

Elevating *All* Zero-Shot Sketch-Based Image Retrieval Through Multimodal Prompt Learning

Mainak Singha¹, Ankit Jha^{1,2}, Divyam Gupta¹, Pranav Singla¹, and
Biplab Banerjee¹

¹ Indian Institute of Technology Bombay, India

² INRIA, Grenoble, France

{mainaksingha.iitb, ankitjha16, divsg1803, pranavsingla.c minds.iitb,
getbiplab}@gmail.com

Abstract. We address the challenges inherent in sketch-based image retrieval (SBIR) across various settings, including zero-shot SBIR, generalized zero-shot SBIR, and fine-grained zero-shot SBIR, by leveraging the vision-language foundation model CLIP. While recent endeavors have employed CLIP to enhance SBIR, these approaches predominantly follow uni-modal prompt processing and overlook to exploit CLIP’s integrated visual and textual capabilities fully. To bridge this gap, we introduce SPLIP, a novel multi-modal prompt learning scheme designed to operate effectively with frozen CLIP backbones. We diverge from existing multi-modal prompting methods that treat visual and textual prompts independently or integrate them in a limited fashion, leading to suboptimal generalization. SPLIP implements a bi-directional prompt-sharing strategy that enables mutual knowledge exchange between CLIP’s visual and textual encoders, fostering a more cohesive and synergistic prompt processing mechanism that significantly reduces the semantic gap between the sketch and photo embeddings. In addition to pioneering multi-modal prompt learning, we propose two innovative strategies for further refining the embedding space. The first is an adaptive margin generation for the sketch-photo triplet loss, regulated by CLIP’s class textual embeddings. The second introduces a novel task, termed conditional cross-modal jigsaw, aimed at enhancing fine-grained sketch-photo alignment by implicitly modeling sketches’ viable patch arrangement using knowledge of unshuffled photos. Our comprehensive experimental evaluations across multiple benchmarks demonstrate the superior performance of SPLIP in all three SBIR scenarios. Project page: <https://mainaksingha01.github.io/SplIP/>.

Keywords: CLIP · Sketch Based Image Retrieval · Prompt Learning

1 Introduction

Hand-drawn sketches adeptly convey abstract ideas with their simple yet evocative lines. The advent of touchscreen mobile devices has propelled sketch-based

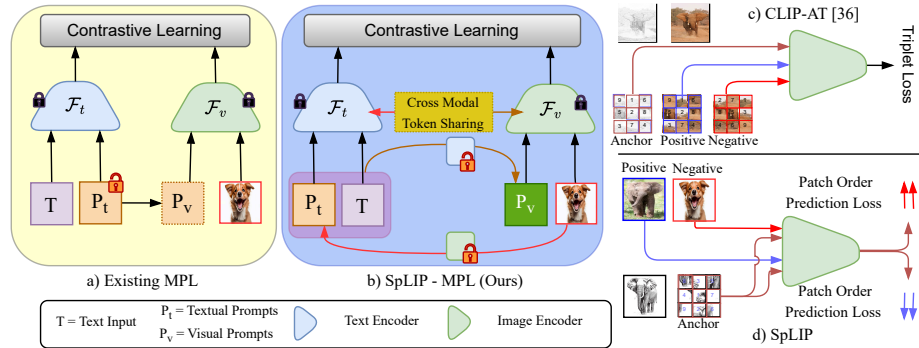


Fig. 1: (a, b) The difference between the existing multi-modal prompt learning (MPL) vs ours. As opposed to the literature [20, 21, 37], we propose to enhance the generalizability of both the textual and visual prompt embeddings with mutual knowledge sharing in a principled layer-wise fashion. **(c, d) Proposed conditional cross-modal jigsaw vs the literature [39].** As against a triplet loss connecting sketch-photo with the same permutation while contrasting against the photo with a different permutation, we propose better learning of patch arrangements by positively associating a permuted sketch with its intact photo counterpart through a novel objective.

image retrieval (SBIR) [25, 61] into the limelight, offering myriad practical uses. SBIR retrieves photos from a vast repository based on the same category as a query sketch. However, despite representing the same class, sketches and photos often differ noticeably due to their distinct domains. Traditional methods [14, 43] address domain heterogeneity, assuming full visibility of test classes during training, showing promise in effective retrieval. Yet, a more realistic challenge emerges when test set categories remain unseen during training, termed zero-shot SBIR (ZS-SBIR) [5, 11, 28]. Contrarily, generalized ZS-SBIR (GZS-SBIR) [11] encompasses retrieval photos of known and novel classes for novel-class sketch queries during inference, heightening complexity. Additionally, instance-level fine-grained ZS-SBIR (FG-ZS-SBIR) [31, 58] focuses on precise shape matching, intensifying challenges compared to category-level SBIR.

The crux of all the SBIR settings lies in learning an embedding space for the sketches and photos that is *unbiased* to the training data, *discriminative* given data from both the photo and sketch modalities and hence, *domain-agnostic*.

Leading SBIR frameworks leveraging ConvNets and ViTs [5, 11, 52, 55] face intrinsic semantic constraints by their architecture. On the contrary, the advent of multi-modal foundational models, notably CLIP [34] and Align [19], has markedly enhanced visual comprehension by integrating visual and textual information. These VLMs have excelled in cross-domain learning, but their SBIR integration remains limited. Initiatives like [39] have aimed to tailor CLIP to ZS-SBIR and FG-ZS-SBIR, focusing on visual prompts and patch shuffling to align sketches and photos at both micro and macro levels. Other efforts [29, 59] seek to finetune CLIP’s embeddings for SBIR, indicating an increasing desire to leverage these multi-modal platforms beyond original purposes.

Highlighting research gaps: Despite the success, these strategies often rely on simplistic, one-dimensional prompting, failing to harness CLIP’s dual-pathway synergy fully. This overlooks CLIP’s visual-textual fusion’s rich, complementary insights, leading to suboptimal ZS-SBIR performance. Addressing this, there’s an urgent call for novel approaches that dynamically utilize this combined knowledge to surpass the semantic limitations of existing models, thus expanding SBIR potentials.

Moving away from singular prompting, multi-modal methods recommended by [20, 21, 37] employ simultaneous prompting across CLIP’s pathways, thus narrowing the semantic gap in embeddings. Nonetheless, these approaches struggle with one-sided prompt integration and the dismissal of static textual elements, causing a reduction in semantic depth. Particularly, the adaptability of the textual pathway is compromised, remaining insensitive to visual nuances, even as [63] underscores the importance of enhancing textual adaptability in CLIP. We advocate for a more cohesive and dynamic knowledge interchange between the visual and textual domains to enhance the flexibility of (G)ZS-SBIR. Furthermore, although the patch shuffling strategy by [39] enhances the shape equivalence of images and sketches, especially for FG-ZS-SBIR, indiscriminate matching of patch-permuted sketch-photo versions without regard to the natural qualities of objects could lead to overfitting. We recommend matching sketch patch permutations with semantically corresponding, unshuffled photos and the reverse, helping embeddings to discern how patch arrangements correspond with the entirety of an image’s objects.

Proposed solution: Our proposed model, SP-LIP, tackles these challenges by implementing a novel deep multimodal prompting approach (Fig. 1), facilitating efficient knowledge exchange between text and image branches of frozen CLIP. In CLIP’s text encoder, we enhance static textual tokens by incorporating additional latent visual tokens at each layer, departing from previous random initialization methods [20, 37]. Likewise, we enrich the image embeddings within the CLIP vision backbone with information from these image-conditioned textual token embeddings, summarized layer by layer over all the semantic categories. This *bidirectional* information sharing mitigates the semantic gap in the obtained embeddings, thus contributing positively towards zero-shot inference.

As we advance, we introduce a novel conditional cross-modal jigsaw task to strengthen the linkage between photo and sketch pairs for all the considered SBIR variants, predominantly for FG-ZS-SBIR. This method requires solving a jigsaw puzzle [30] using an anchor sketch, assisted by a positive and a negative image from the alternate modality while ensuring that the positive image significantly aids in this process, thereby enhancing model generalization through the learning of patch arrangement for reconstructing a complete image. This approach deviates from previous works [33, 39] that either deal with mixed-modal images for a naive jigsaw solver or employ triplet objectives with uniform or varying patch permutations for identifying positive and negative pairs, thus failing to bridge the gap between local and global contexts effectively (Fig. 1). Finally, we integrate the gold-standard cross-modal triplet loss for SBIR and

a sketch/photo-text-based classification loss of CLIP, introducing an adaptive margin scheme for the triplet objective derived from the semantic class-prompt embeddings. Our salient contributions are, therefore,

- Introducing SPLIP, a novel deep multi-modal prompt tuning model within the realm of CLIP tailored for ZS-SBIR and FG-ZS-SBIR tasks. It introduces a more systematic cross-modal mutual guidance within CLIP’s text and vision encoders. To the best of knowledge, this is the first endeavor of multi-modal prompting for solving ZS-SBIR variants.

- We enhance the conventional cross-modal triplet loss objective in ZS-SBIR by incorporating an adaptive margin scheme, leveraging CLIP’s textual prompt embeddings. Additionally, we introduce a novel conditional cross-modal jigsaw task aimed at refining the association between sketch and photo pairs.

- We conduct extensive experiments on three benchmark datasets: Sketchy-Ext [11, 44, 57], TU-Berlin-Ext [13, 27], and QuickDraw-Ext [7, 17], covering (G)ZS-SBIR and FG-ZS-SBIR settings. SPLIP consistently outperforms existing competitors, achieving significant improvements in all the metrics.

2 Related Works

2.1 Sketch-based image retrieval (SBIR)

(Generalized) Zero-shot SBIR: SBIR tasks involve retrieving photos corresponding to specific categories from a diverse collection of multi-category images based on a given sketch query, demanding a thorough understanding of the joint sketch-photo manifold. The literature is rich in fully-supervised SBIR endeavors utilizing deep and hand-crafted descriptors and involving different learning mechanisms, including both generative and discriminative approaches [18, 24, 32, 43, 44, 50, 65]. Recently, [4] introduced a data-free training strategy for SBIR, relaxing the need for curating sketch-photo pairs.

