

ViTAA: Visual-Textual Attributes Alignment in Person Search by Natural Language

Zhe Wang^{1*}, Zhiyuan Fang^{2*}, Jun Wang¹, and Yezhou Yang²

¹ Beihang University
{wangzhewz,wangj203}@buaa.edu.cn

² Arizona State University
{zfang29,yz.yang}@asu.edu

Abstract. Person search by natural language aims at retrieving a specific person in a large-scale image pool that matches given textual descriptions. While most of the current methods treat the task as a holistic visual and textual feature matching one, we approach it from an **attribute-aligning perspective** that allows grounding specific attribute phrases to the corresponding visual regions. We achieve success as well as a performance boost by a robust feature learning that the referred identity can be accurately bundled by multiple attribute cues. To be concrete, our Visual-Textual Attribute Alignment model (dubbed as ViTAA) learns to disentangle the feature space of a person into sub-spaces corresponding to attributes using a light auxiliary attribute segmentation layer. It then aligns these visual features with the textual attributes parsed from the sentences via a novel contrastive learning loss. We validate our ViTAA framework through extensive experiments on tasks of person search by natural language and by attribute-phrase queries, on which our system achieves state-of-the-art performances. Codes and models are available at <https://github.com/Jarr0d/ViTAA>.

Keywords: Person Search by Natural Language, Person Re-identification, Vision and Language, Metric Learning

1 Introduction

Recently, we have witnessed numerous practical breakthroughs in person modeling related tasks, *e.g.*, pedestrian detection [2,5,47], pedestrian attribute recognition [30,40] and person re-identification [13,26,59]. Person search [25,15] as an aggregation of the aforementioned tasks thus gains increasing research attention. Comparing with searching by image queries, person search by natural language [25,24,6,52] makes the retrieving procedure more user-friendly with increased flexibility due to a supporting of open-form natural language queries.

* Equal contribution. This work was done when Z. Wang was a visiting scholar at Active Perception Group, Arizona State University.

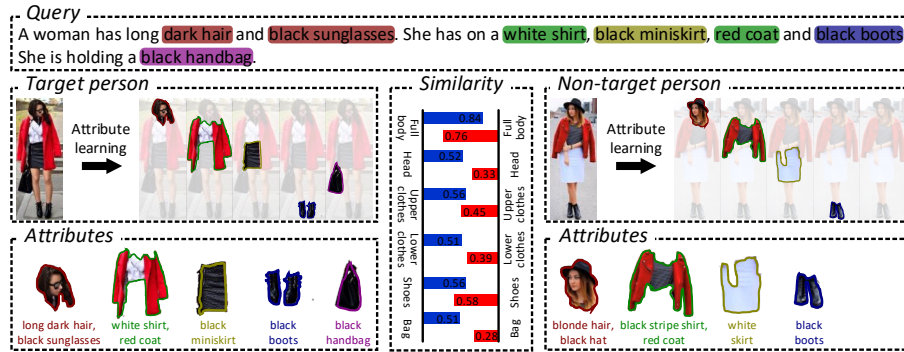


Fig. 1. In a case when two persons exhibit similar appearance attributes, it is hard to discriminate them merely by full-body appearance. Instead of matching the textual descriptions with the images at global-level, we decompose both image and text into attribute components and conduct a fine-grained matching strategy.

Meanwhile, learning robust visual-textual associations becomes increasingly critical, which calls an urgent demand for a representation learning schema that is able to fully exploit both modalities.

Relevant studies in person modeling related research points out the critical role of the discriminative representations, especially of the local fragments in both image and text. For the former, [38,58] learn the pose-related features from the key points map of human, while [20,27] leverage the body-part features by auxiliary segmentation-based supervision. For the latter, [25,24,55] decompose the complex sentences into noun phrases, and [52,23] directly adopt the attribute-specific annotations to learn fine-grained attribute related features. Moving forward, attribute specific features from image and text are even requisite for person search by natural language task, and how to effectively couple them stays an open question. We seek insight from a fatally flawed case that lingers in most of the current visual-language systems in Figure 1, termed as “malpositioned matching”. For example, tasks like textual grounding [35,33], VQA [1], and image-text retrieval [36,18] are measuring the similarities or mutual information across modalities in a holistic fashion by answering: are the feature vectors of image and text match with each other? That way, when users input “a girl in white shirt and black skirt” as retrieval query, the model is not able to distinguish the nuances of the two images as shown in Figure 1, where the false positive one actually shows “black shirt and white skirt”. As both the distinct color visual cues (“white” and “black”) exist in the images, overall matching without the ability of referring them to specific appearance attributes prevents the model from discriminating them as needed. Such cases exist extensively in almost all cross-modal tasks, which pose an indispensable challenge for a system to tackle with the ability of fine-grained interplay between image and text.

Here, we put forward a novel **Visual-Textual Attributes Alignment** model (ViTAA). For feature extraction, we fully exploit both visual and textual at-

tribute representations. Specifically, we leverage segmentation labels to drive the attribute-aware feature learning from the input image. As shown in Figure 3, we design multiple local branches, each of which is responsible to predict one particular attribute visual feature. This process is guided by the supervision on segmentation annotations, so the features are intrinsically aligned through the label information. We then use a generic natural language parser to extract attribute-related phrases, which at the same time remove the complex syntax in natural language and redundant non-informative descriptions. Building upon this, we adopt a contrastive learning schema to learn a joint embedding space for both visual and textual attributes. Meanwhile, we also notice that there may exist common attributes across different person identities (*e.g.*, two different persons may wear similar “*black shirt*”). To thoroughly exploits these cases during training, we propagate a novel sampling method to mine surrogate positive examples which largely enriches the sampling space, and also provides us with valid informative samples for the sake of overcoming convergence problem in metric learning.

