

AXM-Net: Cross-Modal Context Sharing Attention Network for Person Re-ID

Ammarah Farooq

Muhammad Awais

Josef Kittler

Syed Safwan Khalid
CVSSP, University of Surrey

{ammarah.farooq, m.a.rana, j.kittler, s.khalid}@surrey.ac.uk

Abstract

Cross-modal person re-identification (Re-ID) is critical for modern video surveillance systems. The key challenge is to align inter-modality representations according to semantic information present for a person and ignore background information. In this work, we present AXM-Net, a novel CNN based architecture designed for learning semantically aligned visual and textual representations. The underlying building block consists of multiple streams of feature maps coming from visual and textual modalities and a novel learnable context sharing semantic alignment network. We also propose complementary intra modal attention learning mechanisms to focus on more fine-grained local details in the features along with a cross-modal affinity loss for robust feature matching. Our design is unique in its ability to implicitly learn feature alignments from data. The entire AXM-Net can be trained in an end-to-end manner. We report results on both person search and cross-modal Re-ID tasks. Extensive experimentation validates the proposed framework and demonstrates its superiority by outperforming the current state-of-the-art methods by a significant margin.

1. Introduction

Person re-identification (Re-ID) is a critical component of intelligent video surveillance systems which aims to retrieve a queried person from a large database of pedestrian images. The database typically contains non-overlapping camera viewpoints with respect to query image. Depending on the type of information provided as a query, the task is referred as person Re-ID or person search in research literature. Person search [22] aims to find a person based on the natural language description of the person while images are provided as query in person Re-ID. In person search, there is no constraint on the camera viewpoints of the person, i.e., in the visual gallery person may have the same pose for which the description is given. Nevertheless, in

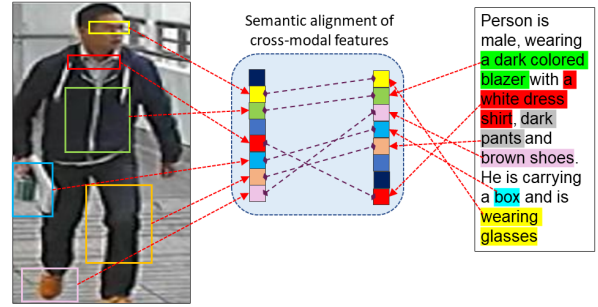


Figure 1. Illustration of semantic alignment for visual and textual features. The semantic information present in the features should be aligned to learn the associations among image parts and textual phrases for robust retrieval.

both tasks it is critical to learn discriminative feature representations which are unique to an individual, as well as well aligned with-in class for good matching.

Recent research literature is packed with numerous deep learning based person re-identification approaches. These methods aim to learn robust person representations [45, 10], apply various attention mechanisms [38, 4], look for cross-domain knowledge transfer [17, 9] and so on. Cross-modal person Re-ID is another important aspect of Re-ID task [40, 11, 12, 25, 17]. The dependency on available image query limits the practical application of a vision based system. For example, in case of criminal search, quite often the CCTV (or normal image) of criminal is not available. Therefore, officers rely on the unique cues of the criminal from the witnesses descriptions often given in term of natural language description. In such cases, with no images available, this descriptive information serves as a query for person Re-ID. Hence, employing a multi-modal Re-ID system can overcome the limitations of image based systems.

This work focuses on cross-modal person Re-ID using visual and textual information of the person while keeping the genuine Re-ID constraints of non overlapping camera viewpoints. There are number of challenges while dealing

with two distinct modalities. First, the structure of information in both modalities is quite different with respect to persons. Presumably, images have persons always standing up-right while the person description can have any order. Second, it is critical to learn a network which is able to extract the semantics in data instead of memorising corresponding image-text pairs for the identities seen during training. Third, the semantic information present in the features, for example, colour and type of clothes person is wearing, activity of the person, accessories etc, should be aligned across the modalities to learn the associations among image parts and textual phrases and disregard background noise (Figure 1).