Researchers have addressed challenges in recognizing unseen test-time classes through ZS-SBIR, facilitating category-level generalization. Extending the traditional ConvNet-based frameworks [7, 11, 28, 57], graph convolutional networks, ViTs, and their combinations have been introduced to learn an unbiased shared feature space [16, 26, 62]. Recent advancements [8, 29, 39, 59] leverage CLIP’s zero-shot inference by integrating text with sketches and photos, outperforming counterparts. A combined loss function of supervised cross-entropy and metric objectives fosters a discriminative embedding space. On the other hand, GZS-SBIR permits the presence of training and test time photos during inference, causing the model to show high bias towards the training data, which has been tackled in the literature from different perspectives [5, 11, 15, 26, 29, 36, 66]. *We take a different route to tackle (G/FG)ZS-SBIR through multi-modal prompting in CLIP, thus modeling the visual-semantic synergy effectively.*

Fine-grained ZS-SBIR: Transitioning from categorical ZS-SBIR, FG-ZS-SBIR aims at identifying specific photos relating to sketches at an instance level. Initially rooted in a deep triplet-ranking Siamese framework [58], FG-ZS-SBIR’s evolution incorporates attention modules [42, 50], hybrid cross-domain mapping [32], and manifold modeling for universality [1]. Enhancements proceed

with an intra-modal triplet goal, solutions for sparse sketch annotations [40], patch similarity via cross-interaction [51], and innovative patch shuffling [33, 39]. *Exploiting patch shuffling for precise shape alignment and utilizing CLIP’s robust capabilities, we propose a unique conditional cross-modal jigsaw challenge. Aimed at refining alignment between permutations of sketch patches and intact photos, this initiative significantly deepens contextual comprehension from a local to a global scale, marking a pivotal advance in FG-ZS-SBIR development.*

2.2 Vision-language models and multi-modal prompt learning

Vision-Language Models (VLMs) like CLIP [34] and VisualBERT [23] have revolutionized computer vision by merging visual and textual data through multimodal learning, creating detailed representations. Leveraging textual features from language models (e.g., BERT [6], GPT [35]) and visual inputs from ConvNets or ViTs [9], these models achieve semantic depth and exhibit strong zero-shot inference for varied tasks [2, 20].

Recent research [2, 3, 47–49, 63, 64] highlights prompt learning as a viable alternative to VLM fine-tuning on downstream tasks. To this end, unlike uni-modal approaches for CLIP, multi-modal deep prompts synergize its visual and textual components. [20] proposed to learn both visual and textual prompts and showed that initializing the visual prompts from the textual counterparts enhances the performance. In addition, [21] and [37] introduced regularization and feature consistency to prevent overfitting and ensure textual variety. *However, existing methods mainly employ unidirectional token sharing, restricting the overall generalizability and semantic depth of the learned embeddings. Contrarily, we propose a bilateral approach to disseminate relevant cross-modal insights across CLIP’s branches, establishing a more effective multi-modal prompting paradigm.*

3 Methodology

SBIR entails retrieving \mathcal{K} photos $\{p_k\}_{k=1}^{\mathcal{K}} \in \mathcal{P}$ from a gallery (\mathcal{G}), given a query-sketch ($s \in \mathcal{S}$) belonging to a specific category out of a total of \mathcal{C} classes. In zero-shot tasks, \mathcal{C} is divided into seen training classes (\mathcal{C}^s) and unseen testing classes (\mathcal{C}^u), where $\mathcal{C} = \mathcal{C}^s \cup \mathcal{C}^u$ and $\mathcal{C}^s \cap \mathcal{C}^u = \emptyset$.

The training dataset $\mathcal{G}^s = (\mathcal{S}^s, \mathcal{P}^s, \mathcal{C}^s)$ comprises sketches \mathcal{S}^s and photos \mathcal{P}^s from \mathcal{C}^s categories. During inference, the gallery $\mathcal{G}^u = (\mathcal{S}^u, \mathcal{P}^u, \mathcal{C}^u)$ containing sketches \mathcal{S}^u and photos \mathcal{P}^u with category labels in \mathcal{C}^u is involved. In contrast to ZS-SBIR, the GZS-SBIR setup considers photos in $\mathcal{P}^s \cup \mathcal{P}^u$ for a given sketch query $s^u \in \mathcal{S}^u$ for retrieval during inference. FG-ZS-SBIR aims at instance-level sketch-photo matching within specific categories [58], in contrast to the broader category-level focus of conventional SBIR methods. Following [39], we consider the multi-category FG-ZS-SBIR setting where paired sketch-photo instances are available from multiple categories.

In the following, we delve into the initialization of text inputs and provide a detailed explanation of the image-text embeddings of CLIP in Section 3.1. Moving forward, our image-driven textual prompting approach is elaborated upon in Section 3.2, while Section 3.3 addresses the proposed visual prompting

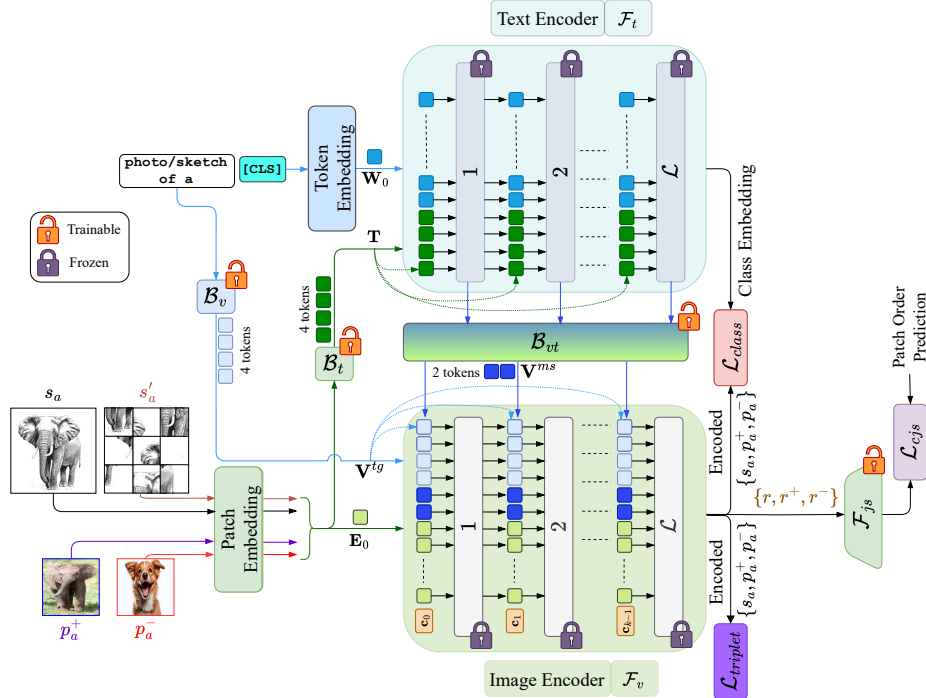


Fig. 2: The model architecture for SPLIT, which capitalizes on CLIP’s static text and vision backbones, \mathcal{F}_t and \mathcal{F}_v , introducing a bidirectional prompt exchange. Image patch embeddings, transformed via \mathcal{B}_t , generate textual tokens \mathbf{T} for different layers of \mathcal{F}_t . Similarly, \mathcal{F}_v layers receive "sketch/photo of a" token embeddings \mathbf{V}^{tg} through \mathcal{B}_v and consolidated prompt tokens from \mathcal{F}_t^l , across all semantic classes in \mathcal{C}^s , \mathbf{V}^{ms} , via \mathcal{B}_{vt} . This setup enriches both textual and visual pathways with diverse information sources. The model also tackles a unique conditional cross-modal jigsaw task, with a decoder \mathcal{F}_{js} processing sketch-photo pairs from $(s_a, s'_a, p_a^+, p_a^-)$ to understand complex relationships through pairwise fused features of s'_a and the remaining counterparts, (r, r^+, r^-) . Training involves a blend of loss functions: photo-sketch triplet loss $\mathcal{F}_{triplet}$, text-image classification loss \mathcal{F}_{class} , and proposed jigsaw loss \mathcal{L}_{cjs} . During inference on \mathcal{G}^u , \mathbf{V}^{ms} is derived leveraging the classes in \mathcal{C}^s , overlooking the need of \mathcal{C}^u , and leading to a nearest-neighbor based ranking of photos for sketch queries in the output of \mathcal{F}_v .

methodology. For fine-grained sketch-photo feature association, we discuss the proposed conditional cross-modality jigsaw task and the related details in Section 3.4. Finally, in Section 3.5, we mention the considered loss objectives. A list of important variables is summarized in the **Supplementary**.

3.1 Initialization of visual-textual embeddings

The pre-trained CLIP model operates on two modalities: text and image. It consists of a transformer-based [53] text encoder (\mathcal{F}_t) and ViT-based [9] image encoder (\mathcal{F}_v). They both contain \mathcal{L} transformer encoder layers. \mathcal{F}_t generates feature representations for text descriptions to capture semantic information. Initially,

it tokenizes the text inputs, consisting of \mathcal{J} words, and projects them into word embeddings $\mathbf{W}_0 = [w_0^1, w_0^2, \dots, w_0^{\mathcal{M}}] \in \mathbb{R}^{\mathcal{M} \times d_t}$, where $[\cdot, \cdot]$ denotes stacking and concatenation, \mathcal{M} refers the number of embedding tokens (77 per class-prompt) and d_t is the dimension of text tokens. In our approach, we use the input texts as ‘‘sketch/photo of a [CLS]’’ for sketches and photos, respectively, for (G)ZS-SBIR, meaning $\mathcal{J} = 5$. Here, the [CLS] token represents class embeddings, completing the prompt embedding \mathbf{W}_0 . For the FG-ZS-SBIR task, we use common text input, ‘‘visual representation of [CLS]’’, for both sketches and photos. We obtain the l -th layer embedding of \mathcal{F}_t , denoted as \mathcal{F}_t^l , as follows,

$$[\mathbf{W}_l] = \mathcal{F}_t^l(\mathbf{W}_{l-1}) \in \mathbb{R}^{\mathcal{M} \times d_t} \quad l = 1, 2, \dots, \mathcal{L} \quad (1)$$

On the image side, the images from \mathcal{P} and \mathcal{S} are partitioned into fixed-size patches and encoded through \mathcal{F}_v . Each patch undergoes projection to generate initial patch embeddings (\mathbf{E}_0), along with a learnable cls token \mathbf{c}_0 , denoted as $[\mathbf{c}_0, \mathbf{E}_0] \in \mathbb{R}^{1+\mathcal{N} \times d_v}$, where \mathcal{N} denotes the number of patches and d_v is the dimension of patch tokens. Henceforth, the output embedded tokens of the l -th layer of \mathcal{F}_v can be expressed as,

$$[\mathbf{c}_l, \mathbf{E}_l] = \mathcal{F}_v^l([\mathbf{c}_{l-1}, \mathbf{E}_{l-1}]) \in \mathbb{R}^{1+\mathcal{N} \times d_v} \quad l = 1, 2, \dots, \mathcal{L} \quad (2)$$

3.2 Proposed vision-guided deep textual prompting

We propose a novel *vision-guided deep textual prompting* approach, where deep prompting determines slightly changing the input raw tokens of each of the layers of \mathcal{F}_t for both text inputs associated with sketches and photos. Our proposal involves incorporating visual information into the tuning process of textual prompts. Specifically, we introduce a *visual-to-textual mapping* block (\mathcal{B}_t), which generates m learnable prompt tokens, denoted ($\mathcal{T}_{1:m}$) collectively as (\mathbf{T}), from the \mathcal{N} visual patch embeddings \mathbf{E}_0 (Fig. 2).