To this end, we argue and show that the benefits of the attribute alignment to person search model go well beyond the obvious. As the images used for person search tasks often contain a large variance on appearance (*e.g.*, varying poses or viewpoints, with/without occlusion, and with cluttered background), the abstracted attribute-specific features could naturally help to resolve the ambiguity in feature representations. Also, searching by appearance attributes innately brings interpretability for the retrieving task and enables the attribute specific retrieval. It is worth mentioning that, there exist few very recent efforts that attempt to utilize the local fragments in both visual and textual modalities [8,49] and hierarchically align them [6,3]. The pairing schema of visual features and textual phrases in these methods are all based on the same identity, where they neglect the cues that exist across different identities. Comparing with them, ours is a more comprehensive modeling method that fully exploits the identical attributes from different persons thus greatly helps the alignment learning.

To validate these speculations, extensive experiments are conducted to our ViTAA model on the task of 1) person search by natural language and 2) by attribute, showing that our proposed model is capable of linking specific visual cues with specific words/phrases. More concretely, our ViTAA achieves a promising results across all these tasks. Further qualitative analysis verifies that our alignment learning successfully learns the fine-grained level correspondence across the visual and textual attributes. To summarize our contributions:

- We design an attribute-aware representation learning framework to extract and align both visual and textual features for the task of person search by natural language. To the best of our knowledge, this is the first to adopt both semantic segmentation as well as natural language parsing to facilitate a semantically aligned representation learning.
- We design a novel cross-modal alignment learning schema based on contrastive learning which can adaptively highlight the informative samples during the alignment learning. Meanwhile, an unsupervised data sampling

method is proposed to facilitate the construction of learning pairs by exploiting more surrogate positive samples across different person identities.

- We evaluate the superiority of ViTAA over other state-of-the-art methods for the person search by natural language task. We also conduct qualitative analysis to demonstrate the interpretability of ViTAA.

2 Related Work

Person Search. Given the form of the querying data, current person search tasks can be categorized into two major thrusts: searching by images (termed as Person Re-Identification), and person search by textual descriptions. Typical person re-identification (re-id) methods [13,26,59] are formulated as retrieving the candidate that shows highest correlation with the query in the image galleries. However, a clear and valid image query is not always available in the real scenario, thus largely impedes the applications of re-id tasks. Recently, researchers alter their attention to re-id by textual descriptions: identifying the target person by using free-form natural languages [25,24,3]. Meanwhile, it also comes with great challenges as requiring the model to deal with the complex syntax from the long and free-form descriptive sentence, and the inconsistent interpretations of low-quality surveillance images. To tackle these, methods like [25,24,4] employ attention mechanism to build relation module between visual and textual representations, while [55,60] propose cross-modal objective functions for joint embedding learning. Dense visual feature is extracted in [32] by cropping the input image for learning a regional-level matching schema. Beyond this, [19] introduces pose estimation information for delicate human body-part parsing.

Attribute Representations. Adopting appropriate feature representations is of crucial importance for learning and retrieving from both image and text. Previous efforts in person search by natural language unanimously use holistic features of the person, which omit the partial visual cues from attributes at the fine-grained level. Multiple re-id systems have focused on the processing of body-part regions for visual feature learning, which can be summarized as: hand-craft horizontal stripes or grid [26,43,46], attention mechanism [37,45], and auxiliary information including keypoints [50,41], human parsing mask [20,27] and dense semantic estimation [56]. Among these methods, the auxiliary information usually provides more accurate partition results on localizing human parts and facilitating body-part attribute representations thanks to the multi-task training or the auxiliary networks. However, only few work [14] pay attention to the accessories (such as the backpack) which could be the potential contextual cues for accurate person retrieval. As the corresponding components to specific visual cues, textual attribute phrases are usually provided as ground-truth labels or can be extracted from sentences through identifying the noun phrases with sentence parsing. Many of them use textual attributes as auxiliary label information to complement the content of image features [23,39,29]. Recently, a few attempts leverage textual attribute as query for person retrieval [6,52]. [52] imposes an attribute-guided attention mechanism to capture the holistic appearance of per-

son. [6] proposes a hierarchical matching model that can jointly learn global category-level and local attribute-level embedding.

Visual-Semantic Embedding. Works in vision and language propagate the notion of visual semantic embedding, with a goal to learn joint feature space for both visual inputs and their correspondent textual annotations [10,53]. Such mechanism plays a core role in a series of cross-modal tasks, *e.g.*, image/video captioning [21,51,7], image retrieval through natural language [55,49,8], and vision question answering [1]. Conventional joint embedding learning framework adopts two-branch architecture [55,10,53,8,9], where one extracts image features and the other encodes textual descriptions. The extracted cross-modal embedding features are learned through carefully designed objective functions.

3 Our Approach

Our network is composed of an image stream and a language stream (see Figure 3), with the intention to encode inputs from both modalities for a visual-textual embedding learning. To be specific, given a person image \mathcal{I} and its textual description \mathcal{T} , we first use the image stream to extract a global visual representation \mathbf{v}_0 , and a stack of local visual representations of N_{att} attributes $\{\mathbf{v}_1, \dots, \mathbf{v}_{N_{att}}\}$, $\mathbf{v}_i \in \mathbb{R}^d$. Similarly, we follow the language stream to extract overall textual embedding \mathbf{t}_0 , then decompose the whole sentence using standard natural language parser [22] into a list of the attribute phrases, and encode them as $\{\mathbf{t}_1, \dots, \mathbf{t}_{N_{att}}\}$, $\mathbf{t}_i \in \mathbb{R}^d$. Our core contribution is the cross-modal alignment learning that matches each visual component \mathbf{v}_a with its corresponding textual phrase \mathbf{t}_a , along with the global representation matching $\langle \mathbf{v}_0, \mathbf{t}_0 \rangle$ for the person search by natural language task.

3.1 The Image Stream

We adopt the sub-network of ResNet-50 (conv1, conv2_x, conv3_x, conv4_x) [17] as the backbone to extract feature maps \mathbf{F} from the input image. Then, we introduce a global branch \mathcal{F}^{glb} , and multiple local branches \mathcal{F}_a^{loc} to generate global visual features $\mathbf{v}_0 = \mathcal{F}^{glb}(\mathbf{F})$, and attribute visual features $\{\mathbf{v}_1 \dots \mathbf{v}_{N_{att}}\}$ respectively, where $\mathbf{v}_a = \mathcal{F}_a^{loc}(\mathbf{F})$. The network architectures are shown in Table 1. On the top of all the local branches is an auxiliary segmentation layer supervising each local branch to generate the segmentation map of one specific attribute category (shown in Figure 3). Intuitively, we argue that the additional auxiliary task acts as a knowledge regulator that diversifies each local branch to present attribute-specific features.