The main idea of this work is to align visual and textual features of a person to enable cross-modal or multi-modal search seamlessly. In order to achieve this, we present AXM-Net, a novel convolutional neural network (CNN) based architecture designed to deal with the challenges mentioned above and capable of learning semantically aligned cross-modal feature representations. The underlying building block consists of multiple streams of feature maps capturing variable amount of local context from visual and textual networks and a novel context sharing semantic alignment network that is learnt based on the critical cues present in both modalities. The output feature maps are, hence, attended according to the fused information. The context sharing semantic alignment network leverages multi-scale, multi-context intra modality information and inter modality semantic information to boost the informative channels and suppress the channels containing noisy/background information.

Apart from exploiting inter modal semantic and contextual information, we also propose modality specific attention mechanisms to effectively extract discriminative and complementary intra modality representations. To be specific, we introduce an auxiliary part-based feature learning branch to vision in order to locally attend different spatial parts for finer details. We also note that contextual information from other parts can provide useful information while attending a given spatial part, therefore, we propose to use part based semi-global context sharing. We use CNN based language modelling, however, the sequential nature of the language modality necessitates learning long-term associations among the person attributes. An annotator can describe person in any sequence of attention to body parts and accessories. For example, a description may contain information about head (hairs, style) at the beginning follows by description for lower body and again carries information about head (wearing hat). Hence, we propose to employ a non-local attention technique to learn these useful associations among textual attributes. Our contributions can be summarised as follows:

- Adaptive cross-modal context sharing semantic alignment

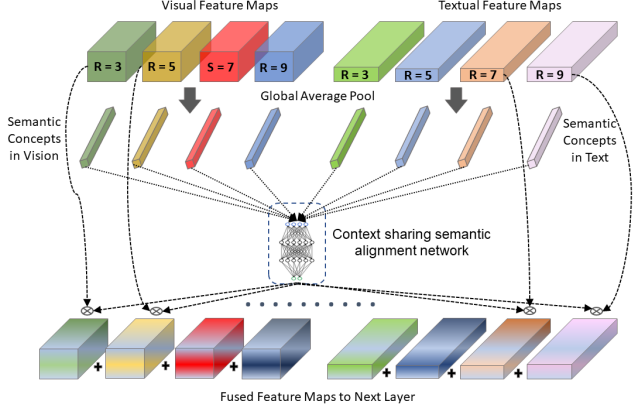


Figure 2. Adaptive cross-modal context sharing semantic alignment (AXM)-Block

block (AXM-Block) is proposed to capture context aware shared semantic concepts and suppress noisy information between the vision and text modality. To our knowledge, this is the first work to employ implicit semantic alignment across modalities in the person Re-ID setting.

- We put forward effective intra-modality attention mechanisms to extract local spatial region based details in vision representations while exploiting semi-global context sharing and to learn inter-dependencies among textual phrases.
- We design a cross-modal affinity loss which rewards and penalises the learning network based on the similarity among the aligned features.
- Extensive experiments demonstrate the superiority of the proposed AXM-Net over the state-of-the-art (SOTA) person search models on the CUHK-PEDES [22] benchmark and cross-modal Re-ID on CrossRe-ID and CUHK-SYSU [11, 12]. We also propose a cross-modal protocol for the famous Market-1501 data [43].

2. Related Work

2.1. Person Re-identification

Recent vision based Re-ID approaches are mainly based on deep learning techniques [29, 19, 32, 10, 45, 13, 31]. These approaches can be divided into the following categories. [39, 38, 4, 23, 30, 2] devise attention based mechanisms on the basis of higher order information present in the features, person’s body mask or pixel/region level features. [6] advocates a diverse attention mechanism along with orthogonality constraint to preserve diversity between the layer activations and weights. Pose based methods [16, 26] use pose detectors to learn aligned features across poses. [24, 14] use GANs for learning pose-invariant features. Another group of works [31, 37, 18] focus on local part based feature learning. The work [45] proposes an omni-scale network to learn fused multi-scale features. Our approach

also uses multi-scale features, however, the design aims to attend the most useful features across modalities based on the context present in these features.

