi-modal representations, outperforming several previous CNN-based [26] or LSTM-based [27, 28] approaches. Compared to previous studies, this paper aims to develop a more efficient self-supervised objective that can be easily implemented in pre-training and provides better representations for real-world applications. Thus, we review research on masked learning strategies and end-to-end multi-modal schemes that inspired us the most.

2.1 Masked Learning Strategies

Masked modeling is the vital self-supervised task in BERT [12] and initially demonstrates outstanding abilities in natural language processing. Researchers have replicated its strength in language models because of its utility in multi-modal and vision tasks. Most VL works [16, 25, 29] transfer masked modeling into visual tokens and use a *regression* task to construct the token feature from nonsense-replace or a *classification* task to predict the token’s attribute. To reduce the difficulty in learning, Kaleido-BERT [2] optimizes masked modeling by employing a Kaleido strategy that facilitates coherent learning for multi-grained semantics. Although this work improves the performance of VL-related tasks in fashion indeed, we argue that the token-patch pre-alignment scheme by using auxiliary tool [30, 31] is still complex and impedes the application to practical settings. Another work [32] introduces the MLIM approach that strengthens the masked image modeling with an *image reconstruction* task, which shares a similar idea to ours. However, our experiments showed that requiring a model to reconstruct the entire image without any reminder is too difficult. Recently, BEiT [33] and MAE [34] utilize a BERT-style pre-training as part of the visual learner, and they discover that models are effective at learning semantics with such a scheme. These two works strengthen our conviction that converting the original masked image modeling (*i.e.*, a

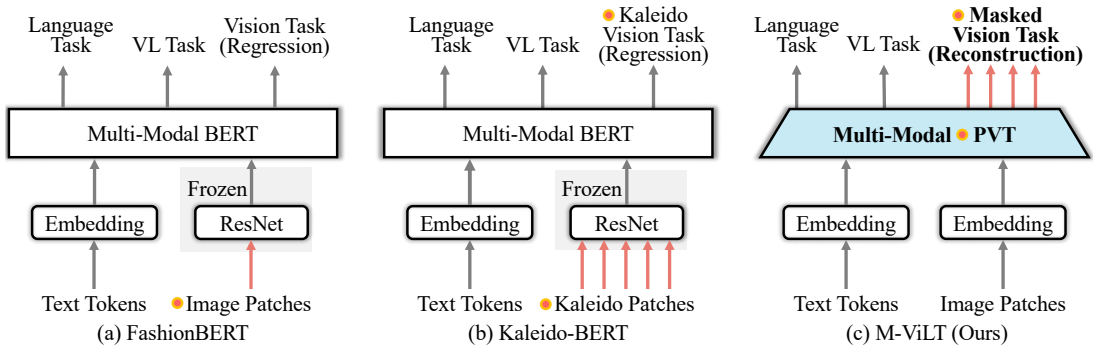


Fig. 2 Comparison of MVLT to cutting-edge fashion-oriented VL frameworks. FashionBERT (a) utilizes a language-based encoder (*i.e.*, BERT) to extract VL representations with single-scale visual input (*i.e.*, image patches). Kaleido-BERT (b) extends it with two upgrades: adds five fixed-scale inputs (*i.e.*, Kaleido patches) to acquire hierarchical visual features and designs Kaleido vision tasks to fully learn VL representations. However, the visual embedding of these models is frozen (*i.e.*, without parameter updating), thus, a lack of domain-specific visual knowledge severely hinders their transferability. Differently, our MVLT (c) adaptively learns hierarchical features by introducing masked vision tasks in an end-to-end framework, significantly boosting the VL-related understanding and generation.

regression task) to a masked image reconstruction task is possible. However, our primary goal is to design a generative pretext task that makes the multi-modal modeling in VL pre-training easier while eliminating the need for using prior knowledge. It will be extremely helpful in our practical application setting with billion-level data.

2.2 End-To-End Multi-Modal Schemes

Pixel-BERT [35] is the first method to consider end-to-end pre-training. It employs 2×2 max-pooling layers to reduce the spatial dimension of image features, with each image being downsampled 64 times. Although this work sets a precedent for end-to-end training, such a coarse and rigid method cannot work well in practical settings because it is simply combined with a ResNet [13] as part of joint pre-training, without considering the loss in speed and performance. Recently, VX2TEXT [36] proposes to convert all modalities into language space and then perform end-to-end pre-training using a relaxation scheme. Though it is exciting to translate all the modalities into a unified latent space, it ignores that the usage of data extracted by pre-trained methods as input to the model cannot be regarded as an end-to-end framework. According to the timeline, ViLT [37] is the first method that indeed investigates an end-to-end framework via replacing region- or grid-based features with patch-based projections. However, without other designs, it cannot obtain

competitive performance since it is just a vanilla extension of ViT [3]. Grid-VLP [38] is similar to ViLT, but it takes a further step by demonstrating that using a pre-trained CNN network as the visual backbone can improve performance on downstream tasks. SOHO [39] takes the entire image as input and creates a visual dictionary to affine the local region. However, this method does not fit fashion-specific applications due to the lack of reliable alignment information. As a result, the vision dictionary may merely learn the location of the background or foreground rather than complex semantics. FashionVLP [40] uses a feedback strategy to achieve better retrieval performance. In practice, they use the well-pretrained knowledge extracted from ResNet and then model the whole, cropped, and landmark representations. Besides, they adopt Faster-RCNN as an object detector for popping out RoI candidates. Besides, some works are designed for end-to-end pre-training [41–43], but they are used for specific tasks and are not directly applicable to our research.

Despite existing methods employing different approaches to construct an end-to-end scheme, solutions that forgo pre-trained methods (*e.g.*, ResNet, BERT) and use raw data (*i.e.*, text, image) as inputs remain under-explored and are needed urgently in multi-modal applications.

Remarks. As shown in Fig. 2, similar to the existing two fashion-based approaches, *i.e.*, FashionBERT (a) and Kaleido-BERT (b), the proposed MVLT (c) is also a patch-based VL learner,

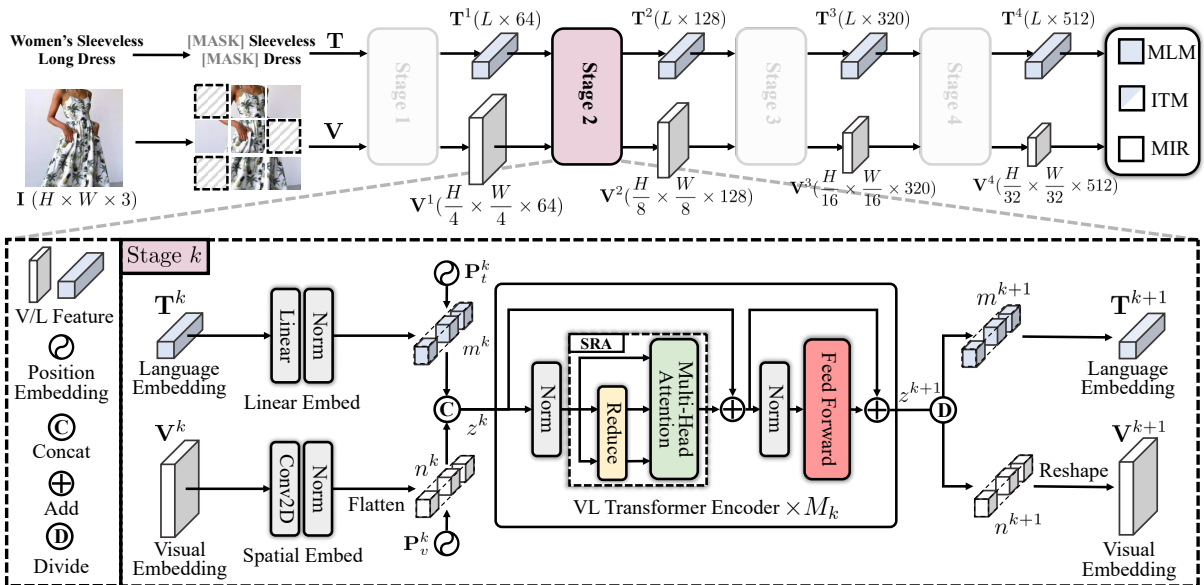


Fig. 3 Pipeline of our MVLT framework. Our overall architecture consists of four stages containing language and visual embeddings and multiple transformer encoders ($\times M_k$). Introducing the masking strategy for three sub-tasks, *i.e.*, masked image reconstruction (MIR), image-text matching (ITM), and masked language modeling (MLM), our MVLT can be trained in an end-to-end manner. More details can be found in Sec. 3.

which extends the pyramid vision transformer [21] to an architecture that adaptively extracts hierarchical representations for fashion cross-modal tasks. It is the first model that solves the end-to-end problem of VL pre-training in fashion, which allows us to simplify the implementation of our MVLT in the fashion industry using a twin-tower architecture [44].

3 Masked Vision-Language Transformer

Our goal is to build an end-to-end VL framework for the fashion domain. The overall pipeline of our MVLT is depicted in Fig. 3. Like PVT, our architecture inherits four stages’ properties and generates features with different sizes. Two keys of the proposed architecture are the multi-modal encoder (Sec. 3.1) and the pre-training objectives (Sec. 3.2).