Our segmentation layer utilizes the architecture of a lightweight MaskHead [16] and can be removed during inference phase to reduce the computational cost. The remaining unsolved problem is that parsed annotations are not available in all person search datasets. To address that, we first train a human parsing network with HRNet [42] as an off-the-shelf tool, where the HRNet is jointly trained on multiple human parsing datasets: MHPv2 [57], ATR [28], and VIPeR [44].



Fig. 2. Attribute annotation generated by human parsing network. Torsos are labeled as background since there is no corresponding textual descriptions.

We then use the attribute category predictions as our segmentation annotations (illustrated in Figure 2). With these annotations, local branches receive the supervision needed from the segmentation task to learn attribute-specific features. Essentially, we are distilling the attribute information from a well-trained human parsing networks to the lightweight segmentation layer through joint training¹.

Discussion. Using attribute feature has the following advantages over the global features. 1) The textual annotations in person search by natural language task describe the person mostly by their dressing/body appearances, where the attribute features perfectly fit the situation. 2) Attribute aligning avoids the “*mal-positioned matching*” cases as shown in Fig. 1: using segmentation to regularize feature learning equips the model to be resilient over the diverse human poses or viewpoints, and also robust to the background noises.

3.2 The Language Stream

Given the raw textual description, our language stream first parses and extracts noun phrases *w.r.t.* each attribute through the Stanford POS tagger [31], and then feeds them into a language network to obtain the sentence-level as well as the phrase-level embeddings. We adopt a bi-directional LSTM to generate the global textual embedding \mathbf{t}_0 and the local textual embedding. Meanwhile, we adopt a dictionary clustering approach to categorize the novel noun phrases in the sentence to specific attribute phrases as in [8]. Concretely, we manually collect a list of words per attribute category, *e.g.*, “*shirt*”, “*jersey*”, “*polo*” to represent the upper-body category, and use the average-pooled word vectors [12] of them as the anchor embedding \mathbf{d}_a , and form the dictionary $\mathbf{D} = [\mathbf{d}_1, \dots, \mathbf{d}_{N_{att}}]$, where N_{att} is the total number of attributes. Building upon that, we assign the noun phrase to the category that has the highest cosine similarity, and form the local textual embedding $\{\mathbf{t}_1 \dots \mathbf{t}_N\}$. Different from previous works like [56, 32], we include accessory as one type of attribute as well, which serves as a crucial matching clue in many cases.

3.3 Visual-Textual Alignment Learning

Once we extract the global and attribute features, the key objective for the next stage is to learn a joint embedding space across the visual and the textual

¹ More details of our human parsing network and segmentation results can be found in the experimental part and the supplementary materials.

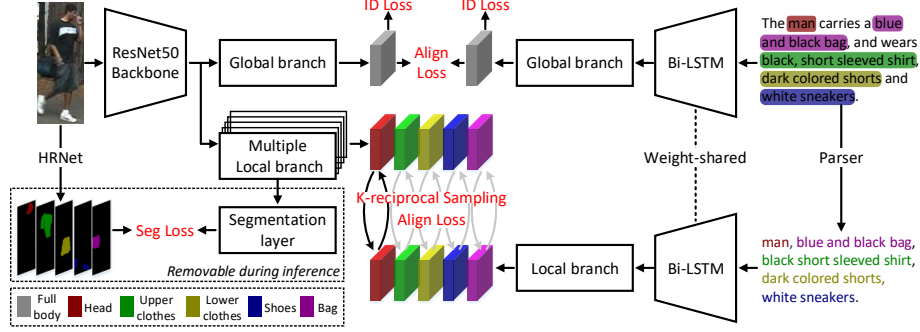


Fig. 3. Illustrative diagram of our ViTAA network, which includes an image stream (left) and a language stream (right). Our image stream first encodes the person image and extract both global and attribute representations. The local branch is additional supervised by an auxiliary segmentation layer where the annotations are acquired by an off-the-shell human parsing network. In the meanwhile, the textual description is parsed and decomposed into attribute atoms, and encoded by a weight-shared Bi-LSTM. We train our ViTAA jointly under global/attribute align loss in an end-to-end manner.

modalities, where the visual cues are tightly matched with the given textual description. Mathematically, we formulate our learning objective as a contrastive learning task that takes input as triplets, *i.e.*, $\langle \mathbf{v}^i, \mathbf{t}^+, \mathbf{t}^- \rangle$ and $\langle \mathbf{t}^i, \mathbf{v}^+, \mathbf{v}^- \rangle$, where i denotes the index of person to identify, and $+/-$ refer to the corresponding feature representations of the person i , and a randomly sampled irrelevant person respectively. We note that features in the triplet can be both at the global-level and the attribute-level. In the following, we discuss the learning schema on $\langle \mathbf{v}^i, \mathbf{t}^+, \mathbf{t}^- \rangle$ which can be extended to $\langle \mathbf{t}^i, \mathbf{v}^+, \mathbf{v}^- \rangle$.

We adopt the cosine similarity as the scoring function between visual and textual features $S = \frac{\mathbf{v}^T \cdot \mathbf{t}}{\|\mathbf{v}\| \cdot \|\mathbf{t}\|}$. For a positive pair $\langle \mathbf{v}^i, \mathbf{t}^+ \rangle$, the cosine similarity S^+ is encouraged to be as large as possible, which we define as *absolute similarity criterion*. While for a negative pair $\langle \mathbf{v}^i, \mathbf{t}^- \rangle$, enforcing the cosine similarity S^- to be minimal may yield an arbitrary constraint over the negative samples \mathbf{t}^- . Instead, we propose to optimize the deviation between S^- and S^+ to be larger than a preset margin, called *relative similarity criterion*. These criterion can be formulated as:

$$S^+ \rightarrow 1 \text{ and } S^+ - S^- > m, \quad (1)$$

where m is the least margin that positive and negative similarity should differ and is set to 0.2 in practice.