2.2. Person Search

The task of person search by natural language descriptions was introduced by [22]. The proposed method was a CNN-RNN network to learn global level cross-modal features. The following works [21, 7] also focused on similar network architectures and a little improvement was observed in performance. Major improvements were shown by [5, 42, 44] where researchers start exploiting global-local associations and improving the feature embedding space. More recent works [16, 36, 35, 1] started employing auxiliary learning branches to explicitly make use of pose key-points, person attributes, segmentation masks, body parts and textual phrases. These approaches brought improvements as compared to using only global features. Our method also takes advantage of local level features from both modalities in a unique and effective way. Moreover, it is the first design to learn visual and textual features using an integrated convolutional block.

2.3. Cross-modal Person Re-identification

As mentioned earlier, cross-modal searches help to overcome limitations of vision only systems. A number of approaches have been proposed based on using infra-red images [33, 25, 41] along with RGB images. For text based person Re-ID, [40] published the pioneering work and reported results on the CUHK03 and Viper datasets under multiple retrieval scenarios. Recent works [11, 12] proposed to jointly optimise the two modalities and applied canonical correlation analysis to enhance similarity between the corresponding features. These works also reported results for larger cross-modal test splits including CrossRe-ID and CUHK-SYSU.

3. Cross-modal AXM-Net Framework

Figure 3 shows the overall network architecture of the proposed AXM-Net, which includes a feature learning backbone network, a global vision branch, a spatial part based vision branch and a globally attended textual branch. The details of each part are presented in the following subsections.

3.1. Vision and Textual Feature Learning Network

For the feature learning network, we take motivation from [8, 45, 20, 28] to exploit multi-scale features for person Re-ID. We propose a novel idea to align representations across modalities based on the semantic information. For this purpose, we introduce an adaptive context sharing semantic alignment block (AXM-block) between vision and

text. The illustrative diagram of this block is shown in Figure 2. Each block consists of multiple feature streams bringing features from various receptive fields (R) for both modalities, capturing different amount of local context, and a context sharing semantic alignment network. The input features of the AXM-block are rich in intra-modal semantics, having coarse-to-fine information from different contextual support regions R . The context sharing semantic alignment network aims to align this information for robust cross-modal retrieval. The architecture of feature learning backbone in the AXM-Net is presented in Table 1.

Adaptive Cross-modal Context Sharing Semantic Alignment Network: Intrinsically, the semantic information present in both modalities is the same as both are describing the same person identity. We propose to combine this information dynamically by employing a learnable cross-modal context sharing semantic alignment neural network. The network amplifies the semantic information by jointly paying attention to salient semantic information from multiple modalities and at the same time reduces the effect of non-informative parts which are not agreed upon across the modalities (c.f. Figure 2). This property of AXM-Net is particularly useful as the textual input does not have background clutter information unlike visual input. Therefore, the effect of background can be reduced by the AXM-Net as shown in the Subsection 4.5. To model the shared semantics across modalities, we first define $X = (V_r, T_r)$ to be the visual and textual features maps at receptive fields $r = 1, 2, \dots, R$. At first, we apply a global average pooling operation to collect the channel-wise global semantics present in each feature map generating $\mathbf{x} = (\mathbf{v}_r, \mathbf{t}_r)$ vectors of length spanning the entire channel dimension. These vectors are then fed to the context and semantic learning mini-network together to get scale vectors to re-scale the information based on the cross-modal semantics.

$$\tilde{\mathbf{x}} = C(\mathbf{x}_r) \quad \forall \mathbf{v}_r, \mathbf{t}_r \text{ in } \mathbf{x} \quad (1)$$

The function $C(x)$ indicates the non-linear mapping to generate a set of scale vectors $\tilde{\mathbf{x}} = (\tilde{\mathbf{v}}_r, \tilde{\mathbf{t}}_r)$ corresponding to each input feature map. These scaling vectors carry the importance of each channel with respect to the shared context between modalities. Finally, feature alignment is performed by taking element-wise product between scaling vectors and the corresponding input feature map.