3.1 Multi-Modal Encoder

As shown in Fig. 3, MVLT admits visual and verbal inputs. On the language side, we first tokenize the caption of a fashion product and use the specific token [MASK] to randomly mask out the

caption tokens with the masking ratio¹ r_l . Following the masking procedure, we obtain a sequence of word tokens. Then, we insert a specific [CLS] token at the head of this sequence. Besides, we pad the sequence to a unified length L using the [PAD] token if the length is shorter than 128. This procedure generates the language input ids $\mathbf{T} \in \mathbb{R}^L = \langle t_1; \dots; t_L \rangle$. On the vision side, we treat $\mathbf{I} \in \mathbb{R}^{H \times W \times 3}$ as visual input, where H and W denote the height and width of the given input. This input is sliced into multiple grid-like patches $\mathbf{V} \in \mathbb{R}^{N \times P \times P \times 3} = \langle v_1; \dots; v_N \rangle$, where $N = \frac{HW}{P^2}$ is the total number of patches and P denotes the patch size. Similarly, the split patches are masked out with mask ratio r_v . We provide more details about the above masking strategy for the language and vision parts in Sec. 3.2.

The above multi-modal inputs are embedded and fed into the consequent four VL interaction stages (*i.e.*, $k \in \{1, 2, 3, 4\}$). In the first stage, we generate the vision and language embeddings, \mathbf{T}^1 and \mathbf{V}^1 respectively, via the given inputs (\mathbf{T} and \mathbf{V}). Regarding the subsequent stages, we consider only the k -th stage, to have concise illustrations. As shown in the bottom part of Fig. 3, we first embed the language embedding $\mathbf{T}^k \in \mathbb{R}^{L \times D_k}$ into

¹We follow the default setting in BERT [12].

Table 1 Hyperparameter of our multi-modal encoders.

Hyperparameter	$k = 1$	$k = 2$	$k = 3$	$k = 4$
Layer number M_k	2	2	2	2
Hidden size D_k	64	128	320	512
Reduction size R_k	4	8	16	32
Kernel size K_k	4	2	2	2
Stride length S_k	4	2	2	2

the language hidden feature $m^k \in \mathbb{R}^{L \times D_{k+1}}$, which is formulated as:

$$m^k = \mathbf{T}^k * \mathbf{W}_t^k + \mathbf{P}_t^k, \quad (1)$$

where $\mathbf{W}_t^k \in \mathbb{R}^{D_k \times D_{k+1}}$ and $\mathbf{P}_t^k \in \mathbb{R}^{L \times D_{k+1}}$ are the learnable linear embedding and position embedding matrices. D_k is the size of the hidden feature embedding.

The visual embeddings are $\mathbf{V}^k \in \mathbb{R}^{\frac{H}{R_k} \times \frac{W}{R_k} \times D_k}$, where R_k denotes the spatial reduction factor of visual embedding. To acquire pyramid visual features, \mathbf{V}^k are then embedded and flattened into the visual hidden feature $n^k \in \mathbb{R}^{(HW/R_{k+1}^2) \times D_{k+1}}$ via a two-dimensional projection (*i.e.*, Conv2D block). In particular, this projection enforces the network to reduce the equivalent spatial dimension from \mathbb{R}^{HW/R_k^2} to $\mathbb{R}^{HW/R_{k+1}^2}$ by utilizing the convolutional kernel $\mathbf{W}_v^k \in \mathbb{R}^{D_k \times K_k \times K_k \times D_{k+1}}$ with kernel size K_k and stride length S_k . This could be formulated as:

$$n^k = \mathbf{Flatten}(\mathbf{V}^k * \mathbf{W}_v^k) + \mathbf{P}_v^k, \quad (2)$$

where $\mathbf{P}_v^k \in \mathbb{R}^{N \times D_{k+1}}$ denotes the position embedding matrix. We then concatenate these two VL hidden features $z^k = \langle m^k; n^k \rangle$ and feed them into multiple (M_k) VL transformer encoders. Each encoder contains the multi-head self-attention layer with spatial reduction (*i.e.*, reduce block), multi-layer perceptron, and layer normalization. Finally, we obtain the encoded multi-modal feature $z^{k+1} = \langle m^{k+1}; n^{k+1} \rangle$ and divide it into a language part $\mathbf{T}^{k+1} = m^{k+1}$ and a visual part $\mathbf{V}^{k+1} = \mathbf{Reshape}(n^{k+1})$, where the **Reshape**(\cdot) operation consists in recovering the spatial dimension of the given feature.

After four VL interaction stages, we generate the four text embeddings $\{\mathbf{T}^k\}_{k=1}^4$ and four pyramid vision embeddings $\{\mathbf{V}^k\}_{k=1}^4$, respectively. Table 1 presents more detailed hyperparameter settings of our method.

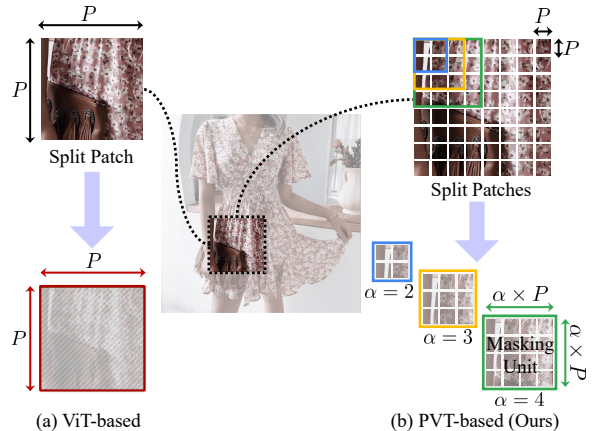


Fig. 4 PVT-based architectures offer more options for designing the masking strategy. The vanilla ViT-based method (a) [37] only selects a fixed-scale patch to mask, *i.e.*, P^2 . However, PVT-based method (b) is more versatile because it combines more fine-grained patches as a basic masking unit, *i.e.*, $(\alpha \times P)^2$, where $\alpha \in \{1, 2, \dots, 8\}$. These masked patches are not overlapped with each other. This characteristic provides a flexible way to learn the suitable semantics by using different values for α . Notably, we adopt a fixed scale factor of masking units in an individual experiment.

3.2 Pre-Training Objectives

To acquire discriminative multi-modal representations, we adopt three pre-training tasks to establish the inter-and intra-relationships between the most primitive VL modalities, including vision (masked image reconstruction, MIR), language (*i.e.*, masked language modeling, MLM), and VL (image-text matching, ITM) modalities.

Objective 1: Masked Image Reconstruction (MIR). As for the general domain, models are enough to learn the coarse-grained semantics from the patch- or region-based objectives and achieve satisfactory results. However, the fashion-specific models require more fine-grained representations, such as a suit with different materials (*e.g.*, wool) or collars (*e.g.*, Windsor), which needs a pixel-to-pixel vision pre-training objective. Inspired by the masked language modeling [12], we attempt to build the pixel-to-pixel relationships from the perspective of generative tasks, which promote the scalability of visual representations. We design the Masked Image Reconstruction (MIR) to accomplish this idea. To help our model learn better by MIR, we utilize the pyramid characteristic of PVT architecture [21] to design a flexible masking strategy. Unlike the ViT-based method (a)

in Fig. 4, PVT-based architecture (b) masks out the input image according to the masking unit matrix that contains small-grained patches. Given the patch sequence $\mathbf{V} = \{v_n\}_{n=1}^N \in \mathbb{R}^{N \times P \times P \times 3}$, the masked-out sequence $\mathbf{V}_{\setminus\Phi}$ is defined as:

$$\begin{aligned} \mathbf{V}_{\setminus\Phi} &= \mathcal{F}_M(\{\mathbf{M}(q; \alpha; \Phi)\}_q^Q, \{v_n\}_{n=1}^N) \\ &= \begin{cases} [\text{ZERO}], & \mathbf{M}(q; \alpha; \Phi) = 1, \\ v_n, & \mathbf{M}(q; \alpha; \Phi) = 0, \end{cases} \end{aligned} \quad (3)$$

where $\mathcal{F}_M(\cdot; \cdot)$ represents a function (or procedure) of our masking strategy, q is the randomly selected area of the masking unit, and [ZERO] means that we use a pixel value of zero² to fill the selected areas. The masking units $\{\mathbf{M}(q; \alpha; \Phi)\}_{q=1}^Q$ are derived from the indicator function:

$$\mathbf{M}(q; \alpha; \Phi) = \mathbf{1}(q) = \begin{cases} 1, & q \in \Phi, \\ 0, & q \notin \Phi, \end{cases} \quad (4)$$

where each value in a set of integers Φ is randomly selected from range $[1, Q]$ with ratio r_v . $Q = \frac{H \times W}{(\alpha \times P)^2}$ is the total number of masking units. For instance in Fig. 4 (b), we can define α from 1 to 8. In our default settings, we set $\alpha = 4$ to capture more fine-grained semantics³.

Since the smooth- ℓ_1 loss is less sensitive to the outliers, we use it as the pre-training objective to reconstruct the whole image via the masked-out sequence $\mathbf{V}_{\setminus\Phi}$. It is defined as:

$$\mathcal{L}_{\text{MIR}} = \begin{cases} 0.5 \times (\mathbf{I}'_{(x,y)} - \mathbf{I}_{(x,y)})^2, & \text{if } \mathbf{I}'_{(x,y)} - \mathbf{I}_{(x,y)} < 1, \\ |\mathbf{I}'_{(x,y)} - \mathbf{I}_{(x,y)}| - 0.5, & \text{otherwise,} \end{cases} \quad (5)$$

where $\mathbf{I}'_{(x,y)}$ and $\mathbf{I}_{(x,y)}$ denote the pixel at coordinate (x, y) in the reconstructed image \mathbf{I}' and the input image \mathbf{I} , respectively. $\mathbf{I}' = \mathcal{F}_{\text{MIR}}(\mathbf{V}_{\setminus\Phi}; \mathbf{W}_{\text{MIR}})$ is parameterized by learnable weights \mathbf{W}_{MIR} . Function $\mathcal{F}_{\text{MIR}}(\cdot; \mathbf{W}_{\text{MIR}})$ denotes a standard four-level U-Net [45] decoder, which admits four pyramidal vision embeddings $\{\mathbf{V}^k\}_{k=1}^4$ as inputs.