In contrastive learning, the general form of the basic objective function are either hinge loss $\mathcal{L}(\mathbf{x}) = \max\{0, 1 - \mathbf{x}\}$ or logistic loss $\mathcal{L}(\mathbf{x}) = \log(1 + \exp(-\mathbf{x}))$. One crucial drawback of hinge loss is that its derivative *w.r.t.* \mathbf{x} is a constant value: $\frac{\partial \mathcal{L}}{\partial \mathbf{x}} = -1$. Since the pair-based construction of training data leads to a polynomial growth of training pairs, inevitably we will have a certain part of the randomly sampled negative texts being less informative during training. Treat-

ing all the redundant samples equally might raise the risk of a slow convergence and/or even model degeneration for the metric learning tasks. While the derivative of logistic loss *w.r.t.* x is: $\frac{\partial \mathcal{L}}{\partial x} = -\frac{1}{e^x + 1}$, which is related with the input value. Hence, we settle with the logistic loss as our basic objective function.

With the logistic loss, the aforementioned criterion can be further derived and rewritten as:

$$(S^+ - \alpha) > 0, \quad -(S^- - \beta) > 0, \quad (2)$$

where $\alpha \rightarrow 1$ denotes the lower bound for positive similarity and $\beta = (\alpha - m)$ denotes the upper bound for negative similarity. Together with logistic loss function, our final *Alignment loss* can be unrolled as:

$$\mathcal{L}_{align} = \frac{1}{N} \sum_{i=1}^N \left\{ \log \left[1 + e^{-\tau_p(S_i^+ - \alpha)} \right] + \log \left[1 + e^{\tau_n(S_i^- - \beta)} \right] \right\}, \quad (3)$$

where τ_p and τ_n denote the temperature parameters that adjust the slope of gradient. The partial derivatives are calculated as:

$$\frac{\partial \mathcal{L}_{align}}{\partial S_i^+} = \frac{-\tau_p}{1 + e^{\tau_p(S_i^+ - \alpha)}}, \quad \frac{\partial \mathcal{L}_{align}}{\partial S_i^-} = \frac{\tau_n}{1 + e^{\tau_n(\beta - S_i^-)}}. \quad (4)$$

Thus, we show that Eq. 3 outputs continuous gradients and will assign higher weights to more informative samples accordingly.

K-reciprocal Sampling. One of the premise of visual-textual alignment is to fully exploit the informative positive and negative samples $\mathbf{t}^+, \mathbf{t}^-$ to provide valid supervisions. However, most of the current contrastive learning methods [48, 54] construct the positive pairs by selecting samples belonging to the same class and simply treat the random samples from other classes as negative. This is viable when using only global information at coarse level during training, but may not be able to handle the case as illustrated in Figure 1 where a fine-grained level comparison is needed. This practice is largely depending on the average number of samples for each attribute category to provide comprehensive positive samples. With this insight, we propose to further enlarge the searching space of positive samples from the cross-id incidents.

For instance, as in Figure 1, though the two ladies are with different identities, they share the extremely alike shoes which can be treated as the positive samples for learning. We term these kinds of samples with identical attributes but belong to different person identities as the “surrogate positive samples”. Kindly including the common attribute features of the surrogate positive samples in positive pairs makes much more sense than the reverse. It is worth noting that, this is unique only to our attribute alignment learning phase because attributes can only be compared at the fine-grained level. Now the key question is, how can we dig out the surrogate positive samples since we do not have direct cross-ID attribute annotations? Inspired by the re-ranking techniques in re-id community [11, 61], we propose k-reciprocal sampling as an unsupervised method to generate the surrogate labels at the attribute-level. How does the proposed method sample from a batch of visual and textual features? Straightforwardly,

for each attribute a , we can extract a batch of visual and textual features from the feature learning network and mine their corresponding surrogate positive samples using our sampling algorithm. Since we are only discussing the input form of $\langle \mathbf{v}^i, \mathbf{t}^+, \mathbf{t}^- \rangle$, our sampling algorithm is actually mining the surrogate positive textual features for each \mathbf{v}^i . Note that, if the attribute information in either modality is missing after parsing, we can simply ignore them during sampling.

Algorithm 1: *K-reciprocal Sampling Algorithm*

Input:

$\mathcal{V}_a = \{\mathbf{v}_a^i\}_{i=1}^N$ is a set of visual feature for attribute a
 $\mathcal{T}_a = \{\mathbf{t}_a^i\}_{i=1}^N$ is a set of textual feature for attribute a

Output:

\mathcal{P}_a is a set of surrogate positive sample for attribute a

```

1 for each  $\mathbf{v}_a \in \mathcal{V}_a$  do
2   find the top- $K$  nearest neighbours of  $\mathbf{v}_a$  w.r.t.  $\mathcal{T}_a$ :  $\mathcal{K}_{\mathcal{T}_a}$ ;
3    $\mathcal{S} \leftarrow \emptyset$ ;
4   for each  $\mathbf{t}_a \in \mathcal{K}_{\mathcal{T}_a}$  do
5     find the top- $K$  nearest neighbours of  $\mathbf{t}_a$  w.r.t.  $\mathcal{V}_a$ :  $\mathcal{K}_{\mathcal{V}_a}$ ;
6     if  $\mathbf{v}_a \in \mathcal{K}_{\mathcal{V}_a}$  then
7        $\mathcal{S} = \mathcal{S} \cup \mathbf{t}_a$ 
8     end
9   end
10   $\mathcal{P}_a = \mathcal{P}_a \cup \mathcal{S}$ 
11 end

```

3.4 Joint Training

The entire network is trained in an end-to-end manner. We adopt the widely-used cross-entropy loss (ID Loss) to assist the learning of the discriminative features of each instance, as well as pixel-level cross-entropy loss (Seg Loss) to classify the attribute categories in the auxiliary segmentation task. For the cross-modal alignment learning, we design the Alignment Loss on both the global-level and the attribute-level representations. The overall loss function thus emerges:

$$\mathcal{L} = \mathcal{L}_{id} + \mathcal{L}_{seg} + \mathcal{L}_{align}^{glo} + \mathcal{L}_{align}^{attr}. \quad (5)$$

4 Experiment

4.1 Experimental Setting.

Datasets. We conduct experiments on the CUHK-PEDES [25] dataset, which is currently the only benchmark for person search by natural language. It contains 40,206 images of 13,003 different persons, where each image comes with two human-annotated sentences. The dataset is split into 11,003 identities with 34,054 images in the training set, 1,000 identities with 3,078 images in validation, and 1,000 identities with 3,074 images in testing set.