$$(\tilde{V}, \tilde{T}) = \sum_{r=1}^R (V_r \odot \tilde{\mathbf{v}}_r, T_r \odot \tilde{\mathbf{t}}_r) \quad (2)$$

3.2. Attentive Local Feature Learning for Vision

From visual perspective, the person Re-ID task depends on the discriminative local cues characterising each individual. The global visual branch focuses onto the person as

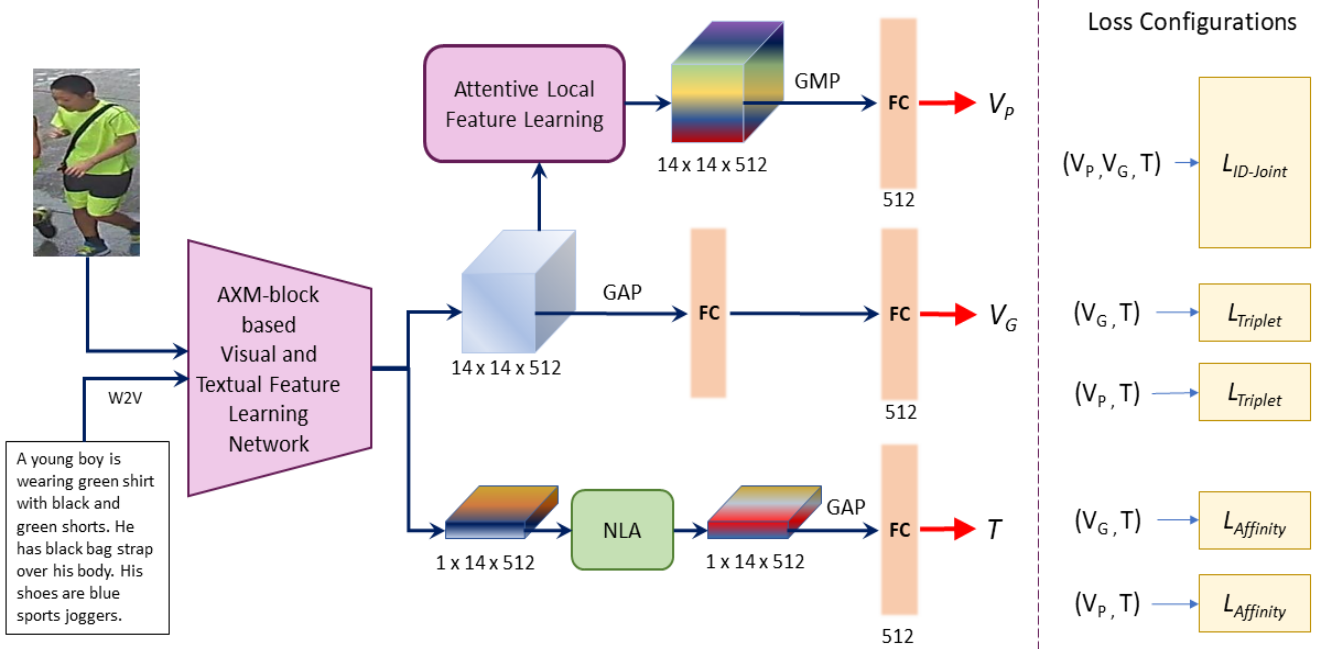


Figure 3. Illustrative diagram of our cross-modal AXM-Net, which generates global visual feature V_G , local spatial part based visual feature V_P and textual feature T . Softmax loss $\mathcal{L}_{ID_{joint}}$ is shared for all the features. Matching losses are trained pairwise with the textual feature for each visual feature.