Precisely, each layer in \mathcal{F}_t receives (\mathbf{T}) in the corresponding input space. In the first layer, \mathbf{T} (aka \mathbf{T}_0) replaces m tokens of \mathbf{W}_0 , and the input embedding for the first layer of \mathcal{F}_t becomes $[\mathbf{T}_0; \mathbf{W}_0]$. Here, $[a; b]$ denotes stacking after replacing a similar number of tokens of b with all of the tokens of a . Consequently, for the l -th layer, $\mathbf{T}_l = \mathbf{T} = \mathcal{B}_t(\mathbf{E}_0)$.

As already pointed out, our approach differs from existing literature [20, 37] in that our learnable tokens in the textual prompts capture visual distributions, as opposed to the random initialization approach followed by [20, 37].

Finally, the output operation for the l -th layer can be expressed as,

$$[_, \mathbf{W}_l] = \mathcal{F}_t^l([\mathbf{T}_{l-1}; \mathbf{W}_{l-1}]) \in \mathbb{R}^{\mathcal{M} \times d_t} \quad l = 1, 2, \dots, \mathcal{L} \quad (3)$$

3.3 Proposed text-guided deep visual prompting

Our visual prompting approach leverages textual information within CLIP via two distinct mechanisms. We harness the initial tokenized text input excluding the [CLS] token, denoted as \mathcal{W}' comprising of $(\mathcal{J} - 1)$ tokens. These are employed as *semantic domain knowledge* (\mathbf{V}^{sg}), which is then mirrored by an equivalent stack of $(\mathcal{J} - 1)$ learnable tokens ($\mathcal{V}_{1:\mathcal{J}-1}^{\text{tg}}$). This mirroring is facilitated via a

textual-to-visual mapping block (\mathcal{B}_v), which operates across all layers of \mathcal{F}_v . For any given layer ($l + 1$) within \mathcal{F}_v , this process can be succinctly described as: $\mathbf{V}_l^{\text{tg}} = \mathbf{V}^{\text{tg}} = \mathcal{B}_v(\mathcal{W}')$, effectively embedding semantic textual insights into the visual domain for enhanced model comprehension and interaction (Fig. 2).

In addition, we enforce token sharing from each layer of \mathcal{F}_t to the respective layer input of \mathcal{F}_v . Note that, from $l \geq 1$, \mathbf{W}_l implicitly includes visual knowledge as per our proposed textual prompting. Secondly, as opposed to the visual prompting approach proposed in [20], which only shares the learnable textual tokens with the visual branch, we propose to consider all the tokens from \mathbf{W}_l over all the classes present in \mathcal{C}^s to be included into the input space of \mathcal{F}_v^l . This class-agnostic knowledge-sharing approach helps diminish the semantic gap in the learned embeddings.

We proceed by transferring the output (\mathbf{W}_l) of the l -th layer of \mathcal{F}_t to a *vision-text conjunction block* (\mathcal{B}_{vt}), which then generates scale-specific inputs ($\mathbf{V}_{l-1}^{\text{ms}}$) for the corresponding l -th layer of \mathcal{F}_v consisting of n learnable tokens ($\mathcal{V}_{1:n}^{\text{ms}}$). As mentioned, \mathcal{B}_{vt} takes all of the class-defined text feature tokens of \mathbf{W}_l , as input. Evidently, unlike \mathbf{V}^{tg} , the prompt tokens of \mathbf{V}^{ms} are not similar to each other for every layers. We note that (\mathcal{B}_{vt}) is shared across all the encoder layers of \mathcal{F}_t . For the l -th layer of \mathcal{F}_v , \mathbf{V}^{ms} can be defined as,

$$\mathbf{V}_{l-1}^{\text{ms}} = \{\mathcal{V}_{k_{l-1}}^{\text{ms}} \in \mathbb{R}^{d_v}\}_{k=1}^n = \mathcal{B}_{vt}(\mathcal{F}_t^l([\mathbf{T}_{l-1}; \mathbf{W}_{l-1}])) \in \mathbb{R}^{n \times d_v}, \quad l = 1, 2, \dots, \mathcal{L} \quad (4)$$

Finally, we concat the generated prompt tokens of \mathbf{V}^{tg} and \mathbf{V}^{ms} for each of the layers, expressed as a common *visual prompt* (\mathbf{V}) i.e. for inputting to the l -th layer of \mathcal{F}_v : $\mathbf{V}_{l-1} = [\mathbf{V}_{l-1}^{\text{tg}}, \mathbf{V}_{l-1}^{\text{ms}}]$. Hence, the processing at the l -th layer of \mathcal{F}_v is mentioned as,

$$[\mathbf{c}_{l, -}; \mathbf{E}_l] = \mathcal{F}_v^l([\mathbf{c}_{l-1}, \mathbf{V}_{l-1}; \mathbf{E}_{l-1}]) \in \mathbb{R}^{1+\mathcal{N} \times d_v} \quad l = 1, 2, \dots, \mathcal{L} \quad (5)$$

3.4 Proposed conditional cross-modal jigsaw solver for fine-grained sketch-photo feature association

Furthermore, we introduce the task of *conditional cross-modal jigsaw* to enhance the intricate connections between photos and sketches belonging to identical classes (or specific instances in the context of FG-ZS-SBIR). This technique marks a departure from previous methods, such as the one by [33], which created a hybrid image by interspersing patches between sketches and photos randomly and pre-training the feature extraction backbone to predict the sequence of patches in this blended image. Similarly, [39] utilized a uniform permutation for sketch-photo pairs to delineate positive pairs, contrasting this with a distinct permutation on the photo to forge the negative pair. However, these methods grapple with the challenge of precisely aligning the shuffled image with its original format, a critical step for accurately learning patch arrangements in detail and maintaining the spatial coherence within the images.

To address these challenges, our proposed method incorporates positive and negative counterparts (p_a^+, p_a^-) from set \mathcal{P}^s when resolving the jigsaw puzzle for a sample $s_a \in \mathcal{S}^s$ given its permuted version s'_a which is obtained through a permutation function δ , given the random permutation $y^{\text{perm}} \in \mathcal{Y}^{\text{perm}}$: $s'_a =$

$\delta(s_a, y^{\text{perm}})$. Here, the information embedded in (s'_a, p_a^+) and (s'_a, p_a^-) is then processed through a transformer-based jigsaw-solver network, \mathcal{F}_{js} , operating atop \mathcal{F}_v , to leverage p_a^+ to inform better the prediction of the permutation arrangements of s'_a compared to p_a^- , achieved through a hinge objective. Collectively, \mathcal{F}_{js} directly resolves the jigsaw puzzle for s'_a when coupled with s_a . These combined losses are called \mathcal{L}_{cjs} . For simplicity, we define the fused inputs from pairs of images to \mathcal{F}_{js} as follows: $r = [\mathcal{F}_v(s_a), \mathcal{F}_v(s'_a)]$, $r^+ = [\mathcal{F}_v(p_a^+), \mathcal{F}_v(s'_a)]$, and $r^- = [\mathcal{F}_v(p_a^-), \mathcal{F}_v(s'_a)]$, respectively.

3.5 Loss functions and inference

Following [39], we train the LayerNorm parameters (θ, ϕ) of \mathcal{F}_v and \mathcal{F}_t , together with $(\mathcal{B}_v, \mathcal{B}_t, \mathcal{B}_{vt}, \mathcal{F}_{js})$ while keeping the other layers fixed.

$\mathcal{L}_{triplet}$: Cross-visual modality triplet loss with proposed adaptive margin: Given the triplet of embeddings $(\mathcal{F}_v(s_a), \mathcal{F}_v(p_a^+), \mathcal{F}_v(p_a^-))$, we formulate a triplet objective aimed at minimizing the distance between $\mathcal{F}_v(s_a)$ and $\mathcal{F}_v(p_a^+)$ while maximizing the distance between $\mathcal{F}_v(s_a)$ and $\mathcal{F}_v(p_a^-)$. Traditional triplet objectives employ a fixed margin, which may not be optimal for zero-shot tasks. In contrast, we propose leveraging the semantic space of CLIP to define the margin, utilizing the embeddings of positive and negative class names from \mathcal{F}_t . Let $\mathcal{F}_t(c^+)$ and $\mathcal{F}_t(c^-)$ represent the semantic class embeddings for p_a^+/s_a and p_a^- , respectively. We define $\mu(c^+, c^-) = \cos(\mathcal{F}_t(c^+), \mathcal{F}_t(c^-))$. Consequently, the margin increases when c^+ and c^- are semantically close. The loss is,

$$\mathcal{L}_{triplet} = \min_{\substack{\mathcal{B}_v, \mathcal{B}_t, \mathcal{B}_{vt}, \\ \theta, \phi}} \sum_{(s_a, p_a^+, p_a^-) \in \mathcal{G}^s} \left[\|\mathcal{F}_v(s_a) - \mathcal{F}_v(p_a^+)\|_2^2 - \|\mathcal{F}_v(s_a) - \mathcal{F}_v(p_a^-)\|_2^2 + \mu(c^+, c^-) \right]_+ \quad (6)$$