²In fact, we set [ZERO] = 10^{-6} to bring better optimization stability and less pattern degradation.

³The vanilla masking strategy in Fig. 4 (a) with $P = 32$ becomes a special case of our masking strategy in Fig. 4 (b) when $\alpha = 8, P = 4$.

Objective 2: Image-Text Matching (ITM).

The appended classification embedding in the last language embedding \mathbf{T}^4 is used to couple the representations from VL modalities. We utilize the function $\mathcal{F}_{\text{ITM}}(\cdot; \mathbf{W}_{\text{ITM}})$ to denote a full-connected (FC) and softmax layers, parameterized by the weights \mathbf{W}_{ITM} . \mathcal{F}_{ITM} outputs a two-class probability vector $\mathbf{p}_{\text{ITM}} = \mathcal{F}_{\text{ITM}}(\langle \mathbf{T}, \mathbf{V} \rangle; \mathbf{W}_{\text{ITM}})$, representing whether the input fashion image and caption match (*i.e.*, positive pair) or not (*i.e.*, negative pair). The positive pairs are selected from the same fashion product category, whereas the negative pairs are chosen at random from different entries. The binary cross-entropy loss function finally constrains this task:

$$\begin{aligned} \mathcal{L}_{\text{ITM}} &= -\mathbb{E}_{\langle \mathbf{T}, \mathbf{V} \rangle} [\mathbf{y}_{\text{ITM}} \log(\mathbf{p}_{\text{ITM}}) \\ &\quad + (1 - \mathbf{y}_{\text{ITM}}) \log(1 - \mathbf{p}_{\text{ITM}})], \end{aligned} \quad (6)$$

where \mathbf{y}_{ITM} denotes the ground-truth label, *i.e.*, 1 for matched pairs and 0 for unmatched pairs.

Objective 3: Masked Language Modeling (MLM).

Following [46], we randomly use the specific token [MASK] to replace the original text tokens. The target of the MLM is to predict the text content for the masked tokens using the unmasked tokens and patches. Given a tokenized sequence $\mathbf{T} = \{t_1, \dots, t_L\}$, the masked-out sequence is denoted by $\mathbf{T}_{\setminus i} = \{t_1, \dots, [\text{MASK}]_i, \dots, t_L\}$. We use the cross-entropy loss to model this objective:

$$\mathcal{L}_{\text{MLM}} = -\mathbb{E}_{\mathbf{T}} [\log(\mathbf{p}_{\text{MLM}})], \quad (7)$$

where $\mathbf{p}_{\text{MLM}} = \mathcal{F}_{\text{MLM}}(\mathbf{T}_{\setminus i}; \mathbf{W}_{\text{MLM}})$ denotes the predicted probability for each masked-out token [MASK]_{*i*} using $\mathbf{T}_{\setminus i}$. The function $\mathcal{F}_{\text{MLM}}(\cdot; \mathbf{W}_{\text{MLM}})$ represents the parameters \mathbf{W}_{MLM} of a classifier. The final pre-training objective of the proposed MVLT is a combination of the three objectives:

$$\mathcal{L}_{\text{total}} = w_1 \times \mathcal{L}_{\text{MIR}} + w_2 \times \mathcal{L}_{\text{ITM}} + w_3 \times \mathcal{L}_{\text{MLM}}. \quad (8)$$

3.3 Downstream Tasks

For a fair comparison, we follow the same training/inference protocols as in [1, 2] and also adopt the Fashion-Gen 2018 [47] benchmark as the base of our experiments. This dataset contains 67,666

Table 2 Retrieval (*i.e.*, TIR and ITR) and recognition (*i.e.*, M-CR and S-CR) performances on Fashion-Gen dataset. \uparrow means the larger, the better. Here, $\text{Sum}\mathcal{R}=(\mathcal{R}@1+\mathcal{R}@5+\mathcal{R}@10) \times 100$ and $\text{Sum}\mathcal{C}=(\mathcal{A} + \text{macro-}\mathcal{F}) \times 100$. “N/A” means the score is not available. “Diff” means the numerical difference between the performance of the second-ranked competitor and our MVLT.

Task	Metric	VSE	VSE++	SCAN	PFAN	ViLBERT	ImageBERT	FashionBERT	VL-BERT	OSCAR	Kaleido-BERT	MVLT	
		arXiv ₁₄	BMVC ₁₈	ECCV ₁₈	arXiv ₁₉	NeurIPS ₁₉	arXiv ₂₀	SIGIR ₂₀	ICLR ₂₀	ECCV ₂₀	CVPR ₂₁	OUR ₂₂	Diff
TIR	$\mathcal{R}@1$	\uparrow 4.350%	\uparrow 4.600%	\uparrow 4.300%	\uparrow 6.200%	\uparrow 21.12%	\uparrow 24.78%	\uparrow 26.75%	\uparrow 22.63%	\uparrow 25.10%	\uparrow <u>33.88%</u>	34.60%	+0.72%
	$\mathcal{R}@5$	\uparrow 12.76%	\uparrow 16.89%	\uparrow 13.00%	\uparrow 20.79%	\uparrow 37.23%	\uparrow 45.20%	\uparrow 46.48%	\uparrow 36.48%	\uparrow 49.14%	\uparrow <u>60.60%</u>	78.00%	+17.40%
	$\mathcal{R}@10$	\uparrow 20.91%	\uparrow 28.99%	\uparrow 22.30%	\uparrow 31.52%	\uparrow 50.11%	\uparrow 55.90%	\uparrow 55.74%	\uparrow 48.52%	\uparrow 56.68%	\uparrow <u>68.59%</u>	89.50%	+20.91%
	Sum \mathcal{R}	\uparrow 38.02	\uparrow 50.48	\uparrow 39.6	\uparrow 58.51	\uparrow 108.46	\uparrow 125.88	\uparrow 128.97	\uparrow 128.97	\uparrow 107.63	\uparrow 130.92	\uparrow <u>163.07</u>	202.1
ITR	$\mathcal{R}@1$	\uparrow 4.010%	\uparrow 4.590%	\uparrow 4.590%	\uparrow 4.290%	\uparrow 20.97%	\uparrow 22.76%	\uparrow 23.96%	\uparrow 19.26%	\uparrow 23.39%	\uparrow <u>27.99%</u>	33.10%	+5.11%
	$\mathcal{R}@5$	\uparrow 11.03%	\uparrow 14.99%	\uparrow 16.50%	\uparrow 14.90%	\uparrow 40.49%	\uparrow 41.89%	\uparrow 46.31%	\uparrow 39.90%	\uparrow 44.67%	\uparrow <u>60.09%</u>	77.20%	+17.11%
	$\mathcal{R}@10$	\uparrow 22.14%	\uparrow 24.10%	\uparrow 26.60%	\uparrow 24.20%	\uparrow 48.21%	\uparrow 50.77%	\uparrow 52.12%	\uparrow 46.05%	\uparrow 52.55%	\uparrow <u>68.37%</u>	91.10%	+22.73%
	Sum \mathcal{R}	\uparrow 37.18	\uparrow 43.68	\uparrow 47.69	\uparrow 43.39	\uparrow 109.67	\uparrow 115.42	\uparrow 122.39	\uparrow 105.21	\uparrow 120.61	\uparrow <u>156.45</u>	201.4	+44.95
M-CR	\mathcal{A}	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow 90.77%	\uparrow 91.25%	\uparrow N/A	\uparrow 91.79%	\uparrow <u>95.07%</u>	98.26%	+3.19%
	macro- \mathcal{F}	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow 0.699	\uparrow 0.705	\uparrow N/A	\uparrow <u>0.727</u>	\uparrow 0.714	0.896	+0.169
	Sum \mathcal{C}	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow 160.67	\uparrow 161.75	\uparrow N/A	\uparrow 164.49	\uparrow <u>166.47</u>	187.86	+21.39
S-CR	\mathcal{A}	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow 80.11%	\uparrow 85.27%	\uparrow N/A	\uparrow 84.23%	\uparrow <u>88.07%</u>	93.57%	+5.50%
	macro- \mathcal{F}	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow 0.575	\uparrow 0.620	\uparrow N/A	\uparrow 0.591	\uparrow <u>0.636</u>	0.829	+0.193
	Sum \mathcal{C}	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow N/A	\uparrow 137.61	\uparrow 147.27	\uparrow N/A	\uparrow 143.33	\uparrow <u>151.67</u>	176.47	+24.80

fashion products (*i.e.*, 60,147 entries for training and 7,519 entries for testing) and their associated product descriptions. Each product corresponds to an image set (including 1 ~ 6 samples) at various viewing angles. As a result, we utilize 260,480 and 35,528 image-text pairs as training and testing partitions, respectively. For a fair comparison, we test MVLT and compared models on Fashion-Gen using the following four fashion-related VL downstream tasks.