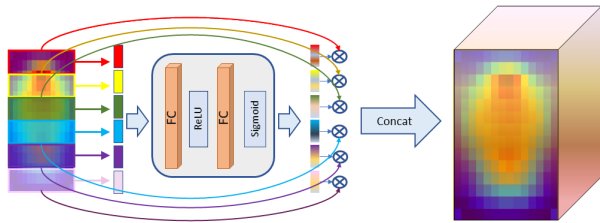


Figure 4. Attentive local feature learning for visual features

a whole, while the AXM-block preserves semantics across modalities. However, it is important to pay attention to local cues within the modality. Therefore, we propose a part based attention network to enhance information in each local image part shown in Figure 4. We divide the feature maps into multiple horizontal strips [31]. Each strip is locally attended channel-wise by a learnable multi-layer network. The network looks into each local channel strip and preserves the most informative channels for a given spatial region. The attended region strips are concatenated back into the input-sized features. We notice that context from other spatial region plays a critical role when attending a particular spatial region. In order to utilise the context across the spatial regions the parameters are shared across the local attention network. Eventually, the learnt visual representations are spatially and channel-wise attended with respect to local regions, where the effect of background

in each strip has already been minimised by the semantic alignment network in AXM-block.

3.3. Non-local Attention for Textual Features

The nature of the free-form natural language descriptions makes it challenging to find relations among the attributes. It is important to capture long-term relationships among the verbal attributes present in the textual representations. Inspired by [34, 15], we propose non-local attention (NLA) to model these dependencies by directly computing interactions between a pair of textual features. Figure 5 shows the schematic for the NLA block. We take feature maps from the last convolution block of the feature learning network as input to NLA. Intuitively, these features represent the important semantic attributes present in the person’s description. Each spatial index represents a response from a local spatial region (phrases). The NLA block takes the feature from each spatial position and computes its affinity with all other features. An attention map is generated based on these affinities and applied to the input feature maps.

3.4. Objective Function

The loss function for training is the sum of softmax loss of the three feature branches, triplet-loss [3] and the proposed affinity loss. Specifically, the softmax weights are shared among all branches to enforce intra-identity feature alignment across the different branches and modalities and is denoted as $\mathcal{L}_{ID_{joint}}$. The retrieval losses are optimised

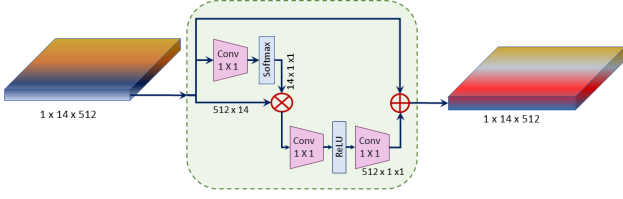


Figure 5. Non-local attention for textual features

Stage	Module	Vision Output	Textual Output
input	-	$224 \times 224, 3$	$1 \times 56, 300$
conv1	7×7 conv, stride=2, 3×3 max pool (1×3 conv, stride=1) on text	$56 \times 56, 64$	$1 \times 56, 64$
conv2	AXM block $\times 2$	$28 \times 28, 256$	$1 \times 28, 256$
conv3	AXM block $\times 2$	$14 \times 14, 384$	$1 \times 14, 384$
conv4	AXM block $\times 2$	$14 \times 14, 512$	$1 \times 14, 512$
conv5	1×1 conv	$14 \times 14, 512$	$1 \times 14, 512$
GAP	global average pool	$1 \times 1, 512$	$1 \times 1, 512$

Table 1. Architecture of visual and textual feature learning network.

in a pairwise manner for each visual feature with textual feature.

$$\mathcal{L}_{Total} = \mathcal{L}_{ID_{joint}} + \mathcal{L}_{trip}(V_g, T) + \mathcal{L}_{trip}(V_p, T) + \mathcal{L}_{aff}(V_g, T) + \mathcal{L}_{aff}(V_p, T) \quad (3)$$

Given tuples (V_a, T_+, T_-) , (T_a, V_+, V_-) , \mathcal{L}_{trip} is a margin (α) based ranking loss, defined over similarity(S) of cross-modal positive and negative feature pairs as follows:

$$\mathcal{L}_{trip} = \max[0, \alpha - (S(V_a, T_+) - (S(V_a, T_-)) + \max[0, \alpha - (S(V_+, T_a) - (S(V_-, T_a)) \quad (4)$$