\mathcal{L}_{class} : Text-image classification loss: A supervised contrastive loss is introduced to correctly classify sketches and photos in \mathcal{G}^s based on the class-wise text prompts outlined in Section 3.1. In Eq. 7, (\mathbf{I}, \mathbf{Y}) denotes a sketch or a photo along with the associated one-hot label vector: $\mathbf{I} \in \{\mathcal{P}^s, \mathcal{S}^s\}$, and $\mathbf{Y} = [y^1, \dots, y^{|\mathcal{C}^s|}]$.

$$\mathcal{L}_{class} = \min_{\substack{\mathcal{B}_v, \mathcal{B}_t, \mathcal{B}_{vt}, \\ \theta, \phi}} \sum_{(\mathbf{I}, \mathbf{Y}) \in \mathcal{G}^s} - \sum_{c=1}^{|\mathcal{C}^s|} y^c \log(p(y^c | \mathbf{I})) \quad (7)$$

We compute $p(y^{c'} | \mathbf{I})$ as follows, where $\text{Prompt}_{y^{c'}} = \text{sketch/photo of a } [CLS_{y^{c'}}]$ represents input text sentence, and τ is a hyper-parameter.

$$p(y^{c'} | \mathbf{I}) = \frac{\exp(\cos(\mathcal{F}_v(\mathbf{I}), \mathcal{F}_t(\text{Prompt}_{y^{c'}})) / \tau)}{\sum_{c=1}^{|\mathcal{C}^s|} \exp(\cos(\mathcal{F}_v(\mathbf{I}), \mathcal{F}_t(\text{Prompt}_{y^c})) / \tau)} \quad (8)$$

\mathcal{L}_{cjs} : Proposed conditional cross-modal jigsaw loss: Given (r, r^+, r^-) and y^{perm} , \mathcal{F}_{js} follows the approach outlined in [30] to address the jigsaw task. This involves predicting the permutation index corresponding to y^{perm} in a list $\mathbf{\Pi}^{1 \times |\mathcal{Y}^{\text{perm}}|}$ that indexes all possible permutations in $\mathcal{Y}^{\text{perm}}$, where $\mathbf{1}_{\mathbf{\Pi}(y^{\text{perm}})}$ represents the one-hot encoding for the index of y^{perm} .

The loss function \mathcal{L}_{cjs} serves two objectives. Firstly, utilizing r and $\mathbf{1}_{\mathbf{\Pi}(y^{\text{perm}})}$, we aim to minimize the cross-entropy loss \mathcal{L}_{ce} , enabling \mathcal{F}_{js} to learn the jigsaw

task effectively. Additionally, we introduce a hinge-loss objective to constrain r^+ to yield a lower $\mathcal{L}_{ce}(\mathcal{F}_{js}(r^+), \mathbf{1}_{\Pi(y^{\text{perm}})})$ compared to $\mathcal{L}_{ce}(\mathcal{F}_{js}(r^-), \mathbf{1}_{\Pi(y^{\text{perm}})})$.

$$\mathcal{L}_{cjs} = \min_{\mathcal{B}_v, \mathcal{B}_t, \mathcal{B}_{vt}, \mathcal{F}_{js}, \theta, \phi} \sum_{(s_a, s'_a, p_a^+, p_a^-) \in \mathcal{G}^s} [\mathcal{L}_{ce}(\mathcal{F}_{js}(r), \mathbf{1}_{\Pi(y^{\text{perm}})}) + \mathcal{L}_{margin}] \quad (9)$$

where \mathcal{L}_{margin} is defined as,

$$\mathcal{L}_{margin} = [\mathcal{L}_{ce}(\mathcal{F}_{js}(r^+), \mathbf{1}_{\Pi(y^{\text{perm}})}) - \mathcal{L}_{ce}(\mathcal{F}_{js}(r^-), \mathbf{1}_{\Pi(y^{\text{perm}})})]_+ \quad (10)$$

\mathcal{L}_{total} : **Total loss for training:** The total loss can be denoted as follows with α and β denoting the relative loss weights.

$$\mathcal{L}_{total} = \mathcal{L}_{triplet} + \alpha * \mathcal{L}_{class} + \beta * \mathcal{L}_{cjs} \quad (11)$$

Inference: During inference, we eliminate the dependence on test class labels from \mathcal{C}^u . For generating the tokens \mathbf{V}^{ms} , where only the class embeddings are required, we propose to utilize class names from \mathcal{C}^s , mirroring the approach used during training, where all the class names are used to define the inputs to \mathcal{B}_{vt} . This strategy aims to map test classes into the discriminative space defined by the training classes, effectively combating model bias. The retrieval process of photos, in response to sketch queries, both from \mathcal{G}^u , is executed within \mathcal{F}_v 's visual embedding space using a nearest-neighbor ranking mechanism.

4 Experiments and Results

- **Datasets:** We evaluate SPLIP on three benchmark datasets for categorical (G)ZS-SBIR: Sketchy-Ext [11, 57], TU-Berlin-Ext [27], and QuickDraw-Ext [7], following the established training-validation protocols [5, 11]. For FG-ZS-SBIR, which necessitates precise sketch-photo matching, we incorporate the Sketchy dataset [44]. Further details, including the splits of Sketchy-1-Ext [11] & Sketchy-2-Ext [57] are described in the **Supplementary**.

- **Implementation details, training and evaluation protocols:** For \mathcal{F}_v , we select the ViT-B/32 backbone of CLIP, while the Transformer-based text encoder is considered for \mathcal{F}_t . Besides, \mathcal{B}_t employs three linear layers to convert \mathbf{E}_0 into $m = 4$ learnable textual tokens. In contrast, \mathcal{B}_v utilizes a singular linear layer to adapt to the visual dimension, producing four learnable visual tokens for a matching batch size. Meanwhile, \mathcal{B}_{vt} incorporates linear layers and employs a bottleneck architecture consisting of two layers (Linear-ReLU-Linear) for the creation of $n = 2$ layer-specific visual tokens. Additionally, \mathcal{F}_{js} is designed with two encoder layers, followed by a classifier, to accurately decode the patch arrangements of s'_a given $(p_a^+/p_a^-/s_a)$.

The training process spans 60 epochs, initiating with a warm-up learning rate of 0.001 and utilizing the Adam optimizer [22] alongside a scheduler. The batch size is configured to 192 for both Sketchy-Ext and TU-Berlin-Ext datasets, whereas a batch size of 64 is adopted for QuickDraw-Ext. α and β are fixed through grid search. Following the literature [7, 55], our evaluation for ZS-SBIR

Table 1: Comparison for categorical ZS-SBIR.

Methods		Sketchy-1-Ext [11]		Sketchy-2-Ext [57]		TU-Berlin-Ext [27]		QuickDraw-Ext [7]		
		mAP@all	P@100	mAP@200	P@200	mAP@all	P@100	mAP@all	P@200	
CNN	CVAE [57]	ECCV'18	19.6	28.4	22.5	33.3	00.5	00.1	00.3	00.3
	CC-DG [31]	CVPR'19	31.1	46.8	-	-	24.7	39.2	-	-
	Doodle [7]	CVPR'19	-	-	46.1	37.0	11.0	12.1	07.5	06.8
	SEM-PCYC [11]	CVPR'19	34.9	46.3	45.9	37.0	29.7	42.6	17.7	18.4
	SAKE [28]	ICCV'19	54.7	69.2	49.7	59.8	47.5	59.9	-	-
	Styleguide [12]	TM'20	37.6	48.4	35.8	40.0	25.4	35.5	-	-
	OCEAN [66]	ICME'20	-	-	-	-	33.3	46.7	-	-
	DSN [56]	IJCAI'21	58.3	70.4	-	-	48.4	59.1	-	-
	TCN [54]	TPAMI'21	61.6	76.3	51.6	60.8	49.5	61.6	14.0	29.8
	BDA [5]	NC'22	43.7	51.4	55.6	45.8	37.4	50.4	15.4	35.5
	Sketch3T [41]	CVPR'22	-	-	57.9	64.8	50.7	67.1	-	-
	ViT	TVT [52]	AAAI'22	64.8	79.6	53.1	61.8	48.4	66.2	14.9
PSKD [55]		ACM MM'22	68.8	78.6	56.0	64.5	50.2	66.2	15.0	29.8
SaA [36]		-	67.1	76.2	53.5	63.0	49.5	60.8	14.8	-
ZSE-RN [26]		CVPR'23	69.8	79.7	52.5	62.4	54.2	65.7	14.5	21.6
ZSE-Ret [26]		CVPR'23	73.6	80.8	50.4	60.2	56.9	63.7	14.2	20.2
CLIP	CLIP-AT [39]	CVPR'23	-	-	72.3	72.5	65.1	73.2	20.2	38.8
	TLT [59]	MMTA'23	77.9	84.3	66.1	73.0	61.5	69.5	27.8	-
	Sherry [8]	-	74.1	83.5	61.6	69.5	54.1	66.4	18.0	29.8
	MARL [29]	WACV'24	-	-	69.1	75.5	70.5	77.7	32.7	42.5
SPLIP			80.2	86.7	76.4	77.3	73.1	78.2	34.2	44.6

considers the top 200 retrieved photos, where we report the mean Average Precision score (mAP@all) and precision at 200 (P@200). Aligning with the recent trends, however, we specifically report precision at 100 (P@100) for the TU-Berlin-Ext dataset and mAP at 200 (mAP@200) for the Sketchy dataset. For FG-ZS-SBIR, accuracy is evaluated by considering only a single category at a time [33], denoted as $\text{Acc}@K$ for Sketchy. This metric reflects the percentage of sketches that have their true matched photo within the top- K list, with our focus being on the Top-1 and Top-5 accuracy metrics [39].