Table 1. Detailed architecture of our global and local branches in image stream. *#Branch.* denotes the number of sub-branches.

Layer name	Parameters	Output size	#Branch
\mathcal{F}^{glb}	$\begin{bmatrix} 3 \times 3, 2048 \\ 3 \times 3, 2048 \end{bmatrix} \times 2$	24×8	1
	Average Pooling	1×1	
\mathcal{F}^{loc}	$\begin{bmatrix} 3 \times 3, 256 \\ 3 \times 3, 256 \end{bmatrix} \times 2$	24×8	5
	Max Pooling	1×1	

Evaluation protocols. Following the standard evaluation setting, we adopt Recall@K (K=1, 5, 10) as the retrieval criteria. Specifically, given a text description as query, Recall@K (R@K) reports the percentage of the images where at least one corresponding person is retrieved correctly among the top-K results.

Implementation details. For the global and local branches in image stream, we use the Basicblock as described in [17], where each branch is randomly initialized (detailed architecture is shown in Table 1). We use horizontally flipping as data augmenting and resize all the images to 384×128 . We use the Adam solver as the training optimizer with weight decay set as 4×10^{-5} , and involves 64 image-language pairs per mini-batch. The learning rate is initialized at 2×10^{-4} for the first 40 epochs during training, then decayed by a factor of 0.1 for the remaining 30 epochs. The whole experiment is implemented on a single Tesla V100 GPU machine. The hyperparameters in Eq. 3 are empirically set as: $\alpha = 0.6, \beta = 0.4, \tau_p = 10, \tau_n = 40$.

Pedestrian attributes parsing. Based on the analysis of image and natural language annotations in the dataset, we warp both visual and textual attributes into 5 categories: head (including descriptions related to hat, glasses, and face), clothes on the upper body, clothes on the lower body, shoes and bags (including backpack and handbag). We reckon that these attributes are visually distinguishable from both modalities. In Figure 2, we visualize the segmentation maps generated by our human parsing network, where attribute regions can be properly segmented and associated with correct labels.

4.2 Comparisons with the State-of-The-Arts

Result on CUHK-PEDES dataset. We summarize the performance of Vi-TAA and compare it with state-of-the-art methods in Table 2 on the CUHK-PEDES test set. Methods like GNA-RNN [25], CMCE [24], PWM-ATH [4] employ attention mechanism to learn the relation between visual and textual representation, while Dual Path [60], CMPM+CMPC [55] design objective function for better joint embedding learning. These methods only learn and utilize the “global” feature representation of both image and text. Moreover, MIA [32] ex-

Table 2. Person search results on the CUHK-PEDES test set. Best results are in bold.

Method	Feature	R@1	R@5	R@10
GNA-RNN [25]	global	19.05	-	53.64
CMCE [24]	global	25.94	-	60.48
PWM-ATH [4]	global	27.14	49.45	61.02
Dual Path [60]	global	44.40	66.26	75.07
CMPM+CMPC [55]	global	49.37	-	79.27
MIA [32]	global+region	53.10	75.00	82.90
GALM [19]	global+keypoint	54.12	75.45	82.97
ViTAA	global+attribute	55.97	75.84	83.52

exploits “*region*” information by dividing the input image into several horizontal stripes and extracting noun phrases from the natural language description. Similarly, GALM [19] leverage “*keypoint*” information from human pose estimation as an attention mechanism to assist feature learning and together with a noun phrases extractor implemented on input text. Though the above two utilize the local-level representations, neither of them learns the associations between visual features with textual phrases. From Table 2, we observe that ViTAA shows a consistent lead on all metrics (R@1-10), outperforming the GALM [19] by a margin of 1.85%, 0.39%, 0.55% and claims the new state-of-the-art results. We note that though the performance could be considered as incremental, the shown improvement on the R@1 performance is challenging. It suggests that the alignment learning of ViTAA contributes to the retrieval task directly. We further report the ablation studies on the effect of different components, and exhibit the attribute retrieval results quantitatively and qualitatively.

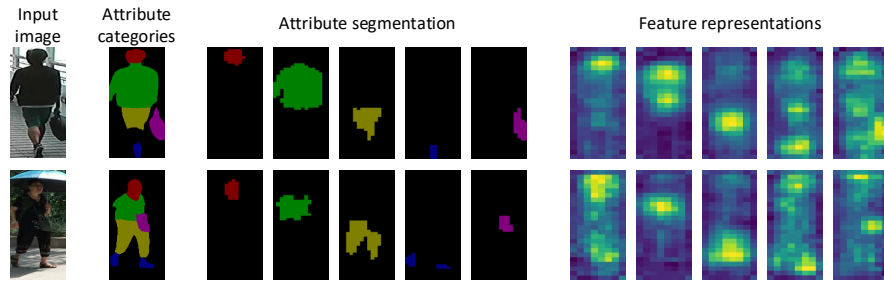
4.3 Ablation Study

We carry out comprehensive ablations to evaluate the contribution of different components and the training configurations.

Comparisons over different component combinations. To compare the individual contribution of each component, we set the baseline model as the one trained with only ID loss. In Table 3, we report the improvement of the proposed components (segmentation, global-alignment, and attribute-alignment) on the basis of the baseline model. From the table, we have the following observations and analyses: First, using segmentation loss only brings marginal improvement because the visual features are not aligned with their corresponding textual features. Similarly, we observe the same trend when the training is combined with only attribute-alignment loss where the visual features are not properly segmented, thus can not be associated for retrieval. An incremental gain is obtained by combining these two components. Next, compared with attribute-level, global-level alignment greatly improves the performance under all criteria, which demonstrates the efficiency of the visual-textual alignment schema. The cause of the performance gap is that: the former is learning the attribute similarity across

Table 3. The improvement of components added on baseline model. Glb-Align and Attr-Align represent global-level and attribute-level alignment respectively.