Cross-modal Affinity Loss: In order to enhance retrieval performance, we also propose a simple yet effective affinity loss. It is based on the affinity between image-text feature pairs. The image features and the corresponding textual features should be aligned and have as high similarity as possible, ideally +1 (in case of cosine similarity) and non corresponding image and textual features should unaligned (uncorrelated) and have as low similarity as possible, ideally 0. Therefore, given the representations for an image-text pair, affinity is measured in terms of cosine similarity score between image feature V and text feature T . The affinity loss is implemented as a binary cross entropy criterion, defined over $\{V_i, T_j, y_{ij}\}$ where $y_{ij} = 1$ for matching pairs and $y_{ij} = 0$ for non-matching ones.

$$\mathcal{L}_{aff} = -[y_{ij} \cdot \log(\sigma(S(V_i, T_j))) + (1 - y_{ij}) \cdot \log(\sigma(1 - S(V_i, T_j)))] \quad (5)$$

where σ is sigmoid function applied to features similarity.

4. Experiments and Results

4.1. Implementation Details

We follow a two stage training strategy to learn the network. In the first stage, we focus on learning the textual branch, vision part branch and fully connected layers from scratch, while keeping the vision backbone fixed to pretrained ImageNet weights. For the first stage the training follows the standard classification paradigm considering each person as an individual class and only using $\mathcal{L}_{ID_{joint}}$. We also apply label smoothing to our cross entropy loss. We use batch size 64, weight decay $5e-4$ and initial learning rate 0.01 with stochastic gradient descent optimisation. Images are resized to 224×224 . Each textual description is mapped to its 300 dimensional word2vec [27] embedding and resized as $1 \times 56 \times 300$ where 56 is the maximum sentence length. We adopted random flipping, random erasing for images and random circular shift of sentences as data augmentation. To achieve computational efficiency, we employ depth-wise separable convolutions at each layer. The retrieval performance is measured based on the cosine similarity between the features and reported in terms of Rank@1 and mean average precision (mAP).

4.2. Datasets

CUHK-PEDES The CUHK person description data [22] is the only large-scale benchmark available for cross-modal person search. It has 13003 person IDs with 40,206 images and 80,440 descriptions. There are 11003,1000,1000 pre-defined IDs for training, validation and test sets. The training and test set include 34054/68126 and 3074/6156 images/descriptions respectively.

CrossRe-ID Dataset For cross-modal Re-ID, we evaluate the models on the protocol introduced by [12] on the test split of CUHK-PEDES data. The gallery and query splits have been carefully separated across viewpoints. The descriptions are also varying across viewpoints. The dataset includes 824 unique IDs. There are 1511/3022 and 1096/2200 images/descriptions in gallery and query sets respectively.

CUHK-SYSU We evaluate our model on the test protocol provided by [11]. There are 5532 IDs for training and 2099 IDs for testing. The corresponding descriptions have been extracted from CUHK-PEDES data. The final gallery and query splits contain 5070/10140 and 3271/6550 images/descriptions respectively.

Market-1501 We also propose a test split for the Market Re-ID dataset. It is a part of the CUHK-PEDES data. We separate the train and test IDs according to the original protocol. Then we merge the train IDs with the rest of the CUHK-PEDES data to form the training set and evaluate on the test set. The test set is again split into non-overlapping gallery and query sets. Hence, the training set

Model	Feature Type	Rank@1	Rank@5	Rank@10	mAP
GNA-RNN [22]	global	19.05	-	53.64	-
IATV [21]		25.94	-	60.48	-
PWM [7]		27.14	49.45	61.02	-
DPCE [44]		44.40	66.26	75.07	-
GLA [5]		43.58	66.93	76.26	-
CMPC + CPM [42]		49.37	-	79.27	-
Baseline: Multi-scale features + joint ID		52.78	72.33	80.29	49.04
PMA [16]	global + keypoints	53.81	73.54	81.23	-
IMG-Net [36]	global + parts	56.48	76.89	85.01	-
ViTAA [35]	global + attribute	55.97	75.84	83.52	-
CMAAM [1]		56.68	77.18	84.86	-
AXM-Net + joint ID	global+ part	59.11	77.46	83.80	54.24
AXM-Net + joint ID + affinity		59.81	77.43	84.27	54.89
AXM-Net + joint ID + triplet		60.12	77.84	84.95	55.35
AXM-Net + joint ID + triplet + affinity		61.19	78.24	84.97	56.13