4.1 Comparison to the literature

- **(G)ZS-SBIR and FG-ZS-SBIR:** Table 1 provides a comparative analysis of SPLIP against methods utilizing ConvNet, ViT, and CLIP backbones across the extended versions of the two splits of the Sketchy, TU-Berlin, and QuickDraw datasets, respectively in categorical ZS-SBIR. Remarkably, SPLIP outperforms all models based on ConvNet and ViT backbones by a substantial margin across all metrics. For instance, on the challenging QuickDraw-Ext dataset, SPLIP surpasses the leading ConvNet-based method, BDA [5], by 18.8% and the premier ViT-based model, PSKD [55], by 19.2% in the mAP@all metric. In comparison with CLIP-based models, SPLIP achieves an improvement over MARL by 1.5% in mAP values, with analogous enhancements observed across all datasets, indicating at least a 2% boost in mAP scores. In Fig. 3, we show qualitative comparisons between our proposed SPLIP and [39] on ZS-SBIR, highlighting our improved retrieval results for different categories.

Further insights into the performance on more stringent GZS-SBIR tasks are provided in Table 2, where SPLIP significantly demonstrates its ability to alleviate model bias towards training classes. It outmatches the second-best approach by 4.8% on Sketchy-Ext and 4.1% on TU-Berlin-Ext in mAP scores, unequiv-

ocally evidencing enhanced generalization capabilities. Table 3 showcases the

Table 2: Comparison for the GZS-SBIR.

	Methods		Sketchy-2-Ext [57]		TU-Berlin-Ext [27]	
			mAP@200	P@200	mAP@all	P@100
CNN	SEM-PCYC [11]	CVPR'19	-	-	19.2	29.8
	OCEAN [66]	ICME'20	-	-	31.2	34.1
	BDA [5]	NC'22	22.6	33.7	25.1	35.7
ViT	SaA [36]	-	-	-	29.0	38.1
	ZSE-Ret [26]	CVPR'23	-	-	46.4	48.5
	ZSE-RN [26]	CVPR'23	-	-	43.2	46.0
	STL [15]	AAAI'23	63.4	53.8	40.2	49.8
CLIP	CLIP-AT [39]	CVPR'23	55.6	62.7	60.9	63.8
	MARL [29]	WACV'24	62.3	68.5	62.6	67.8
	SPLIP		68.2	74.5	66.7	70.3

Table 3: Comparison for the FG-ZS-SBIR on the Sketchy dataset.

Methods	Sketchy [44]	
	Acc@1	Acc@5
CrossGrad [45]	13.40	34.90
CC-DG [31]	22.60	49.00
SketchPVT [40]	30.24	51.65
CLIP-AT [39]	28.68	62.34
MARL [29]	29.96	58.53
SPLIP	33.45	66.71

FG-ZS-SBIR performance on the Sketchy dataset, comparing top-1 and top-5 retrieval accuracies. SPLIP attains a top-1 accuracy of 33.45%, leading the subsequent best method by 3.21%. This represents a nearly 4% improvement over the findings in [39]. Our foundational model, spotlighting multimodal prompting and $\mathcal{L}_{triplet} + \mathcal{L}_{class}$, outperforms the achievements of [39], highlighting the critical importance of multimodal prompting alone. Additionally, the adoption of \mathcal{L}_{cjs} enhances results by almost 3.4% in Acc@1.

- **Across dataset ZS-SBIR:** Leveraging on the across-dataset generalizability capabilities of CLIP, we assess the effectiveness of our proposed SPLIP on the ZS-SBIR across dataset setting by [26], where the model is trained on the Sketchy-Ext dataset and evaluated on 21 unseen classes of TU-Berlin-Ext and 11 unseen classes of Quickdraw-Ext datasets, respectively. In Table 10, our findings demonstrate that SPLIP outperforms CLIP-AT [39] in terms of mAP@all and P@100 metrics by 14.2% and 12.9% on the TU-Berlin-Ext dataset, and by 15.1% and 13.6% on the Quickdraw-Ext dataset.

4.2 Ablation analysis

Besides the following, ablations on the number of training samples, the ratio of known and unknown classes, etc., are mentioned in the **Supplementary**.

- **Analysis of the loss components:** In Table 5, we explore the influence of various loss components as described in Section 3.5 on the ZS-SBIR and FG-ZS-SBIR tasks, utilizing the Sketchy dataset for evaluation. Incorporation of $\mathcal{L}_{triplet}$ demonstrates a notable performance improvement, which is further sharply enhanced by \mathcal{L}_{class} . Adding \mathcal{L}_{cjs} contributes to an approximate 3 – 4% increase in results for both tasks. Specifically, we find that both the components of \mathcal{L}_{cjs} show improvements, *i.e.*, using $\mathcal{L}_{class} + \mathcal{L}_{triplet} + \mathcal{L}_{margin}$ improves the

Table 4: Comparison of ZS-SBIR across datasets while training with the Sketchy-Ext [57] dataset and tested on TU-Berlin-Ext and Quickdraw-Ext datasets. * represents the results reproduced by us.

Methods	TU-Berlin-Ext		Quickdraw-Ext	
	mAP@all	P@100	mAP@all	P@100
CC-DG [31]	30.8	43.4	15.6	22.7
DSN [56]	35.6	46.9	14.9	17.8
SAKE [28]	38.9	50.6	17.4	24.2
ZSE-RN [26]	47.6	59.0	22.8	33.8
CLIP-AT* [39]	56.4	63.1	30.7	45.0
SPLIP	70.6	76.0	45.8	58.6

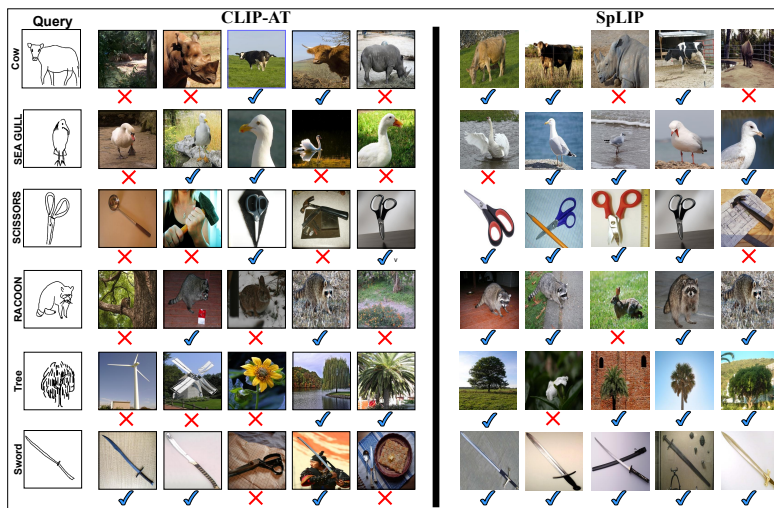


Fig. 3: Qualitative comparisons between SP-LIP and [39] in categorical ZS-SBIR. Improved retrieval outcomes for more ambiguous classes can be observed with SP-LIP.

results of $\mathcal{L}_{class} + \mathcal{L}_{triplet}$ by $\approx 1 - 2\%$, and the use of full \mathcal{L}_{cjs} offers further improvements by 2% for ZS-SBIR and $\approx 2 - 3\%$ for FG-ZS-SBIR.

- **Ablation with learnable blocks of SP-LIP:** In Table 6, we elucidate the significance of each learnable module within SP-LIP, \mathcal{B}_v , \mathcal{B}_t , \mathcal{B}_{vt} , and LayerNorm parameters, for both ZS-SBIR and FG-ZS-SBIR tasks. Note that excluding one of these blocks means the respective prompt token sharing is disabled. Consistent with findings from [39], rendering the parameters of LayerNorm modules learnable results in a performance boost of approximately 2%. Conversely, the absence of cross-modal token sharing detrimentally impacts outcomes, whether in one or both directions. Specifically, we note a decline in precision by at least 4% for ZS-SBIR and 3% for FG-ZS-SBIR tasks, underscoring the critical role of token sharing in enhancing model efficacy across both retrieval challenges.

- **Shallow to deep prompting in SP-LIP:** Additionally, we assess the effect of employing deep prompting across separate textual and visual branches, alongside our integrated multi-modal approach, as depicted in Figure 4 (Left). This evaluation reveals a general trend where mAP values ascend with the inclusion of more layers from $\mathcal{F}_v/\mathcal{F}_t$ for all three scenarios. Notably, the enhancements achieved through multi-modal prompting distinctly surpass those of the uni-modal counterparts, emphasizing the superior efficacy of cohesively leveraging both textual and visual cues.

- **SP-LIP coupled with the multi-modal prompting of [20], the fine-grained metric loss of [39], and with a fixed margin for $\mathcal{L}_{triplet}$:** Our approach, leveraging extensive bi-directional information sharing between visual and textual modalities, excels beyond current multi-modal prompting methods, including those by [20]. Contrary to [20]’s uni-directional and layer-restricted token sharing, our comprehensive integration strategy showcases superior efficacy, as depicted in Fig. 4 (Right) (marked as B1).

Table 5: Ablation of loss terms (Equation 11) on the Sketchy dataset.

Loss	ZS-SBIR		FG-ZS-SBIR	
	mAP@200	P@200	Acc@1	Acc@5
\mathcal{L}_{class}	57.5	58.1	17.23	39.41
$\mathcal{L}_{triplet}$	71.6	72.7	26.54	59.92
$\mathcal{L}_{class} + \mathcal{L}_{triplet}$	73.1	73.9	30.07	62.95
$\mathcal{L}_{class} + \mathcal{L}_{triplet} + \mathcal{L}_{margin}$	74.5	75.1	31.23	63.54
$\mathcal{L}_{class} + \mathcal{L}_{triplet} + \mathcal{L}_{cjs}$	76.4	77.3	33.45	66.71

Table 6: Ablation of learnable blocks of SPLIP on Sketchy dataset.

Method	ZS-SBIR		FG-ZS-SBIR	
	mAP@200	P@200	Acc@1	Acc@5
w/o LayerNorm	74.2	74.9	31.24	64.14
w/o B_v	71.9	73.6	29.76	62.58
w/o B_t	72.8	73.5	30.41	63.39
w/o B_{vt}	70.5	72.0	28.92	62.04
w/o ($B_v + B_{vt}$)	68.8	69.7	26.54	59.92
w/o ($B_v + B_t + B_{vt}$)	62.5	68.7	28.49	59.35
SPLIP (Ours)	76.4	77.3	33.45	66.71

We introduced an experiment to underscore the significance of matching patch-permuted sketches with unshuffled photos, thus capturing essential patch-level alignments for better modality congruence. Replacing \mathcal{L}_{cjs} and \mathcal{F}_{js} with [39]’s patch shuffling metric, where a positive pair of sketch-photo shares identical patch permutations and a negative pair does not, our method markedly outperforms the SPLIP combined with [39] across SBIR tasks, highlighting \mathcal{F}_{js} ’s robustness in SPLIP (Fig. 4 (Right)) (marked as B2).