Task 1: Text-Image Retrieval (TIR). The TIR task requires the model to find a text with the highest similarity value with different query images. In particular, we take a product title and its corresponding image as a positive image-text pair, while the negative pairs are randomly selected from a pool of mismatched images. To increase our experiment’s difficulty, we constrain a set of image-text candidates (*i.e.*, a positive pair and 100 negative pairs) in the same sub-category, making them as similar as possible.

Task 2: Image-Text Retrieval (ITR). As the reverse process of the TIR task, the ITR task aims to retrieve a matching image given a sequence of text entries of fashion description, where these bidirectional retrieval tasks (*i.e.*, TIR and ITR) become a prominent member of cross-modal research. Similar to the above selection strategy in the TIR, we prepare a set of candidate image-text pairs, including a positive pair and 100 negative pairs from the same sub-category. We evaluate the zero-shot learning ability of our MVLT without further fine-tuning for these two retrieval tasks. We utilize three accuracy metrics

(*i.e.*, $\mathcal{R}@1$, $\mathcal{R}@5$, and $\mathcal{R}@10$) for the evaluation by ranking a series of predicted probabilities.

Task 3: Category Recognition (M-CR and S-CR). This task has two parts: main-category recognition (M-CR) and sub-category recognition (S-CR). These tasks act as the fundamental role of practical e-commerce applications that offer the specific category of the queried product. We expect that the model should possess the ability to recognize differences under different granularity levels: 48 main-categories and 122 sub-categories, such as {M-CR = SWEATERS, S-CR = CREWNECKS}. After the class embedding in the last language embedding \mathbf{T}^4 , we add two independent FC layers to generate the final probabilities for two different recognition tasks. This procedure requires additional fine-tuning with recognition labels. We utilize two recognition-related metrics to evaluate performance: accuracy (\mathcal{A}) and macro F-measure (macro- \mathcal{F}).

Task 4: Masked Image Generation (MIG). MIG task can be viewed as a pixel-wise reconstruction task. Each patch in the image is randomly masked with the probability r_v (refer to the pre-training task MIR in Sec. 3.2). Then, we ask the model to recreate the whole image using the uncovered areas as visual clues.

4 Experiments

This section will detail our experiment to determine the factors leading to the success of the proposed MVLT.

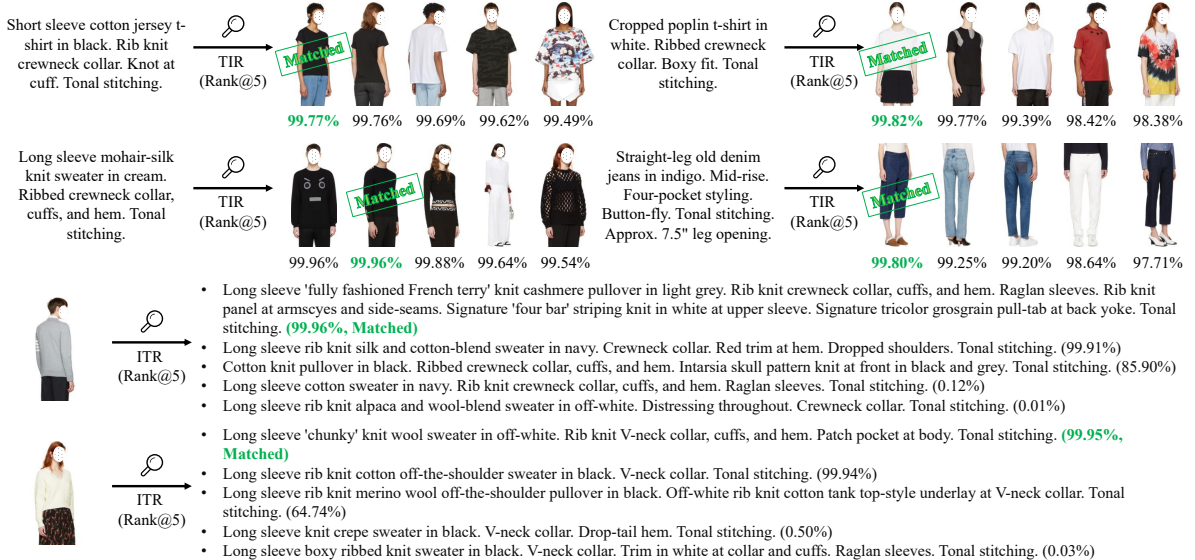


Fig. 5 Visualization results on the TIR and ITR tasks in terms of top-five ranked probabilities predicted by our MVLT. “Matched” indicates the ground-truth image-text pair.

4.1 Settings

This part provides the hyperparameter settings for our training procedure: **i) Pre-training.** We utilize PyTorch to implement our method, which is accelerated by 8 Tesla V100 GPUs. We adopt AdamW optimizer with a momentum value of 0.9, a mini-batch size of 1200 (*i.e.*, 150 per GPU), a weight decay of 10^{-4} . To avoid over-fitting, we initialize MVLT on ImageNet pre-trained weights [21]. The learning rate is initially set to 2.5×10^{-3} and is changed using a cosine learning schedule. For the visual side, the input image is resized to $H=W=256$ and split into the multiple sub-patches with a size of $P = 4$. For the language side, all the product captions are tokenized and padded to tokens with a unified length of $L = 128$, including classification, caption, and padding tokens. The mask probabilities for vision and language are set to $r_v = 0.5$ and $r_l = 0.15$, respectively. We empirically set weighting factors $\{w_1 = 10, w_2 = 1, w_3 = 1\}$ to balance the orders of magnitude of different loss values. **ii) Fine-tuning.** We transfer the pre-trained VL representation to each downstream application via fine-tuning in an end-to-end manner, whose settings are consistent with the pre-training process.

4.2 Results

As described in Sec. 3.3, we provide the details of four downstream fashion-related tasks. Experimental results show that our MVLT outperforms all competitors, including VSE [48], VSE++ [49], SCAN [26], PFAN [50], ViLBERT [16], ImageBERT [15], FashionBERT [1], VL-BERT [25], OSCAR [29], and Kaleido-BERT [2], which demonstrate the superiority for handling the VL understanding and generation tasks.

TIR and ITR. As shown in Table 2, our MVLT surpass the best method (*i.e.*, Kaleido-BERT-CVPR₂₁) on the TIR task by margins of **+17.40%**, **+20.91%** across the $\mathcal{R}@5$, $\mathcal{R}@10$. As for ITR, our method delivers more competitive results, with improvements of **+17.11%**, **+22.73%** on the $\mathcal{R}@5$, $\mathcal{R}@10$ metrics, respectively. In any case, these results strongly support that our model is powerful enough to match vision and language. They also show how **a) MIR** and **b) end-to-end pre-training** are useful in fashion. We believe that MVLT would set a precedent in many industrial applications because it is a simple, cost-effective, and powerful architecture. Besides, we present the visualization results of these two retrieval tasks in Fig. 5.

M-CR and S-CR. Compared with BERT-based architectures [1, 2, 15, 29], we also achieve top-1 performances in these two tasks, demonstrating

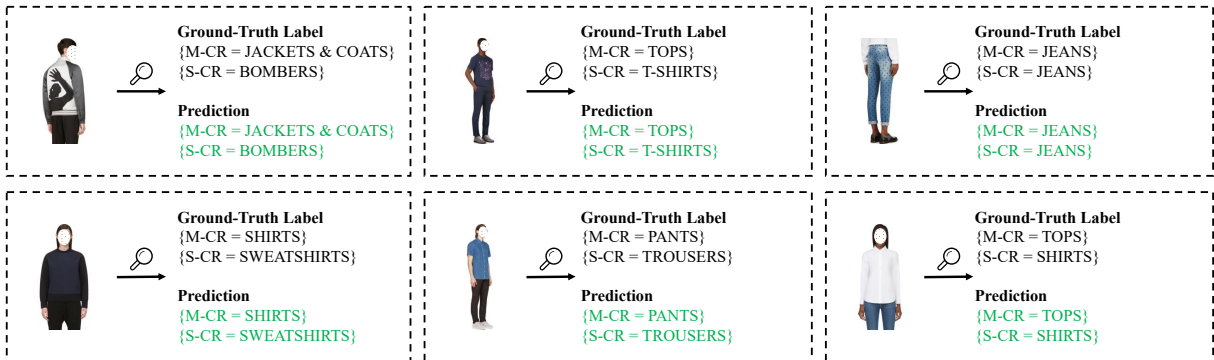


Fig. 6 The visualization of main-/sub-category recognition results on Fashion-Gen. The green predictions hit the targets.



Fig. 7 Visualization of samples generated by our MVLT. The gray blocks represent the masked regions.

our method have an excellent VL understanding capability. Moreover, compared with the best method Kaleido-BERT, our architecture improves by **0.193** in macro- \mathcal{F} metric for the S-CR task. In addition, the mean improvements in terms of the SumC metric (*i.e.*, M-CR: **+21.39** and S-CR: **+24.80**) are very significant. Since this metric is very sensitive to data distribution, it demonstrates MVLT has super-strong robustness. We also present the recognition results of M-CR and S-CR in Fig. 6.

MIG. As shown in Fig. 7, we showcase reconstructed images on the validation part of Fashion-Gen 2018 (a) and our e-commercial website (b). As seen, the reconstruction performance is truly remarkable. Since it requires our method to learn the fashion semantics truly, such results demonstrate the generative ability of our approach.