Model Component			R@1	R@5	R@10	
Segmentation	Attr-Align	Glb-Align				
✓	✓		29.68	51.84	61.57	
			30.93	52.71	63.11	
			31.40	54.09	63.66	
✓	✓	✓	39.26	61.22	68.14	
✓	✓		52.27	73.33	81.61	
			55.97	75.84	83.52	

**Fig. 4.** From left to right, we exhibit the raw input person images, attribute labels generated by the pre-trained HRNet, attribute segmentation result from our segmentation layer, and their corresponded feature maps from the local branches.

different person identities while the latter is concentrating on the uniqueness of each person. At the end, by combining all the loss terms yields the best performance, validating that our global-alignment and attribute-alignment learning are complimentary with each other.

Visual attribute segmentation and representations. In Figure 6, we visualize the segmentation maps from the segmentation layer and the feature representations of the local branches. It evidently shows that, even transferred using only a lightweight structure, the auxiliary person segmentation layer produces accurate pixel-wise labels under different human pose. This suggests that person parsing knowledge has been successfully distilled our local branches, which is crucial for the precise cross-modal alignment learning. On the right side of Figure 6, we showcase the feature maps of local branch per attribute.

K-reciprocal sampling. We investigate how the value of K impacts the pair-based sampling and learning process. We evaluate the R@1 and R@10 performance under different K settings in Figure 5(a). Ideally, the larger the K is, the more potential surrogate positive samples will be mined, while this also comes with the possibility that more non-relevant examples (false positive examples) might be incorrectly sampled. Result in Figure 5(a) agrees with our analysis: best R@1 and R@10 is achieved when K is set to 8, and the performances are

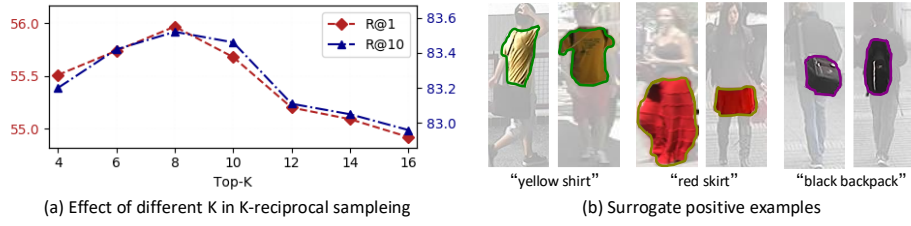


Fig. 5. (a) R@1 and R@10 results across different K value in the proposed surrogate positive data sampling method. (b) Some examples of the surrogate positive data with different person identities.



Fig. 6. Examples of person search results on CUHK-PEDES. We indicate the true/false matching results in green/red boxes.

persistently declining as K goes larger. In Figure 5(b), we provide visual examinations of the surrogate positive pairs that mined by our sampling method. The visual attributes from different persons serve as valuable positive samples in our alignment learning schema.

Qualitative analysis. We present the qualitative examples of person retrieval results to provide a more in-depth examination. As shown in Figure 6, we illustrate the top-10 matching results using the given query. In the successful case (top), ViTAA precisely capture all attributes in the target person. It is worth noting that the wrong answers still capture the relevant attributes: “sweater with black, gray and white stripes”, “tan pants”, and “carrying a bag”. For the failure case (bottom), though the retrieved images are incorrect, we observe that all the attributes described in the query are there in almost all retrieved results.

Table 4. Upper-body clothing attribute retrieve results. Attr is the short form of attribute and “upblack” denotes the upper-body in black.

Market1501				DukeMTMC			
Attr.	#Target	R@1	mAP	Attr.	#Target	R@1	mAP
upblack	1759	99.2	44.0	upblack	11047	100.0	82.2
upwhite	3491	44.7	64.8	upwhite	1095	91.3	35.4
upred	1354	92.3	54.8	upred	821	100.0	44.6
uppurple	363	100.0	61.2	uppurple	65	0.0	9.0
upyellow	1195	72.7	75.6	upgray	2012	81.7	29.8
upgray	1755	49.7	55.2	upblue	1577	77.9	31.3
upblue	1100	70.1	32.4	upgreen	417	89.2	24.5
upgreen	949	87.9	50.4	upbrown	345	18.3	15.2

4.4 Extension: Attribute Retrieval

In order to validate the ability of associating the visual attribute with the text phrase, we further conduct attribute retrieval experiment on the datasets of Market-1501 [59] and DukeMTMC [34], where 27 and 23 human related attributes are annotated per image by [29]. In our experiment, we use our pre-trained ViTAA on CUHK-PEDES without any further finetuning, and conduct the retrieval task using the attribute phrase as the query under R@1 and mAP metrics. In our experiment, we simply test on the *upper-body clothing* attribute category, and post the retrieval results in Table 4. We introduce the details of our experiment in the supplementary materials. From Table 4, it clearly shows that ViTAA achieves great performances on almost all sub-attributes. This further strongly supports our argument that ViTAA is able to associate the visual attribute features with textual attribute descriptions successfully.

5 Conclusion

In this work, we present a novel ViTAA model to address the person search by natural language task from the perspective of an attribute-specific alignment learning. In contrast to the existing methods, ViTAA fully exploits the common attribute information in both visual and textual modalities across different person identities, and further builds strong association between the visual attribute features and their corresponding textual phrases by using our alignment learning schema. We show that ViTAA achieves state-of-the-art results on the challenging benchmark CUHK-PEDES and demonstrate its promising potential that further advances the person search by natural language domain.

Acknowledgements. Vising scholarship support for Z. Wang from the China Scholarship Council #201806020020 and Amazon AWS Machine Learning Research Award (MLRA) support are greatly appreciated. Any opinions, findings, and conclusion or recommendations expressed in this material are those of the authors and do not necessarily reflect the view of the sponsors.