Opting for a dynamic μ value in $\mathcal{L}_{triplet}$, over a static $\mu = 0.2$, significantly boosts performance across all tasks, demonstrating our method’s nuanced adaptability and effectiveness (Fig. 4 (Right)) (marked as B3).

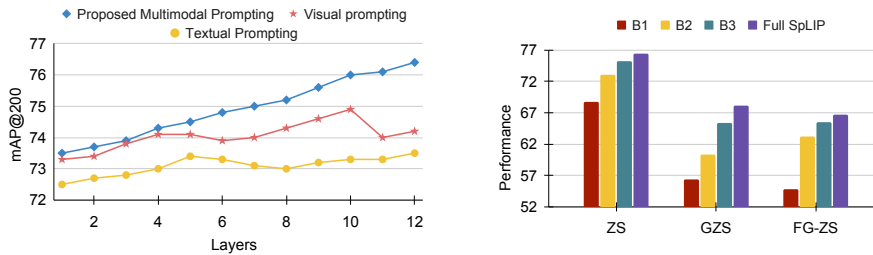


Fig. 4: **Left:** Effects of including more layers into our proposed deep prompting. We show the effects of prompting individually for CLIP’s vision and textual branches, followed by the proposed multi-modal prompting. **Right:** B1 refers effects of replacing the multi-modal prompting of ours by that of [20], B2 refers to replacing the conditional cross-modal jigsaw solver by the patch permutation-based triplet objective proposed in [39], and B3 denotes SPLIP with a fixed μ of 0.2 for $\mathcal{L}_{triplet}$.

5 Takeaways

This paper unveils SPLIP, a novel CLIP-based framework tailored for ZS, GZS, and FG-ZS SBIR tasks. Our methodology is distinguished by three innovations: a multi-modal prompt learning strategy enhancing cross-modal knowledge sharing and embedding learning within text-sketch-photo triads; the use of CLIP’s textual class-name embeddings to dynamically adjust margins for the sketch-photo triplet loss; and a unique conditional cross-modal jigsaw task designed to fine-tune patch-level sketch-photo associations, by repurposing the stand-alone jigsaw task with a metric objective. SPLIP’s efficacy is validated across various benchmarks, consistently showcasing its dominance in all tasks. Our future work will extend to broader vision-language domains, enriching our understanding of visual semantics. The authors sincerely acknowledge the tremendous support from AWL Inc, Japan.

References

1. Bhunia, A.K., Sain, A., Shah, P.H., Gupta, A., Chowdhury, P.N., Xiang, T., Song, Y.Z.: Adaptive fine-grained sketch-based image retrieval. In: European Conference on Computer Vision. pp. 163–181. Springer (2022) [4](#)
2. Bose, S., Jha, A., Fini, E., Singha, M., Ricci, E., Banerjee, B.: Stylip: Multi-scale style-conditioned prompt learning for clip-based domain generalization. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 5542–5552 (2024) [5](#)
3. Bulat, A., Tzimiropoulos, G.: Lasp: Text-to-text optimization for language-aware soft prompting of vision & language models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 23232–23241 (2023) [5](#)
4. Chaudhuri, A., Bhunia, A.K., Song, Y.Z., Dutta, A.: Data-free sketch-based image retrieval. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 12084–12093 (2023) [4](#)
5. Chaudhuri, U., Chavan, R., Banerjee, B., Dutta, A., Akata, Z.: Bda-sketret: Bi-level domain adaptation for zero-shot sbir. *Neurocomputing* **514**, 245–255 (2022) [2](#), [4](#), [10](#), [11](#), [12](#)
6. Devlin, J., Chang, M.W., Lee, K., Toutanova, K.: Bert: Pre-training of deep bidirectional transformers for language understanding. arXiv preprint arXiv:1810.04805 (2018) [5](#)
7. Dey, S., Riba, P., Dutta, A., Lladós, J., Song, Y.Z.: Doodle to search: Practical zero-shot sketch-based image retrieval. In: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. pp. 2179–2188 (2019) [4](#), [10](#), [11](#), [21](#)
8. Dong, S., Zhu, M., Wang, N., Yang, H., Gao, X.: Adapt and align to improve zero-shot sketch-based image retrieval. arXiv preprint arXiv:2305.05144 (2023) [4](#), [11](#)
9. Dosovitskiy, A., Beyer, L., Kolesnikov, A., Weissenborn, D., Zhai, X., Unterthiner, T., Dehghani, M., Minderer, M., Heigold, G., Gelly, S., Uszkoreit, J., Houlsby, N.: An image is worth 16x16 words: Transformers for image recognition at scale. CoRR [abs/2010.11929](https://arxiv.org/abs/2010.11929) (2020), <https://arxiv.org/abs/2010.11929> [5](#), [6](#)
10. Dowson, D., Landau, B.: The fréchet distance between multivariate normal distributions. *Journal of multivariate analysis* **12**(3), 450–455 (1982) [20](#)
11. Dutta, A., Akata, Z.: Semantically tied paired cycle consistency for zero-shot sketch-based image retrieval. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 5089–5098 (2019) [2](#), [4](#), [10](#), [11](#), [12](#), [20](#)
12. Dutta, T., Singh, A., Biswas, S.: Styleguide: Zero-shot sketch-based image retrieval using style-guided image generation. *IEEE Transactions on Multimedia* **23**, 2833–2842 (2020) [11](#)
13. Eitz, M., Hays, J., Alexa, M.: How do humans sketch objects? *ACM Transactions on graphics (TOG)* **31**(4), 1–10 (2012) [4](#), [20](#)
14. Eitz, M., Hildebrand, K., Boubekur, T., Alexa, M.: Sketch-based image retrieval: Benchmark and bag-of-features descriptors. *IEEE transactions on visualization and computer graphics* **17**(11), 1624–1636 (2010) [2](#)
15. Ge, C., Wang, J., Qi, Q., Sun, H., Xu, T., Liao, J.: Semi-transductive learning for generalized zero-shot sketch-based image retrieval. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 37, pp. 7678–7686 (2023) [4](#), [12](#)

16. Gupta, S., Chaudhuri, U., Banerjee, B., Kumar, S.: Zero-shot sketch based image retrieval using graph transformer. In: 2022 26th International Conference on Pattern Recognition (ICPR). pp. 1685–1691. IEEE (2022) [4](#)
17. Ha, D., Eck, D.: A neural representation of sketch drawings. arXiv preprint arXiv:1704.03477 (2017) [4](#), [21](#)
18. Hu, R., Collomosse, J.: A performance evaluation of gradient field hog descriptor for sketch based image retrieval. *Computer Vision and Image Understanding* **117**(7), 790–806 (2013) [4](#)
19. Jia, C., Yang, Y., Xia, Y., Chen, Y.T., Parekh, Z., Pham, H., Le, Q., Sung, Y.H., Li, Z., Duerig, T.: Scaling up visual and vision-language representation learning with noisy text supervision. In: International Conference on Machine Learning. pp. 4904–4916. PMLR (2021) [2](#)
20. Khattak, M.U., Rasheed, H., Maaz, M., Khan, S., Khan, F.S.: Maple: Multi-modal prompt learning. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 19113–19122 (2023) [2](#), [3](#), [5](#), [7](#), [8](#), [13](#), [14](#)
21. Khattak, M.U., Wasim, S.T., Naseer, M., Khan, S., Yang, M.H., Khan, F.S.: Self-regulating prompts: Foundational model adaptation without forgetting. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 15190–15200 (2023) [2](#), [3](#), [5](#)
22. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv preprint arXiv:1412.6980 (2014) [10](#)
23. Li, L.H., Yatskar, M., Yin, D., Hsieh, C.J., Chang, K.W.: Visualbert: A simple and performant baseline for vision and language. arXiv preprint arXiv:1908.03557 (2019) [5](#)
24. Li, Y., Hospedales, T., Song, Y.Z., Gong, S.: Fine-grained sketch-based image retrieval by matching deformable part models. *BMVC 2014 - Proceedings of the British Machine Vision Conference 2014* (01 2014) [4](#)
25. Li, Y., Li, W.: A survey of sketch-based image retrieval. *Machine Vision and Applications* **29**(7), 1083–1100 (2018) [2](#)
26. Lin, F., Li, M., Li, D., Hospedales, T., Song, Y.Z., Qi, Y.: Zero-shot everything sketch-based image retrieval, and in explainable style. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 23349–23358 (2023) [4](#), [11](#), [12](#), [23](#)
27. Liu, L., Shen, F., Shen, Y., Liu, X., Shao, L.: Deep sketch hashing: Fast free-hand sketch-based image retrieval. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 2862–2871 (2017) [4](#), [10](#), [11](#), [12](#), [20](#)
28. Liu, Q., Xie, L., Wang, H., Yuille, A.L.: Semantic-aware knowledge preservation for zero-shot sketch-based image retrieval. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 3662–3671 (2019) [2](#), [4](#), [11](#), [12](#), [23](#)
29. Lyou, E., Lee, D., Kim, J., Lee, J.: Modality-aware representation learning for zero-shot sketch-based image retrieval. In: Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision. pp. 5646–5655 (2024) [2](#), [4](#), [11](#), [12](#)
30. Noroozi, M., Favaro, P.: Unsupervised learning of visual representations by solving jigsaw puzzles. In: European conference on computer vision. pp. 69–84. Springer (2016) [3](#), [9](#)
31. Pang, K., Li, K., Yang, Y., Zhang, H., Hospedales, T.M., Xiang, T., Song, Y.Z.: Generalising fine-grained sketch-based image retrieval. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 677–686 (2019) [2](#), [11](#), [12](#), [23](#)