4.3 Ablation Studies

Mask Ratio. Table 3 (a) present four variants for different mask probability r_v (*i.e.*, 0.10 (A1), 0.30 (A2), 0.70 (A3), 0.90 (A4)) and our choice: 0.50 (**Final**). The $\mathcal{R}@5$ rises steadily with the masking probability until it reaches the sweet spot (75.70% \rightarrow 78.00%); then it reach

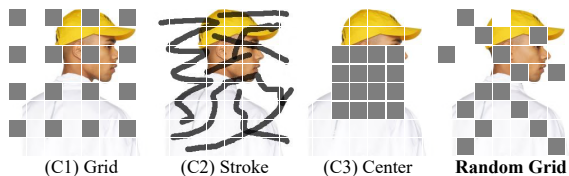


Fig. 8 We designed four strategies to mask fashion images. The random grid performs the best.

performance plummets (73.80%). We argue that increasing the r_v will make MIR more complex, allowing MVLT to learn better semantics in a more restricted situation. However, masking out too much region will naturally result in losing valid visual information, leading to bad results.

Masked Unit Size. Thanks to PVT’s flexibility, we can easily try different sizes of masked patches. As shown in Table 3 (b), we derive four variants with masked unit size α (*i.e.*, 1 (B1), 2 (B2), 8 (B3), 16 (B4)) to compare with our setting: 4 (**Final**). We found the performance is sensitive to this factor. It makes sense, revealing how vital it is to learn a robust fashion-related representation with a moderate granularity.

Masking Style. As shown in Fig. 8, we designed four types of masking strategies for the MIR task, whose quantitative differences are presented in

Table 3 Ablation studies of five key pre-training factors on our MVLT. More relevant analyses refer to Sec. 4.3.

App.	Metric	(a) Mask Ratio (r_v)				(b) Masking Unit Size (α)				(c) Masking Style			(d) Pre-Training Tasks			(e) Pre-Train	MVLT (Final)
		(A1)	(A2)	(A3)	(A4)	(B1)	(B2)	(B3)	(B4)	(C1)	(C2)	(C3)	(D1)	(D2)	(D3)	(E1)	
		0.10	0.30	0.70	0.90	1	2	8	16	Grid	Stroke	Center	ITM	ITM+MIR	ITM+MLM	w/o PVT	
TIR	$\mathcal{R}@1$	31.10%	33.50%	30.50%	30.70%	31.90%	30.30%	30.00%	32.20%	32.20%	31.40%	30.40%	30.40%	32.20%	32.90%	29.00%	34.60%
	$\mathcal{R}@5$	75.70%	76.00%	75.50%	73.80%	75.30%	75.60%	73.90%	76.90%	75.30%	76.10%	75.10%	74.10%	76.00%	76.20%	72.20%	78.00%
	$\mathcal{R}@10$	88.60%	88.70%	88.00%	88.60%	89.60%	88.60%	88.20%	88.60%	88.50%	89.20%	87.20%	83.50%	87.20%	88.60%	86.60%	89.50%
	Sum \mathcal{R}	195.40	198.20	194.00	193.10	196.80	194.50	192.10	197.70	196.00	196.70	192.70	188.00	195.40	197.70	187.80	202.10
	Diff	-6.70	-3.90	-8.10	-9.00	-5.30	-7.60	-10.00	-4.40	-6.10	-5.40	-9.40	-14.10	-6.70	-4.40	-14.30	-
ITR	$\mathcal{R}@1$	30.00%	29.90%	29.90%	28.50%	29.00%	29.70%	29.00%	28.90%	31.40%	31.10%	30.10%	29.30%	30.40%	28.40%	25.60%	33.10%
	$\mathcal{R}@5$	75.70%	74.90%	76.50%	75.00%	76.90%	77.10%	74.20%	77.30%	77.40%	74.50%	73.90%	70.80%	75.50%	76.30%	71.50%	77.20%
	$\mathcal{R}@10$	88.80%	89.00%	89.20%	88.20%	89.40%	87.70%	88.00%	89.90%	89.60%	88.50%	87.80%	86.80%	87.80%	88.80%	85.90%	91.10%
	Sum \mathcal{R}	194.50	193.80	195.60	191.70	195.30	194.50	191.20	196.10	198.40	194.10	191.80	186.90	193.70	193.50	183.00	201.40
	Diff	-6.90	-7.60	-5.80	-9.70	-6.10	-6.90	-10.20	-5.30	-3.00	-7.30	-9.60	-14.50	-7.70	-7.90	-18.40	-
M-CR	\mathcal{A}	98.16%	97.87%	98.09%	98.06%	98.03%	98.04%	98.11%	98.01%	98.12%	98.07%	98.04%	96.49%	97.11%	98.08%	97.92%	98.26%
	macro- \mathcal{F}	0.870	0.860	0.890	0.870	0.870	0.880	0.850	0.870	0.869	0.877	0.870	0.806	0.853	0.876	0.879	0.896
	Sum \mathcal{C}	185.16	183.87	187.09	185.06	185.03	186.04	183.11	185.01	185.02	185.77	185.04	177.09	182.41	185.68	185.82	187.86
	Diff	-2.70	-3.99	-0.77	-2.80	-2.83	-1.82	-4.75	-2.85	-2.84	-2.09	-2.82	-10.77	-5.45	-2.18	-2.04	-
S-CR	\mathcal{A}	93.10%	93.34%	93.36%	93.23%	93.29%	93.34%	93.32%	93.32%	93.37%	93.21%	93.59%	89.64%	90.87%	93.29%	92.90%	93.57%
	macro- \mathcal{F}	0.800	0.810	0.820	0.810	0.810	0.810	0.800	0.799	0.794	0.814	0.830	0.703	0.728	0.809	0.790	0.829
	Sum \mathcal{C}	173.10	174.34	175.36	174.23	174.29	174.34	173.32	173.22	172.77	174.61	176.59	159.94	163.67	174.19	171.90	176.47
	Diff	-3.37	-2.13	-1.11	-2.24	-2.18	-2.13	-3.15	-3.25	-3.70	-1.86	+0.12	-16.53	-12.80	-2.28	-4.57	-

Table 3 (c), *i.e.*, grid (C1), stroke (C2), center (C3) and our random grid (**Final**) masking strategies. As can be seen, the random grid masking (Final) yields the best results, while the other three perform poorly. We believe this is because, in comparison to the grid (C1) and center (C3), random grid masking (Final) can help MVLT construct comprehensive representations. As our strategy (Final) does, the stroke (C2) also randomly masks the image given, yet it more or less leaves unmasked visual cues in the sub-patches. Our strategy enables the model to easily predict the masked region because semantics in the image are well preserved, enhancing the model’s robustness to learning in-sight knowledge. **Pre-Training Objectives.** As shown in Table 3 (d), we derive four different variants to investigate the contribution of each objective, including ITM (D1), ITM+MIR (D2), ITM+MLM (D3), and our ITM+MIR+MLM (**Final**). When comparing D3 to D1 and D2 in the TIR task, we can see that D3 has a better performance in $\mathcal{R}@5$ metric: 74.10% (D1) < 76.00% (D2) < 76.20% (D3). We conclude MLM task can help the model thoroughly learn the language knowledge, so it provides a more precise query to recall better-matching images. In the ITR task, we find the similar conclusion when comparing (D2) to (D1) and D3 in $\mathcal{R}@5$ metric: 70.80% (D1) < 75.50% (D2) < 76.30% (D3). It indicates that better visual learning leads to an accurate image query to match the most appropriate caption.

Table 4 Ablation study for the contribution of loading PVT’s weights pre-trained on ImageNet [51].

	TIR		ITR		M-CR		S-CR	
	$\mathcal{R}@5$	$\mathcal{R}@10$	$\mathcal{R}@5$	$\mathcal{R}@10$	\mathcal{A}	macro- \mathcal{F}	\mathcal{A}	macro- \mathcal{F}
w/o PVT	72.20%	86.60%	71.50%	85.90%	97.92%	0.879	92.90%	0.790
w/ PVT	78.00%	89.50%	77.20%	91.10%	98.26%	0.896	93.57%	0.829
Diff	+5.80%	+2.90%	+5.70%	+5.20%	+0.34%	+1.7%	+0.67%	+3.9%

Loading Pre-Trained Weight. As seen in Table 4, we add an experiment to demonstrate it is very important to load the PVT’s weight pre-trained on ImageNet [51]. If not, it is obvious that our MVLT will suffer fierce drops (*i.e.*, ITR: 77.20% \rightarrow 71.50% in $\mathcal{R}@5$, S-CR: 93.57% \rightarrow 92.90% in \mathcal{A}). It is reasonable because a method pre-trained on large-scale general datasets can be more applicable in a specific field. It has already learned information such as color, texture, shape, *etc.*

4.4 More Discussions

How does MVLT perform in general domains? To further investigate the potential abilities in general settings, we here discuss two extended questions. *a) Can the general models be directly transferred to the fashion domain?* Inspired by the huge impact of general vision-language models, as in Table 5, we further investigate the zero-shot performance of two typical general models (*i.e.*, ViLBERT [16] and CLIP [52]). This has once again demonstrated the necessity and superiority of MVLT pre-trained on the specific domains. *b) Can MVLT also work*

Table 5 The comparison of zero-shot retrieval results on the Fashion-Gen dataset.