References

1. Antol, S., Agrawal, A., Lu, J., Mitchell, M., Batra, D., Lawrence Zitnick, C., Parikh, D.: Vqa: Visual question answering. In: Proceedings of the IEEE international conference on computer vision. pp. 2425–2433 (2015)
2. Benenson, R., Omran, M., Hosang, J., Schiele, B.: Ten years of pedestrian detection, what have we learned? In: European Conference on Computer Vision. pp. 613–627. Springer (2014)
3. Chen, D., Li, H., Liu, X., Shen, Y., Shao, J., Yuan, Z., Wang, X.: Improving deep visual representation for person re-identification by global and local image-language association. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 54–70 (2018)
4. Chen, T., Xu, C., Luo, J.: Improving text-based person search by spatial matching and adaptive threshold. In: 2018 IEEE Winter Conference on Applications of Computer Vision (WACV). pp. 1879–1887 (March 2018)
5. Dollár, P., Wojek, C., Schiele, B., Perona, P.: Pedestrian detection: A benchmark. In: Computer Vision and Pattern Recognition, 2009. CVPR 2009. IEEE Conference on. pp. 304–311. IEEE (2009)
6. Dong, Q., Gong, S., Zhu, X.: Person search by text attribute query as zero-shot learning. In: The IEEE International Conference on Computer Vision (ICCV) (October 2019)
7. Fang, Z., Gokhale, T., Banerjee, P., Baral, C., Yang, Y.: Video2commonsense: Generating commonsense descriptions to enrich video captioning. arXiv preprint arXiv:2003.05162 (2020)
8. Fang, Z., Kong, S., Fowlkes, C., Yang, Y.: Modularized textual grounding for counterfactual resilience. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
9. Fang, Z., Kong, S., Yu, T., Yang, Y.: Weakly supervised attention learning for textual phrases grounding. arXiv preprint arXiv:1805.00545 (2018)
10. Frome, A., Corrado, G.S., Shlens, J., Bengio, S., Dean, J., Ranzato, M., Mikolov, T.: Devise: A deep visual-semantic embedding model. In: Advances in neural information processing systems. pp. 2121–2129 (2013)
11. Garcia, J., Martinel, N., Micheloni, C., Gardel, A.: Person re-identification ranking optimisation by discriminant context information analysis. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1305–1313 (2015)
12. Goldberg, Y., Levy, O.: word2vec explained: deriving mikolov et al.’s negative-sampling word-embedding method. arXiv preprint arXiv:1402.3722 (2014)
13. Gong, S., Cristani, M., Yan, S., Loy, C.C.: Person Re-Identification. Springer Publishing Company, Incorporated (2014)
14. Guo, J., Yuan, Y., Huang, L., Zhang, C., Yao, J.G., Han, K.: Beyond human parts: Dual part-aligned representations for person re-identification. In: The IEEE International Conference on Computer Vision (ICCV) (October 2019)
15. Han, C., Ye, J., Zhong, Y., Tan, X., Zhang, C., Gao, C., Sang, N.: Re-id driven localization refinement for person search. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 9814–9823 (2019)
16. He, K., Gkioxari, G., Dollár, P., Girshick, R.: Mask r-cnn. In: Proceedings of the IEEE international conference on computer vision. pp. 2961–2969 (2017)
17. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 770–778 (2016)

18. Jeon, J., Lavrenko, V., Manmatha, R.: Automatic image annotation and retrieval using cross-media relevance models. In: Proceedings of the 26th annual international ACM SIGIR conference on Research and development in information retrieval. pp. 119–126 (2003)
19. Jing, Y., Si, C., Wang, J., Wang, W., Wang, L., Tan, T.: Pose-guided joint global and attentive local matching network for text-based person search. arXiv preprint arXiv:1809.08440 (2018)
20. Kalayeh, M.M., Basaran, E., Gökmen, M., Kamasak, M.E., Shah, M.: Human semantic parsing for person re-identification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1062–1071 (2018)
21. Karpathy, A., Fei-Fei, L.: Deep visual-semantic alignments for generating image descriptions. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 3128–3137 (2015)
22. Klein, D., Manning, C.D.: Fast exact inference with a factored model for natural language parsing. In: Advances in neural information processing systems. pp. 3–10 (2003)
23. Layne, R., Hospedales, T.M., Gong, S.: Attributes-based re-identification. In: Person Re-Identification, pp. 93–117. Springer (2014)
24. Li, S., Xiao, T., Li, H., Yang, W., Wang, X.: Identity-aware textual-visual matching with latent co-attention. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 1890–1899 (2017)
25. Li, S., Xiao, T., Li, H., Zhou, B., Yue, D., Wang, X.: Person search with natural language description. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1970–1979 (2017)
26. Li, W., Zhao, R., Xiao, T., Wang, X.: Deepreid: Deep filter pairing neural network for person re-identification. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 152–159 (2014)
27. Liang, X., Gong, K., Shen, X., Lin, L.: Look into person: Joint body parsing & pose estimation network and a new benchmark. *IEEE transactions on pattern analysis and machine intelligence* **41**(4), 871–885 (2018)
28. Liang, X., Liu, S., Shen, X., Yang, J., Liu, L., Dong, J., Lin, L., Yan, S.: Deep human parsing with active template regression. *Pattern Analysis and Machine Intelligence, IEEE Transactions on* (12), 2402–2414 (Dec 2015)
29. Lin, Y., Zheng, L., Zheng, Z., Wu, Y., Hu, Z., Yan, C., Yang, Y.: Improving person re-identification by attribute and identity learning. *Pattern Recognition* (2019)
30. Liu, X., Zhao, H., Tian, M., Sheng, L., Shao, J., Yi, S., Yan, J., Wang, X.: Hydraplus-net: Attentive deep features for pedestrian analysis. In: Proceedings of the IEEE international conference on computer vision. pp. 350–359 (2017)
31. Manning, C.D., Surdeanu, M., Bauer, J., Finkel, J., Bethard, S.J., McClosky, D.: The Stanford CoreNLP natural language processing toolkit. In: Association for Computational Linguistics (ACL) System Demonstrations. pp. 55–60 (2014)
32. Niu, K., Huang, Y., Ouyang, W., Wang, L.: Improving description-based person re-identification by multi-granularity image-text alignments. arXiv preprint arXiv:1906.09610 (2019)
33. Plummer, B.A., Wang, L., Cervantes, C.M., Caicedo, J.C., Hockenmaier, J., Lazebnik, S.: Flickr30k entities: Collecting region-to-phrase correspondences for richer image-to-sentence models. In: Proceedings of the IEEE international conference on computer vision. pp. 2641–2649 (2015)
34. Ristani, E., Solera, F., Zou, R., Cucchiara, R., Tomasi, C.: Performance measures and a data set for multi-target, multi-camera tracking. In: European Conference on Computer Vision. pp. 17–35. Springer (2016)