Table 2. Comparison with SOTA models on the CUHK-PEDES dataset

includes 11253 IDs with 34191/68396 images/descriptions. The test set has 750 IDs. The gallery and query sets have 1816/3632 and 1173/2346 images/descriptions respectively. The results are evaluated under single-shot and single-query scenarios for this split.

4.3. Comparison with State-of-the-Art Methods

4.3.1 Results on Person Search

We summarise the performance of the proposed AXM-Net with the state-of-the-art methods on person search in Table 2. The methods have been grouped according to the type of feature representations used for learning. We implement the baseline network with multi-scale features for both modalities and a joint classifier layer. The baseline method performs best among all global level techniques which signifies the benefit of having multi-scale features but still it falls behind the other complex methods based on various feature types. Although [36] uses horizontal image feature parts and intra-modal attention, this attention is performed on global feature maps followed by extracting multiple part features. However, our attentive local feature learning block attends local information of the person and its design implicitly induces both channel and spatial attentions. We incorporate context from all spatial parts to attend a given spatial part. We use single fused part level feature for ID classification as compared to the previous methods, which is computationally efficient in terms of parameters required by the linear layers in each auxiliary branch. Other recent techniques [1, 35], employ multi-label classification loss over the extracted attribute features for both modalities. The features are obtained by an attribute prediction network and pixel-wise segmentation network respectively. It is worth

mentioning here that our method is simple but powerful and enhances the semantic alignment between modalities without any explicit supervision from segmented body parts, pose key-points or attributes. The proposed AXM-Net with simple $\mathcal{L}_{ID_{joint}}$ loss outperforms current SOTA by a large margin, achieving over 59% Rank@1 performance. By using affinity and triplet loss together, we set the new SOTA Rank@1 of 61.19%. Note that parameter free affinity loss enhances the retrieval and is competitive to the triplet loss.

4.3.2 Results on Cross-modal Re-ID

We follow the evaluation protocol of [12] for cross-modal re-identification. In Tables 3, 4, $V \rightarrow V$ indicates image based search, $T \rightarrow V$ indicates description to image search and $VT \rightarrow V$ indicates using both modalities for query and vision as gallery. We report the detailed results including Rank@5 and Rank@10 for all the datasets in the supplementary materials. For CrossRe-ID and CUHK-SYSU, we compare the results with the recently reported joint training technique followed by applying canonical correlation analysis to embed cross-modal features [12]. It is mentioned as JT+CCA in Tables 3. For both datasets, the proposed AXM-Net outperformed the previous method by a significant margin in all retrieval scenarios. Note that, now the $T \rightarrow V$ indicates a description of a person from a different pose, and the gallery images have different poses. We can witness the potential of semantic alignment across-modalities in this challenging case. The improvement in Rank@1 performance shows the capability of the AXM-Net in matching viewpoint-invariant semantic details.

Retrieval results on the proposed Market-1501 test split are reported in Table 4. We trained the OSNet [45] on the

Model	CrossRe-ID						CUHK-SYSU					
	V → V		T → V		VT → V		V → V		T → V		VT → V	
	Rank@1	mAP	Rank@1	mAP	Rank@1	mAP	Rank@1	mAP	Rank@1	mAP	Rank@1	mAP
JT + CCA [12]	86.77	88.90	33.61	39.40	88.59	87.95	74.13	77.16	11.37	15.78	77.68	75.8
AXM-Net + joint ID + affinity	95.14	96.04	44.66	50.49	95.26	95.22	86.00	87.75	19.93	24.82	88.72	87.02
AXM-Net + joint ID + triplet	95.02	96.00	47.33	52.58	95.75	95.41	86.24	88.02	20.93	26.04	87.86	86.40
AXM-Net + joint ID + affinity + triplet	94.29	98.9	46.48	52.21	94.05	93.93	85.86	87.70	21.44	26.77	88.62	86.73