	TIR			ITR		
	$\mathcal{R}@1\uparrow$	$\mathcal{R}@5\uparrow$	$\mathcal{R}@10\uparrow$	$\mathcal{R}@1\uparrow$	$\mathcal{R}@5\uparrow$	$\mathcal{R}@10\uparrow$
ViLBERT (Zero-shot)	7.18%	18.73%	29.84%	8.99%	15.34%	26.14%
CLIP (Zero-shot)	16.30%	40.60%	55.60%	13.60%	43.10%	57.60%
MVLT (OUR)	34.60%	78.00%	89.50%	33.10%	77.20%	91.10%

Table 6 Retrieval results on the MS-COCO 2014 dataset. \dagger means using an extra feature extractor (*e.g.*, Faster RCNN).

	TIR task (5K Test)			ITR task (5K Test)		
	$\mathcal{R}@1\uparrow$	$\mathcal{R}@5\uparrow$	$\mathcal{R}@10\uparrow$	$\mathcal{R}@1\uparrow$	$\mathcal{R}@5\uparrow$	$\mathcal{R}@10\uparrow$
Unicoder-VL \dagger	48.40%	76.70%	85.90%	62.30%	87.10%	92.80%
UNITER-Base \dagger	50.30%	78.50%	87.20%	64.40%	87.40%	93.10%
ViLT-Base/32	41.30%	72.00%	82.50%	61.80%	86.20%	92.60%
MVLT (OUR)	49.66%	79.88%	87.50%	65.38%	90.04%	93.60%

well in the general domain? We further verify the potential ability of our MVLT on the general domain. Table 6 reports the performance on MS-COCO 2014 dataset [53], where MVLT follows the same training standards as in [37]. It shows that MVLT achieves promising results compared to the latest models (*i.e.*, Unicoder-VL [54], UNITER [17], and ViLT [37]) without extra training data and special retrieval losses during the training. It indicates that MVLT is also a promising solution when extended to general scenes.

Why do pyramid architecture and MIR benefit? As mentioned in the introduction, there are two understudied problems in the fashion domain. *To solve the transferability problem*, pyramidal architecture [21] takes raw data as input without complex pre-processing, which essentially alleviates the applied burden in industry. Besides, MIR does not need human annotations like classification tags, bounding boxes, or pixel-wise segmentation labels. *For the granularity problem* [55], the pyramidal architecture [21] provides multi-scale features with rich semantics. Combined with the MIR task, our framework can represent multi-grained fashion knowledge (*e.g.*, dress, V-neck). These features are helpful and urgently required in this field.

A VL model that performs well for semantic understanding tasks (*e.g.*, retrieval [56], classification) can serve as a good foundation and be easily applied to downstream tasks (*e.g.*, text-to-image synthesis [57], image captioning) by utilizing an additional decoder. We didn’t conduct image captioning experiments because we focused on basic representation learning in fashion this time.

MVLT *v.s.* MAE [34]. MAE learns general representations by allowing the model to explore pixel-to-pixel associations. So MVLT and MAE are similar in this regard. However, our MVLT is the first that introduces the vision reconstruction-like pre-training for multi-modal research (*e.g.*, fashion domain).

5 Conclusion

We present a vision-language framework named MVLT, which provides two contributions in this field: 1) a newly-designed masked image reconstruction (MIR) objective, and 2) an end-to-end pre-training scheme. The experimental and ablation analysis demonstrates the superiority of various matching and generative tasks. MVLT outperforms the cutting-edge method Kaleido-BERT with large margins on retrieval and recognition tasks, which would catalyze the fashion domain. The designed out-of-box method working end-to-end could simplify the workflow (*e.g.*, data pre-processing and model training) for the actual engineering value, which improves development and business efficiency on large-scale e-commerce websites by approximately 50%.

In the future, we will continue to investigate an extremely efficient method in this field using famous technologies like hashing [58], network pruning, and knowledge distil to alleviate the storage and computing limitations in real-world e-commerce applications.

Acknowledgments

This work is funded by Toyota Motor Europe via the research project TRACE-Zürich. The authors also would like to thank the anonymous reviewers and editor for their helpful comments on this manuscript.

Conflicts of Interests

The authors declared that they have no conflicts of interest to this work. We declare that we do not have any commercial or associative interest that represents a conflict of interest in connection with the work submitted.

References

- [1] D. Gao, L. Jin, B. Chen, M. Qiu, P. Li, Y. Wei, Y. Hu, and H. Wang, “Fashionbert: Text and image matching with adaptive loss for cross-modal retrieval,” in *International ACM SIGIR Conference on Research and Development in Information Retrieval*. New York, NY, USA: Association for Computing Machinery, 2020, pp. 2251–2260, DOI: [10.1145/3397271.3401430](https://doi.org/10.1145/3397271.3401430).
- [2] M. Zhuge, D. Gao, D.-P. Fan, L. Jin, B. Chen, H. Zhou, M. Qiu, and L. Shao, “Kaleidobert: Vision-language pre-training on fashion domain,” in *Conference on computer vision and pattern recognition*. Nashville, TN, USA: IEEE, 2021, pp. 12 642–12 652, DOI: [10.1109/CVPR46437.2021.01246](https://doi.org/10.1109/CVPR46437.2021.01246).
- [3] A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16-16 words: Transformers for image recognition at scale,” in *International Conference on Learning Representations*. [Online]: PMLR, 2021.
- [4] Z. Liu, Y. Lin, Y. Cao, H. Hu, Y. Wei, Z. Zhang, S. Lin, and B. Guo, “Swin transformer: Hierarchical vision transformer using shifted windows,” in *International conference on computer vision*. Montreal, Canada: IEEE, 2021, pp. 9992–10 002, DOI: [10.1109/ICCV48922.2021.00986](https://doi.org/10.1109/ICCV48922.2021.00986).
- [5] A. Vaswani, N. Shazeer, N. Parmar, J. Uszkoreit, L. Jones, A. N. Gomez, L. Kaiser, and I. Polosukhin, “Attention is all you need,” in *Advances in neural information processing systems*, vol. 30. Long Beach Convention Center, Long Beach, US: Curran Associates, Inc., 2017.
- [6] T.-X. Sun, X.-Y. Liu, X.-P. Qiu, and X.-J. Huang, “Paradigm shift in natural language processing,” *Machine Intelligence Research*, vol. 19, no. 3, pp. 169–183, 2022, DOI: [10.1007/s11633-022-1331-6](https://doi.org/10.1007/s11633-022-1331-6).
- [7] S. Agarwal, G. Krueger, J. Clark, A. Radford, J. W. Kim, and M. Brundage, “Evaluating clip: towards characterization of broader capabilities and downstream implications,” [Online], 2021, Available: <https://arxiv.org/abs/2108.02818>.
- [8] M. Chen, A. Radford, R. Child, J. Wu, H. Jun, D. Luan, and I. Sutskever, “Generative pretraining from pixels,” in *International conference on machine learning*. PMLR, 2020, pp. 1691–1703.
- [9] J. Lin, R. Men, A. Yang, C. Zhou, M. Ding, Y. Zhang, P. Wang, A. Wang, L. Jiang, X. Jia *et al.*, “M6: A chinese multimodal pretrainer,” in *Proceedings of the 27th ACM SIGKDD Conference on Knowledge Discovery & Data Mining*. Association for Computing Machinery, 2021, p. 3251–3261, DOI: [10.1145/3447548.3467206](https://doi.org/10.1145/3447548.3467206).
- [10] A. Ramesh, M. Pavlov, G. Goh, S. Gray, C. Voss, A. Radford, M. Chen, and I. Sutskever, “Zero-shot text-to-image generation,” in *International conference on machine learning*. PMLR, 2021, pp. 8821–8831.
- [11] H. Wu, Y. Gao, X. Guo, Z. Al-Halah, S. Rennie, K. Grauman, and R. Feris, “Fashion iq: A new dataset towards retrieving images by natural language feedback,” in *Conference on computer vision and pattern recognition*. Nashville, TN, USA: IEEE, 2021, pp. 11 307–11 317, DOI: [10.1109/CVPR46437.2021.01115](https://doi.org/10.1109/CVPR46437.2021.01115).
- [12] J. Devlin, M.-W. Chang, K. Lee, and K. Toutanova, “Bert: Pre-training of deep bidirectional transformers for language understanding,” in *Annual Conference of the North American Chapter of the Association for Computational Linguistics*. Minneapolis, Minnesota: Association for Computational Linguistics, 2019, pp. 4171–4186, DOI: [10.18653/v1/N19-1423](https://doi.org/10.18653/v1/N19-1423).
- [13] K. He, X. Zhang, S. Ren, and J. Sun, “Deep residual learning for image recognition,” in *Conference on computer vision and pattern recognition*. Las Vegas, NV, USA: IEEE,