32. Pang, K., Song, Y.Z., Xiang, T., Hospedales, T.M.: Cross-domain generative learning for fine-grained sketch-based image retrieval. In: *BMVC*. pp. 1–12 (2017) [4](#)
33. Pang, K., Yang, Y., Hospedales, T.M., Xiang, T., Song, Y.Z.: Solving mixed-modal jigsaw puzzle for fine-grained sketch-based image retrieval. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 10347–10355 (2020) [3](#), [5](#), [8](#), [11](#)
34. Radford, A., Kim, J.W., Hallacy, C., Ramesh, A., Goh, G., Agarwal, S., Sastry, G., Askell, A., Mishkin, P., Clark, J., et al.: Learning transferable visual models from natural language supervision. In: *International conference on machine learning*. pp. 8748–8763. PMLR (2021) [2](#), [5](#)
35. Radford, A., Narasimhan, K., Salimans, T., Sutskever, I., et al.: Improving language understanding by generative pre-training (2018) [5](#)
36. Ribeiro, L.S.F., Ponti, M.A.: Sketch-an-anchor: Sub-epoch fast model adaptation for zero-shot sketch-based image retrieval. *arXiv preprint arXiv:2303.16769* (2023) [4](#), [11](#), [12](#)
37. Roy, S., Etemad, A.: Consistency-guided prompt learning for vision-language models. *arXiv preprint arXiv:2306.01195* (2023) [2](#), [3](#), [5](#), [7](#)
38. Russakovsky, O., Deng, J., Su, H., Krause, J., Satheesh, S., Ma, S., Huang, Z., Karpathy, A., Khosla, A., Bernstein, M., et al.: Imagenet large scale visual recognition challenge. *International journal of computer vision* **115**, 211–252 (2015) [20](#)
39. Sain, A., Bhunia, A.K., Chowdhury, P.N., Koley, S., Xiang, T., Song, Y.Z.: Clip for all things zero-shot sketch-based image retrieval, fine-grained or not. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 2765–2775 (2023) [2](#), [3](#), [4](#), [5](#), [8](#), [9](#), [11](#), [12](#), [13](#), [14](#), [20](#), [23](#), [24](#)
40. Sain, A., Bhunia, A.K., Koley, S., Chowdhury, P.N., Chattopadhyay, S., Xiang, T., Song, Y.Z.: Exploiting unlabelled photos for stronger fine-grained sbir. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 6873–6883 (2023) [5](#), [12](#)
41. Sain, A., Bhunia, A.K., Potlapalli, V., Chowdhury, P.N., Xiang, T., Song, Y.Z.: Sketch3t: Test-time training for zero-shot sbir. In: *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. pp. 7462–7471 (2022) [11](#)
42. Sain, A., Bhunia, A.K., Yang, Y., Xiang, T., Song, Y.Z.: Cross-modal hierarchical modelling for fine-grained sketch based image retrieval. *arXiv preprint arXiv:2007.15103* (2020) [4](#)
43. Sain, A., Bhunia, A.K., Yang, Y., Xiang, T., Song, Y.Z.: Stylemeup: Towards style-agnostic sketch-based image retrieval. In: *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*. pp. 8504–8513 (2021) [2](#), [4](#)
44. Sangkloy, P., Burnell, N., Ham, C., Hays, J.: The sketchy database: learning to retrieve badly drawn bunnies. *ACM Transactions on Graphics (TOG)* **35**(4), 1–12 (2016) [4](#), [10](#), [12](#), [20](#)
45. Shankar, S., Piratla, V., Chakrabarti, S., Chaudhuri, S., Jyothi, P., Sarawagi, S.: Generalizing across domains via cross-gradient training. *arXiv preprint arXiv:1804.10745* (2018) [12](#)
46. Shen, Y., Liu, L., Shen, F., Shao, L.: Zero-shot sketch-image hashing. In: *Proceedings of the IEEE conference on computer vision and pattern recognition*. pp. 3598–3607 (2018) [20](#)
47. Singha, M., Jha, A., Banerjee, B.: Gopro: Generate and optimize prompts in clip using self-supervised learning. *arXiv preprint arXiv:2308.11605* (2023) [5](#)

48. Singha, M., Jha, A., Bose, S., Nair, A., Abdar, M., Banerjee, B.: Unknown prompt the only lacuna: Unveiling clip’s potential for open domain generalization. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 13309–13319 (2024) [5](#)
49. Singha, M., Pal, H., Jha, A., Banerjee, B.: Ad-clip: Adapting domains in prompt space using clip. In: Proceedings of the IEEE/CVF International Conference on Computer Vision. pp. 4355–4364 (2023) [5](#)
50. Song, J., Yu, Q., Song, Y.Z., Xiang, T., Hospedales, T.M.: Deep spatial-semantic attention for fine-grained sketch-based image retrieval. In: Proceedings of the IEEE international conference on computer vision. pp. 5551–5560 (2017) [4](#)
51. Sun, H., Xu, J., Wang, J., Qi, Q., Ge, C., Liao, J.: Dli-net: Dual local interaction network for fine-grained sketch-based image retrieval. *IEEE Transactions on Circuits and Systems for Video Technology* **32**(10), 7177–7189 (2022) [5](#)
52. Tian, J., Xu, X., Shen, F., Yang, Y., Shen, H.T.: Tvt: Three-way vision transformer through multi-modal hypersphere learning for zero-shot sketch-based image retrieval. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 36, pp. 2370–2378 (2022) [2](#), [11](#)
53. Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A.N., Kaiser, Ł., Polosukhin, I.: Attention is all you need. *Advances in neural information processing systems* **30** (2017) [6](#)
54. Wang, H., Deng, C., Liu, T., Tao, D.: Transferable coupled network for zero-shot sketch-based image retrieval. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **44**(12), 9181–9194 (2021) [11](#)
55. Wang, K., Wang, Y., Xu, X., Liu, X., Ou, W., Lu, H.: Prototype-based selective knowledge distillation for zero-shot sketch based image retrieval. In: Proceedings of the 30th ACM International Conference on Multimedia. pp. 601–609 (2022) [2](#), [10](#), [11](#)
56. Wang, Z., Wang, H., Yan, J., Wu, A., Deng, C.: Domain-smoothing network for zero-shot sketch-based image retrieval. *arXiv preprint arXiv:2106.11841* (2021) [11](#), [12](#), [23](#)
57. Yelamarthi, S.K., Reddy, S.K., Mishra, A., Mittal, A.: A zero-shot framework for sketch based image retrieval. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 300–317 (2018) [4](#), [10](#), [11](#), [12](#), [20](#)
58. Yu, Q., Liu, F., Song, Y.Z., Xiang, T., Hospedales, T.M., Loy, C.C.: Sketch me that shoe. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 799–807 (2016) [2](#), [4](#), [5](#)
59. Zhang, H., Cheng, D., Jiang, H., Liu, J., Kou, Q.: Task-like training paradigm in clip for zero-shot sketch-based image retrieval. *Multimedia Tools and Applications* pp. 1–18 (2023) [2](#), [4](#), [11](#)
60. Zhang, H., Liu, S., Zhang, C., Ren, W., Wang, R., Cao, X.: Sketchnet: Sketch classification with web images. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 1105–1113 (2016) [20](#)
61. Zhang, X., Li, X., Liu, Y., Feng, F.: A survey on freehand sketch recognition and retrieval. *Image and Vision Computing* **89**, 67–87 (2019) [2](#)
62. Zhang, Z., Zhang, Y., Feng, R., Zhang, T., Fan, W.: Zero-shot sketch-based image retrieval via graph convolution network. In: Proceedings of the AAAI Conference on Artificial Intelligence. vol. 34, pp. 12943–12950 (2020) [4](#)
63. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Conditional prompt learning for vision-language models. In: Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition. pp. 16816–16825 (2022) [3](#), [5](#)

64. Zhou, K., Yang, J., Loy, C.C., Liu, Z.: Learning to prompt for vision-language models. *International Journal of Computer Vision* **130**(9), 2337–2348 (2022) [5](#)
65. Zhou, R., Chen, L., Zhang, L.: Sketch-based image retrieval on a large scale database. In: *Proceedings of the 20th ACM international conference on Multimedia*. pp. 973–976 (2012) [4](#)
66. Zhu, J., Xu, X., Shen, F., Lee, R.K.W., Wang, Z., Shen, H.T.: Ocean: A dual learning approach for generalized zero-shot sketch-based image retrieval. In: *2020 IEEE International Conference on Multimedia and Expo (ICME)*. pp. 1–6. IEEE (2020) [4](#), [11](#), [12](#)

A Contents of the supplementary materials

In this supplementary document, we present detailed information and further experimental results, including:

1. **Dataset description and splits for zero-shot SBIR settings:** In Section B, we provide the detailed descriptions and splits of the datasets used in our proposed work.
2. **List of notations used in proposed SPLIP:** In Table 7, we list the important variables used in SPLIP together with their significance.
3. **Ablation with different number of tokens:** Table 8 signifies the ablation with varying number of tokens produced by \mathcal{B}_t , \mathcal{B}_v and \mathcal{B}_{vt} .
4. **Ablation with varying the hard prompts:** In Table 9, we use different hard prompts for the generation of tokens from \mathcal{B}_v .
5. **Another cross-dataset ZS-SBIR experiment:** We show results of another cross-dataset ZS-SBIR experiment in Table 10.
6. **Ablation on varying the number of training samples and seen classes:** They are reported in Fig. 5-6.
7. **Analysis of the alignment of the photo and sketch domains:** We compare the same against [39]. We show the t-SNE plots for both the visual domains, as produced by [39] and SPLIP (Fig. 7-8), as well as show the domain distances produced by both the methods, in terms of the Fréchet distance [10] (Table 11).
8. **Further analysis on the jigsaw task:** We compare the performance of SPLIP on the jigsaw tasks applied on all (s_a, p_a^+, p_a^-) sampled from the same visual domain, *e.g.* all from photo or all from sketch and the proposed cross-modal jigsaw. The results are reported in Table 12.
9. **Further qualitative results:** They are shown in Fig. 9.

B Dataset descriptions

In this work, we evaluate the SPLIP method’s effectiveness across leading ZS-SBIR datasets, including Sketchy-Ext, TU-Berlin-Ext, and QuickDraw-Ext.

Sketchy-Ext [11, 57]: This dataset enlarges the original Sketchy collection [44] to 73,002 images across 125 categories, with an average of 604 sketches and 584 images per category. For zero-shot analysis, two distinct splits are employed. In split-1, *i.e.*, Sketchy-1-Ext [11], 25 categories are randomly chosen for testing, leaving the remaining 100 for training. On the other hand, Sketchy-2-Ext [57] isolates 21 classes not present in the ImageNet [38] dataset for testing, with the other 104 classes designated for training.