35. Rohrbach, A., Rohrbach, M., Hu, R., Darrell, T., Schiele, B.: Grounding of textual phrases in images by reconstruction. In: European Conference on Computer Vision. pp. 817–834. Springer (2016)
36. Shekhar, R., Jawahar, C.: Word image retrieval using bag of visual words. In: 2012 10th IAPR International Workshop on Document Analysis Systems. pp. 297–301. IEEE (2012)
37. Si, J., Zhang, H., Li, C.G., Kuen, J., Kong, X., Kot, A.C., Wang, G.: Dual attention matching network for context-aware feature sequence based person re-identification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5363–5372 (2018)
38. Su, C., Li, J., Zhang, S., Xing, J., Gao, W., Tian, Q.: Pose-driven deep convolutional model for person re-identification. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 3960–3969 (2017)
39. Su, C., Zhang, S., Xing, J., Gao, W., Tian, Q.: Multi-type attributes driven multi-camera person re-identification. *Pattern Recognition* **75**, 77–89 (2018)
40. Sudowe, P., Spitzer, H., Leibe, B.: Person attribute recognition with a jointly-trained holistic cnn model. In: Proceedings of the IEEE International Conference on Computer Vision Workshops. pp. 87–95 (2015)
41. Suh, Y., Wang, J., Tang, S., Mei, T., Mu Lee, K.: Part-aligned bilinear representations for person re-identification. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 402–419 (2018)
42. Sun, K., Xiao, B., Liu, D., Wang, J.: Deep high-resolution representation learning for human pose estimation. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
43. Sun, Y., Zheng, L., Yang, Y., Tian, Q., Wang, S.: Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 480–496 (2018)
44. Tan, Z., Yang, Y., Wan, J., Hang, H., Guo, G., Li, S.Z.: Attention-based pedestrian attribute analysis. *IEEE transactions on image processing* (12), 6126–6140 (2019)
45. Wang, C., Zhang, Q., Huang, C., Liu, W., Wang, X.: Mancs: A multi-task attentional network with curriculum sampling for person re-identification. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 365–381 (2018)
46. Wang, G., Yuan, Y., Chen, X., Li, J., Zhou, X.: Learning discriminative features with multiple granularities for person re-identification. In: 2018 ACM Multimedia Conference on Multimedia Conference. pp. 274–282. ACM (2018)
47. Wang, Z., Wang, J., Yang, Y.: Resisting crowd occlusion and hard negatives for pedestrian detection in the wild. *arXiv preprint arXiv:2005.07344* (2020)
48. Wen, Y., Zhang, K., Li, Z., Qiao, Y.: A discriminative feature learning approach for deep face recognition. In: European conference on computer vision. pp. 499–515. Springer (2016)
49. Wu, H., Mao, J., Zhang, Y., Jiang, Y., Li, L., Sun, W., Ma, W.Y.: Unified visual-semantic embeddings: Bridging vision and language with structured meaning representations. In: The IEEE Conference on Computer Vision and Pattern Recognition (CVPR) (June 2019)
50. Xu, J., Zhao, R., Zhu, F., Wang, H., Ouyang, W.: Attention-aware compositional network for person re-identification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 2119–2128 (2018)

51. Xu, K., Ba, J., Kiros, R., Cho, K., Courville, A., Salakhudinov, R., Zemel, R., Bengio, Y.: Show, attend and tell: Neural image caption generation with visual attention. In: International conference on machine learning. pp. 2048–2057 (2015)
52. Yin, Z., Zheng, W.S., Wu, A., Yu, H.X., Wan, H., Guo, X., Huang, F., Lai, J.: Adversarial attribute-image person re-identification. In: Proceedings of the Twenty-Seventh International Joint Conference on Artificial Intelligence, IJCAI-18. pp. 1100–1106. International Joint Conferences on Artificial Intelligence Organization (7 2018)
53. You, Q., Zhang, Z., Luo, J.: End-to-end convolutional semantic embeddings. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 5735–5744 (2018)
54. Zhang, X., Fang, Z., Wen, Y., Li, Z., Qiao, Y.: Range loss for deep face recognition with long-tailed training data. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 5409–5418 (2017)
55. Zhang, Y., Lu, H.: Deep cross-modal projection learning for image-text matching. In: Proceedings of the European Conference on Computer Vision (ECCV). pp. 686–701 (2018)
56. Zhang, Z., Lan, C., Zeng, W., Chen, Z.: Densely semantically aligned person re-identification. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 667–676 (2019)
57. Zhao, J., Li, J., Cheng, Y., Sim, T., Yan, S., Feng, J.: Understanding humans in crowded scenes: Deep nested adversarial learning and a new benchmark for multi-human parsing. In: 2018 ACM Multimedia Conference on Multimedia Conference. pp. 792–800. ACM (2018)
58. Zheng, L., Huang, Y., Lu, H., Yang, Y.: Pose invariant embedding for deep person re-identification. *IEEE Transactions on Image Processing* (2019)
59. Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., Tian, Q.: Scalable person re-identification: A benchmark. In: Proceedings of the IEEE international conference on computer vision. pp. 1116–1124 (2015)
60. Zheng, Z., Zheng, L., Garrett, M., Yang, Y., Shen, Y.D.: Dual-path convolutional image-text embedding with instance loss. *arXiv preprint arXiv:1711.05535* (2017)
61. Zhong, Z., Zheng, L., Cao, D., Li, S.: Re-ranking person re-identification with k-reciprocal encoding. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1318–1327 (2017)