- 2016, pp. 770–778, DOI: [10.1109/CVPR.2016.90](https://doi.org/10.1109/CVPR.2016.90).
- [14] S. Ren, K. He, R. Girshick, and J. Sun, “Faster r-cnn: Towards real-time object detection with region proposal networks,” in *Advances in neural information processing systems*, vol. 28. Montreal, Quebec, Canada: Curran Associates, Inc., 2015.
- [15] D. Qi, L. Su, J. Song, E. Cui, T. Bharti, and A. Sacheti, “Imagebert: Cross-modal pre-training with large-scale weak-supervised image-text data,” [Online], 2020, Available: <https://arxiv.org/abs/2001.07966>.
- [16] J. Lu, D. Batra, D. Parikh, and S. Lee, “Vilbert: Pretraining task-agnostic visiolinguistic representations for vision-and-language tasks,” in *Advances in neural information processing systems*, vol. 32. Vancouver, Canada: Curran Associates, Inc., 2019.
- [17] Y.-C. Chen, L. Li, L. Yu, A. El Kholy, F. Ahmed, Z. Gan, Y. Cheng, and J. Liu, “Uniter: Universal image-text representation learning,” in *European conference on computer vision*. Glasgow, United Kingdom: Springer, 2020, pp. 104–120, DOI: [10.1007/978-3-030-58577-8_7](https://doi.org/10.1007/978-3-030-58577-8_7).
- [18] W.-L. Hsiao, I. Katsman, C.-Y. Wu, D. Parikh, and K. Grauman, “Fashion++: Minimal edits for outfit improvement,” in *International conference on computer vision*. Montreal, Canada: IEEE, 2019, pp. 5046–5055, DOI: [10.1109/ICCV.2019.00515](https://doi.org/10.1109/ICCV.2019.00515).
- [19] M. I. Vasileva, B. A. Plummer, K. Dusad, S. Rajpal, R. Kumar, and D. Forsyth, “Learning type-aware embeddings for fashion compatibility,” in *European conference on computer vision*. Munich, Germany: Springer, 2018, pp. 405–421, DOI: [10.1007/978-3-030-01270-0_24](https://doi.org/10.1007/978-3-030-01270-0_24).
- [20] D.-P. Fan, M. Zhuge, and L. Shao, “Domain specific pre-training of cross modality transformer model,” 2022, uS Patent App. 17/186,745.
- [21] W. Wang, E. Xie, X. Li, D.-P. Fan, K. Song, D. Liang, T. Lu, P. Luo, and L. Shao, “Pyramid vision transformer: A versatile backbone for dense prediction without convolutions,” in *International conference on computer vision*. Montreal, Canada: IEEE, 2021, pp. 548–558, DOI: [10.1109/ICCV48922.2021.00061](https://doi.org/10.1109/ICCV48922.2021.00061).
- [22] X. Yang, H. Zhang, D. Jin, Y. Liu, C.-H. Wu, J. Tan, D. Xie, J. Wang, and X. Wang, “Fashion captioning: Towards generating accurate descriptions with semantic rewards,” in *European conference on computer vision*. Glasgow, United Kingdom: Springer, 2020, pp. 1–17, DOI: [10.1007/978-3-030-58601-0_1](https://doi.org/10.1007/978-3-030-58601-0_1).
- [23] Z. Al-Halah and K. Grauman, “From paris to berlin: Discovering fashion style influences around the world,” in *Conference on computer vision and pattern recognition*. Seattle, WA, USA: IEEE, 2020, pp. 10 133–10 142, DOI: [10.1109/CVPR42600.2020.01015](https://doi.org/10.1109/CVPR42600.2020.01015).
- [24] H. Tan and M. Bansal, “Lxmert: Learning cross-modality encoder representations from transformers,” in *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing*. Hong Kong, China: Association for Computational Linguistics, 2019, p. 5100–5111.
- [25] W. Su, X. Zhu, Y. Cao, B. Li, L. Lu, F. Wei, and J. Dai, “Vi-bert: Pre-training of generic visual-linguistic representations,” in *International Conference on Learning Representations*. [Online]: PMLR, 2020.
- [26] K.-H. Lee, X. Chen, G. Hua, H. Hu, and X. He, “Stacked cross attention for image-text matching,” in *European conference on computer vision*. Munich, Germany: Springer, 2018, pp. 212–228, DOI: [10.1007/978-3-030-01225-0_13](https://doi.org/10.1007/978-3-030-01225-0_13).
- [27] Z. Niu, M. Zhou, L. Wang, X. Gao, and G. Hua, “Hierarchical multimodal lstm for dense visual-semantic embedding,” in *International conference on computer vision*. Venice, Italy: IEEE, 2017, pp. 1899–1907, DOI: [10.1109/ICCV.2017.208](https://doi.org/10.1109/ICCV.2017.208).

- [28] J. Xia, M. Zhuge, T. Geng, S. Fan, Y. Wei, Z. He, and F. Zheng, “Skating-mixer: Multimodal mlp for scoring figure skating,” [Online], 2022, Available: <https://arxiv.org/abs/2203.03990>.
- [29] X. Li, X. Yin, C. Li, P. Zhang, X. Hu, L. Zhang, L. Wang, H. Hu, L. Dong, F. Wei *et al.*, “Oscar: Object-semantics aligned pre-training for vision-language tasks,” in *European conference on computer vision*. Glasgow, United Kingdom: Springer, 2020, pp. 121–137, DOI: [10.1007/978-3-030-58577-8_8](https://doi.org/10.1007/978-3-030-58577-8_8).
- [30] M. Zhuge, D.-P. Fan, N. Liu, D. Zhang, D. Xu, and L. Shao, “Salient object detection via integrity learning,” *Transactions on pattern analysis and machine intelligence*, 2022, DOI: [10.1109/TPAMI.2022.3179526](https://doi.org/10.1109/TPAMI.2022.3179526).
- [31] K. Xu, J. Ba, R. Kiros, K. Cho, A. Courville, R. Salakhudinov, R. Zemel, and Y. Bengio, “Show, attend and tell: Neural image caption generation with visual attention,” in *International conference on machine learning*. Lille, France: PMLR, 2015, pp. 2048–2057.
- [32] T. Arici, M. S. Seyfioglu, T. Neiman, Y. Xu, S. Train, T. Chilimbi, B. Zeng, and I. Tutar, “Mlim: Vision-and-language model pre-training with masked language and image modeling,” [Online], 2021, Available: <https://arxiv.org/abs/2109.12178>.
- [33] H. Bao, L. Dong, and F. Wei, “BEiT: BERT Pre-Training of Image Transformers,” in *ICLR*, 2022, Available: <https://openreview.net/forum?id=p-BhZSz59o4>.
- [34] K. He, X. Chen, S. Xie, Y. Li, P. Dollár, and R. Girshick, “Masked autoencoders are scalable vision learners,” in *Conference on computer vision and pattern recognition*. New Orleans, LA, USA: IEEE, 2022, pp. 15 979–15 988, DOI: [10.1109/CVPR52688.2022.01553](https://doi.org/10.1109/CVPR52688.2022.01553).
- [35] Z. Huang, Z. Zeng, B. Liu, D. Fu, and J. Fu, “Pixel-bert: Aligning image pixels with text by deep multi-modal transformers,” [Online], 2020, Available: <https://arxiv.org/abs/2004.00849>.
- [36] X. Lin, G. Bertasius, J. Wang, S.-F. Chang, D. Parikh, and L. Torresani, “Vx2text: End-to-end learning of video-based text generation from multimodal inputs,” in *Conference on computer vision and pattern recognition*. Nashville, TN, USA: IEEE, 2021, pp. 7001–7011, DOI: [10.1109/CVPR46437.2021.00693](https://doi.org/10.1109/CVPR46437.2021.00693).
- [37] W. Kim, B. Son, and I. Kim, “Vilt: Vision-and-language transformer without convolution or region supervision,” in *International conference on machine learning*. PMLR, 2021, pp. 5583–5594.
- [38] M. Yan, H. Xu, C. Li, B. Bi, J. Tian, M. Gui, and W. Wang, “Grid-vlp: Revisiting grid features for vision-language pre-training,” [Online], 2021, Available: <https://arxiv.org/abs/2108.09479>.
- [39] Z. Huang, Z. Zeng, Y. Huang, B. Liu, D. Fu, and J. Fu, “Seeing out of the box: End-to-end pre-training for vision-language representation learning,” in *Conference on computer vision and pattern recognition*. Nashville, TN, USA: IEEE, 2021, pp. 12 971–12 980, DOI: [10.1109/CVPR46437.2021.01278](https://doi.org/10.1109/CVPR46437.2021.01278).
- [40] S. Goenka, Z. Zheng, A. Jaiswal, R. Chada, Y. Wu, V. Hedau, and P. Natarajan, “Fashionvlp: Vision language transformer for fashion retrieval with feedback,” in *Conference on computer vision and pattern recognition*. New Orleans, LA, USA: IEEE, 2022, pp. 14 085–14 095, DOI: [10.1109/CVPR52688.2022.01371](https://doi.org/10.1109/CVPR52688.2022.01371).
- [41] J. Lei, L. Li, L. Zhou, Z. Gan, T. L. Berg, M. Bansal, and J. Liu, “Less is more: Clipbert for video-and-language learning via sparse sampling,” in *Conference on computer vision and pattern recognition*. Nashville, TN, USA: IEEE, 2021, pp. 7327–7337, DOI: [10.1109/CVPR46437.2021.00725](https://doi.org/10.1109/CVPR46437.2021.00725).
- [42] H. Xu, M. Yan, C. Li, B. Bi, S. Huang, W. Xiao, and F. Huang, “E2e-vlp: End-to-end vision-language pre-training enhanced by visual learning,” in *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language*

Processing. Association for Computational Linguistics, 2021, p. 503–513.