TU-Berlin-Ext [27]: An augmentation of the original TU-Berlin [13] dataset, which included 20,000 sketches over 250 categories. This extension incorporates 204,489 natural images [60], averaging 787 images per category, though with notable class-wise image imbalances. Adhering to the partitioning protocol from [46], we allocate 30 classes for testing and 220 for training, ensuring each test class has a minimum of 400 images to mitigate the imbalance.

QuickDraw-Ext [7]: A large-scale dataset designed for ZS-SBIR, extending the Google QuickDraw dataset [17] with detailed photographs. It includes 110 categories, featuring 330,000 sketches (3,000 per category) and 204,000 Flickr-sourced images, each tagged appropriately. The partitioning strategy from [7] is used, segregating 30 test classes not found in ImageNet, with the remaining 80 classes utilized for training.

C List of important variables and their significance

Table 7: Summarizing the variables used.

Notations	Description
$\mathcal{F}_v, \mathcal{F}_t$	Frozen CLIP’s visual and text encoders
\mathcal{F}_{js}	Transformer-based jigsaw-solver network
\mathcal{L}	Number of encoder layers in \mathcal{F}_v and \mathcal{F}_t
s_a, s'_a	Anchor Sketch, Permuted Anchor Sketch
p_a^+, p_a^-	Positive Photo, Negative Photo, given the anchor sketch
$\mathcal{B}_t, \mathbf{T}$	Vision-guided deep textual prompting block, and the obtained tokens
$\mathcal{B}_v, \mathbf{V}^{sg}$	Text-guided deep visual prompting block, and the respective semantic domain knowledge defined tokens
$\mathcal{B}_{vt}, \mathbf{V}^{ms}$	Vision-text conjunction block, along with the respective tokens
\mathcal{M}, d_t	Number of embedding tokens and dimension of text tokens
$\mathcal{J}, \mathbf{W}_0$	Number of words in input texts and corresponding word embeddings
$\mathbf{E}_0, \mathcal{N}$	Initial patch embeddings, number of patch embeddings
$\mathbb{1}$	One-hot label vector
θ, ϕ	LayerNorm parameters
\mathcal{L}_{ce}	Cross-entropy loss
\mathcal{L}_{class}	Text-image classification loss
$\mathcal{L}_{triplet}$	Cross-visual modality triplet loss
\mathcal{L}_{margin}	Margin loss proposed within \mathcal{L}_{cjs}
\mathcal{L}_{cjs}	Conditional cross-modal jigsaw loss

D Ablation with different number of tokens produced by \mathcal{B}_t , \mathcal{B}_v and \mathcal{B}_{vt}

We conducted a comprehensive analysis of our SPLIP framework by adjusting the number of tokens generated from the vision-guided deep textual prompting block (\mathcal{B}_t), the text-guided deep visual prompting block (\mathcal{B}_v), and the vision-text conjunction block (\mathcal{B}_{vt}). In Table 8, the configuration of token settings is denoted in the order: $(\mathcal{B}_t, \mathcal{B}_v, \mathcal{B}_{vt})$. We undertook ablation studies across three SBIR tasks: ZS-SBIR, GZS-SBIR, and FG-ZS-SBIR, to gauge the impact of token variation on performance.

The findings from these ablation studies highlight a trend where an increase in the number of tokens correlates with improved performance. Specifically, the configuration with tokens set to (4, 4, 2) achieved superior results for FG-ZS-SBIR tasks on the Sketchy-Ext dataset. Moreover, it was observed that employing a higher number of prompts tends to diminish the raw feature representation, affecting the overall performance of the system. This observation underscores the

balance required between the number of prompts and the preservation of feature quality for optimal performance.

Table 8: Ablation of number of tokens chosen for all types of prompting. $(.,.,.)$ signifies number of tokens produced by $\mathcal{B}_t, \mathcal{B}_v, \mathcal{B}_{vt}$ respectively.

Tokens $(\mathcal{B}_t, \mathcal{B}_v, \mathcal{B}_{vt})$	ZS-SBIR		GZS-SBIR		FG-ZS-SBIR	
	mAP@200	P@200	mAP@200	P@200	Acc@1	Acc@5
(1,1,1)	71.6	72.4	60.4	65.8	26.54	56.71
(1,4,1)	73.1	73.7	64.5	69.8	29.63	60.92
(4,1,1)	74.5	74.9	66.0	71.9	31.85	63.07
(4,4,1)	76.2	77.0	68.4	74.6	33.34	66.59
(4,4,2)	76.4	77.3	68.2	74.5	33.45	66.71
(4,4,3)	75.9	77.1	68.1	74.3	33.24	66.43
(4,4,4)	75.3	76.4	67.4	73.6	32.85	65.96

E Ablation with varying the hard textual prompts

In Table 9, we examine the impact of employing hard prompting strategies on the performance of the three considered SBIR tasks.

Table 9: Ablation of varying hard textual prompts.

Hard Prompt	ZS-SBIR		GZS-SBIR		FG-ZS-SBIR	
	mAP@200	P@200	mAP@200	P@200	Acc@1	Acc@5
<i>photo/sketch of a</i> [CLS]	76.4	77.3	68.4	74.6	31.07	64.36
<i>an image of a</i> [CLS]	75.6	76.3	67.2	73.6	33.14	65.97
<i>visual representation of</i> [CLS]	75.2	76.1	67.0	73.3	33.45	66.71

F More cross-dataset ZS-SBIR results

Pl. see Table 10.

G Ablation on the number of training samples, and the ratio between seen and unseen classes

Pl. refer to Fig. 5 and Fig. 6 for the same.

Table 10: Comparison of ZS-SBIR across datasets while training with the TU-Berlin-Ext dataset and tested on Sketchy-Ext and Quickdraw-Ext datasets. * represents the results reproduced by us.

Methods	Sketchy-Ext		Quickdraw-Ext	
	mAP@all	P@100	mAP@all	P@100
CC-DG [31]	62.4	69.3	23.1	29.6
DSN [56]	61.3	65.4	21.8	24.6
SAKE [28]	62.6	70.1	23.5	31.8
ZSE-RN [26]	74.6	81.6	27.3	37.6
CLIP-AT* [39]	78.2	85.9	33.6	43.5
SpLIP	84.1	90.3	37.0	47.4

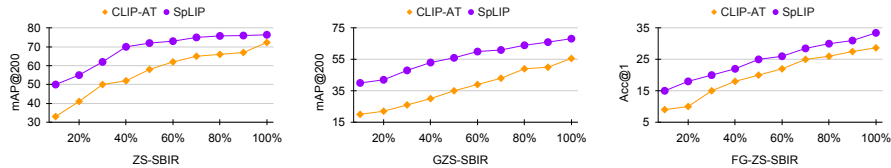


Fig. 5: Performance of [39] and SpLIP with varying training data size for Sketchy-Ext dataset.

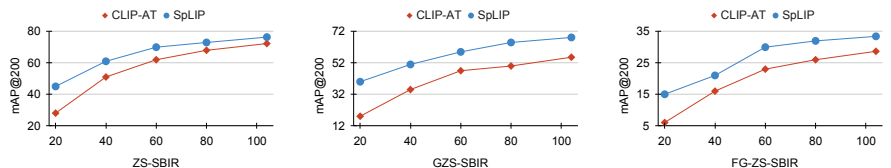


Fig. 6: Performance of [39] and SpLIP with varying number of seen classes while training for Sketchy-Ext dataset.

H Alignment of the sketch and photo domains

We present t-SNE visualizations for both modalities, generated by [39] and SpLIP, in Figures 7 and 8, respectively. These visualizations demonstrate that SpLIP achieves superior class separability. Additionally, we assess the domain alignment by comparing the Fréchet distance between the modalities in the embedding space, as detailed in Table 11. Our results indicate that SpLIP achieves a lower Fréchet distance, signifying enhanced alignment between domains.

I Ablation on the jigsaw task

In Table 12, we examine the performance sensitivity of SpLIP for the cross-modal conditional jigsaw task. This assessment includes a comparison against a

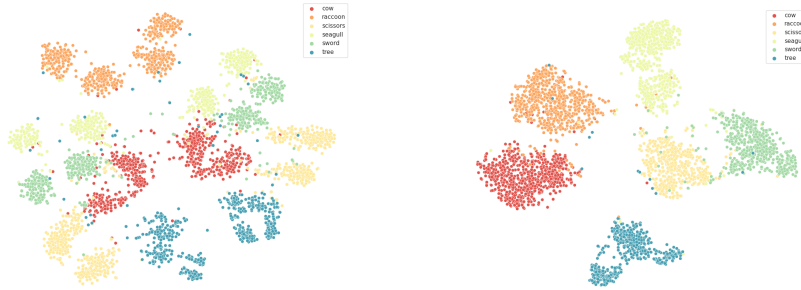


Fig. 7: t-sne of sketch domain, by [39] and SPLIP.



Fig. 8: t-sne of photo domain, by [39] and SPLIP.

Table 11: Fréchet distance between sketch and photo in the embedding space on Sketchy-ext dataset.

Method	ZS-SBIR	GZS-SBIR	FG-ZS-SBIR
CLIP-AT [39]	0.528	0.617	0.725
SPLIP	0.443	0.485	0.634

baseline variant where the anchor, along with its positive and negative examples, originate from a single modality, specifically either sketches or photos. Our findings reveal that cross-modal conditioning outperforms the uni-modal approach by a margin of 3-4% across all datasets examined.

J More qualitative results for FG-ZS-SBIR

Pl see Fig. 9, for results produced by SPLIP (ours) and [39].

Table 12: Ablation on the jigsaw task. We compare the cross-modal triplets against within-modality triplet selection.

\mathcal{L}_{cjs}	ZS-SBIR		GZS-SBIR		FG-ZS-SBIR	
	mAP@200	P@200	mAP@200	P@200	Acc@1	Acc@5
uni-modal	73.2	74.0	64.9	71.3	29.82	63.05
cross-modal	76.4	77.3	68.4	74.6	33.45	66.71



Fig. 9: Retrieved photos for sketch query instances for FG-ZS-SBIR.