- [43] H. Akbari, L. Yuan, R. Qian, W.-H. Chuang, S.-F. Chang, Y. Cui, and B. Gong, “Vatt: Transformers for multimodal self-supervised learning from raw video, audio and text,” in *Advances in neural information processing systems*, vol. 34. Curran Associates, Inc., 2021, pp. 24 206–24 221.
- [44] X. Yi, J. Yang, L. Hong, D. Z. Cheng, L. Heldt, A. Kumthekar, Z. Zhao, L. Wei, and E. Chi, “Sampling-bias-corrected neural modeling for large corpus item recommendations,” in *Proceedings of the 13th ACM Conference on Recommender Systems*. New York, NY, USA: Association for Computing Machinery, 2019, p. 269–277, DOI: [10.1145/3298689.3346996](https://doi.org/10.1145/3298689.3346996).
- [45] O. Ronneberger, P. Fischer, and T. Brox, “U-Net: Convolutional networks for biomedical image segmentation,” in *International Conference on Medical Image Computing and Computer Assisted Intervention*. Munich, Germany: Springer, 2015, pp. 234–241, DOI: [10.1007/978-3-319-24574-4_28](https://doi.org/10.1007/978-3-319-24574-4_28).
- [46] C. Alberti, J. Ling, M. Collins, and D. Reitter, “Fusion of detected objects in text for visual question answering,” in *2019 Conference on Empirical Methods in Natural Language Processing and 9th International Joint Conference on Natural Language Processing, Proceedings of the Conference*. Hong Kong, China: Association for Computational Linguistics, 2019, pp. 2131–2140.
- [47] N. Rostamzadeh, S. Hosseini, T. Boquet, W. Stokowiec, Y. Zhang, C. Jauvin, and C. Pal, “Fashion-gen: The generative fashion dataset and challenge,” in *International conference on machine learning Workshops*, 2018.
- [48] R. Kiros, R. Salakhutdinov, and R. S. Zemel, “Unifying visual-semantic embeddings with multimodal neural language models,” [Online], 2014, Available: <https://arxiv.org/abs/1411.2539>.
- [49] F. Faghri, D. J. Fleet, J. R. Kiros, and S. Fidler, “Vse++: Improving visual-semantic embeddings with hard negatives,” in *British Machine Vision Conference*. Newcastle, UK: BMVA Press, 2018.
- [50] Y. Wang, H. Yang, X. Qian, L. Ma, J. Lu, B. Li, and X. Fan, “Position focused attention network for image-text matching,” [Online], 2019, Available: <https://arxiv.org/abs/1907.09748>.
- [51] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *Conference on computer vision and pattern recognition*. Miami, FL, USA: IEEE, 2009, pp. 248–255, DOI: [10.1109/CVPR.2009.5206848](https://doi.org/10.1109/CVPR.2009.5206848).
- [52] A. Radford, J. W. Kim, C. Hallacy, A. Ramesh, G. Goh, S. Agarwal, G. Sastry, A. Askell, P. Mishkin, J. Clark *et al.*, “Learning transferable visual models from natural language supervision,” in *International conference on machine learning*. PMLR, 2021, pp. 8748–8763.
- [53] T.-Y. Lin, M. Maire, S. Belongie, J. Hays, P. Perona, D. Ramanan, P. Dollár, and C. L. Zitnick, “Microsoft coco: Common objects in context,” in *European conference on computer vision*. Zurich, Switzerland: Springer, 2014, pp. 740–755, DOI: [10.1007/978-3-319-10602-1_48](https://doi.org/10.1007/978-3-319-10602-1_48).
- [54] G. Li, N. Duan, Y. Fang, M. Gong, D. Jiang, and M. Zhou, “Unicoder-vl: A universal encoder for vision and language by cross-modal pre-training,” in *AAAI Conference on Artificial Intelligence*. New York, NY, USA: AAAI Press, 2020, pp. 11 336–11 344.
- [55] L. Y. Wu, D. Liu, X. Guo, R. Hong, L. Liu, and R. Zhang, “Multi-scale spatial representation learning via recursive hermite polynomial networks,” in *Proceedings of the Thirty-First International Joint Conference on Artificial Intelligence*. Messe Wien, Vienna, Austria: International Joint Conferences on Artificial Intelligence Organization, 2022, pp. 1465–1473, DOI: [10.24963/ijcai.2022/204](https://doi.org/10.24963/ijcai.2022/204).

- [56] D. Chen and et al., “Cross-modal retrieval with heterogenous graph embedding,” in *Proceedings of the 30th ACM International Conference on Multimedia*. Lisboa, Portugal: Association for Computing Machinery, 2022, p. 3291–3300, DOI: [10.1145/3503161.3548195](https://doi.org/10.1145/3503161.3548195).
- [57] D. Liu, L. Wu, F. Zheng, L. Liu, and M. Wang, “Verbal-person nets: Pose-guided multi-granularity language-to-person generation,” *Transactions on Neural Networks and Learning Systems*, 2022, DOI: [10.1109/TNNLS.2022.3151631](https://doi.org/10.1109/TNNLS.2022.3151631).
- [58] Z. Zhang, H. Luo, L. Zhu, G. Lu, and H. T. Shen, “Modality-invariant asymmetric networks for cross-modal hashing,” *Transactions on Knowledge and Data Engineering*, 2022, DOI: [10.1109/TKDE.2022.3144352](https://doi.org/10.1109/TKDE.2022.3144352).

Ge-Peng Ji is currently a PhD student at Australian National University, supervised by Professor Nick Barnes, majoring in Engineering and Computer Science. Before that, he received an M. Sc. degree in communication and information systems from Wuhan University, China, in 2021. He has published about 10 peer-reviewed journal and conference papers. In 2021, he received the Student Travel Award from Medical Image Computing and Computer-Assisted Intervention Society.

His research interests lie in computer vision, especially in a variety of dense prediction tasks, such as video analysis, medical image segmentation, camouflaged object segmentation, and saliency detection.

E-mail: gepengai.ji@gmail.com

ORCID iD: 0000-0001-7092-2877

Mingchen Zhuge is a PhD student in KAUST, under the supervision of Prof. Juergen Schmidhuber. He received his M.S. degree in Computer Science from the China University of Geosciences in 2021. For the past two years, In 2019, he won the champion in the ZTE algorithm competition. He has worked as an intern at Alibaba Group and IIAI, as well as a visiting scholar at SUSTech. Besides, he has been invited to serve as a top

conference reviewer for CVPR, ICML, ECCV, NeurIPS, etc.

His primary research interests include multi-modal learning and reinforcement learning.

E-mail: mczhuge@gmail.com

ORCID iD: 0000-0003-2561-7712

Dehong Gao received the PhD degree from The Hong Kong Polytechnic University, Hong Kong, China, in 2014. He is now working as an associate professor in Northwestern Polytechnical University.

His research interests include information retrieval, recommendation, natural language processing and machine learning.

E-mail: gaodehong_polyu@163.com,
dehong.gdh@alibaba-inc.com

ORCID iD: 0000-0002-6636-5702

Deng-Ping Fan received his PhD degree from the Nankai University in 2019. He joined Inception Institute of Artificial Intelligence (IIAI) in 2019. He has published about 50 top journal and conference papers such as TPAMI, IJCV, TIP, TNNLS, TMI, CVPR, ICCV, ECCV, IJCAI, etc. He won the Best Paper Finalist Award at IEEE CVPR 2019, the Best Paper Award Nominee at IEEE CVPR 2020. He was recognized as the CVPR 2019 outstanding reviewer with a special mention award, the CVPR 2020 outstanding reviewer, the ECCV 2020 high-quality reviewer, and the CVPR 2021 outstanding reviewer. He served as a program committee board (PCB) member of IJCAI 2022-2024, a senior program committee (SPC) member of IJCAI 2021, a program committee members (PC) of CAD&CG 2021, a committee member of China Society of Image and Graphics (CSIG), area chair in NeurIPS 2021 Datasets and Benchmarks Track, area chair in MICCAI2020 Workshop.

His research interests include computer vision, deep learning, and visual attention, especially the human vision on co-salient object detection, RGB salient object detection, RGB-D salient object detection, and video salient object detection.

E-mail: dengpfan@gmail.com (Corresponding author)

ORCID iD: 0000-0002-5245-7518

Christos Sakaridis is a postdoctoral researcher at Computer Vision Lab, ETH Zurich. Since 2021,

he is the Principal Engineer in TRACE-Zurich, a project on computer vision for autonomous cars running at Computer Vision Lab and funded by Toyota Motor Europe. Moreover, he is the Team Leader in the EFCL project Sensor Fusion, in which they develop adaptive sensor fusion architectures for high-level visual perception. He obtained his PhD in Electrical Engineering and Information Technology from ETH Zurich in June 2021, working at Computer Vision Lab and supervised by Prof. Luc Van Gool. Prior to joining Computer Vision Lab, he received his MSc in Computer Science from ETH Zurich in 2016 and his Diploma in Electrical and Computer Engineering from the National Technical University of Athens in 2014, conducting his Diploma thesis at CVSP Group under the supervision of Prof. Petros Maragos.

His broad research fields are Computer Vision and Machine Learning. The focus of his research is on high-level visual perception, involving adverse visual conditions, domain adaptation, semantic segmentation, depth estimation, object detection, synthetic data generation, and fusion of multiple sensors including lidar, radar and event cameras, with emphasis on their application to autonomous cars and robots.

E-mail: csakarid@vision.ee.ethz.ch

ORCID iD: 0000-0003-1127-8887

Luc Van Gool received a degree in electromechanical engineering at the Katholieke Universiteit Leuven in 1981. Currently, he is a professor at the Katholieke Universiteit Leuven in Belgium and the ETH in Zurich, Switzerland. He leads computer vision research at both places and also teaches at both. He has been a program committee member of several major computer vision conferences. He received several Best Paper awards, won a David Marr Prize and a Koenderink Award, and was nominated Distinguished Researcher by the IEEE Computer Science committee. He is a co-founder of 10 spin-off companies.

His interests include 3D reconstruction and modeling, object recognition, tracking, gesture analysis, and a combination of those.

E-mail: vangool@vision.ee.ethz.ch

ORCID iD: 0000-0002-3445-5711