Table 3. Performance comparison on cross-modal Re-ID. Query → Gallery

Model	V → V		T → V		VT → V	
	Rank@1	mAP	Rank@1	mAP	Rank@1	mAP
OSNet [45]	48.66	54.39	-	-	-	-
AXM-Net + joint ID + affinity	85.73	88.01	41.86	48.09	90.53	86.56
AXM-Net + joint ID + triplet	85.33	87.72	43.73	49.44	89.86	86.53
AXM-Net + joint ID + affinity + triplet	84.80	87.35	43.20	49.24	90.40	86.64

Table 4. Performance comparison on the proposed cross-modal split for Market-1501 . Query → Gallery

Model	Context Share	Part Attention	NLA	Rank@1
Baseline	-	-	-	52.78
Model: 1	✓	-	-	55.90
Model: 2	-	✓	-	56.99
Model: 3	-	-	✓	53.25
Model: 4	✓	-	✓	56.27
Model: 5	✓	✓	-	58.59
Unified Visual	✓	✓	✓	54.53
Single Stage	✓	✓	✓	55.05
Proposed AXM-Net	✓	✓	✓	59.11

Table 5. Ablation study on the AXM-Net on CUHK-PEDES test set

Attention Weights	Rank@1	Pooling Type	Rank@1
Separate	57.62	Average (GAP)	57.36
Shared	57.93	Max (GMP)	58.82
Feature Drop	Rank@1	No. of FC Layers	Rank@1
Random	57.62	1	58.45
Part	58.33	2	58.06
No drop	59.11	Proposed	59.11

Table 6. Design parameters for the attentive visual part learning branch

proposed split and compared result for vision based scenario. First, we see large performance gap between single and multi-modality learning. Second, the OSNet is trained with much larger data as compared to training set of only Market data. But it is clear from the result that the matching capability of the single modality is limited by number of samples per class seen during training. It gives an important insight that having textual descriptions not only help as query for matching but also decrease the dependency on number of samples of training set. Experimental results for the AXM-Net proves that the proposed system is able to focus on important cues of the person instead of holistically memorising the image-description pairs.

4.4. Component Analysis

In order to assess the contribution of each component in the complete AXM-Net, we perform ablation study on

the proposed framework by adding the components step-by-step. The study is performed on the CUHK-PEDES test set with joint softmax loss $\mathcal{L}_{ID_{joint}}$ and all hyper-parameters are kept the same for training in all settings. Table 5 presents the corresponding results. The **baseline** model has the same architecture as the global visual and textual branch, including multi-scale features for both modalities but no context sharing. We list our observations as follows:

- **Effect of Individual Components. Models: 1, 2, 3** correspond to individual component’s contribution. We note that each component has boosted the retrieval capability compared to the baseline. Each component is essential to the design of cross-modal Re-ID. The context sharing semantic alignment induced by AXM-block enhances the useful semantic information across modalities. The auxiliary visual branch brings the locally informative cues, while non-local attention keeps track of textual dependencies. **Models: 4, 5** also emphasise the complementary effect of various components together.
- **Effect of Unified Visual Feature Branch.** We experiment with unified visual branch by removing the global branch and keeping only the attentive feature branch. It is indicated as **Unified Visual** in Table 5. We observe a performance drop of 4.58% as compared to the proposed design. Separate learning branches for global and part level features complement each other in our design.
- **Effect of Single Stage versus Two Stages of Learning.** Working with cross-modal networks is also challenging in terms of the learning policy being adopted. As mentioned earlier, we used two stage training for network learning. We also test a single stage policy in which we tune all parameters together with the same initialisation setting. **Single stage** model clearly shows the difference between the two policies. We notice that aligning the textual branch with vision in the first stage and then training them jointly is the best strategy.
- **Design Parameters for Vision Part Branch.** We con-

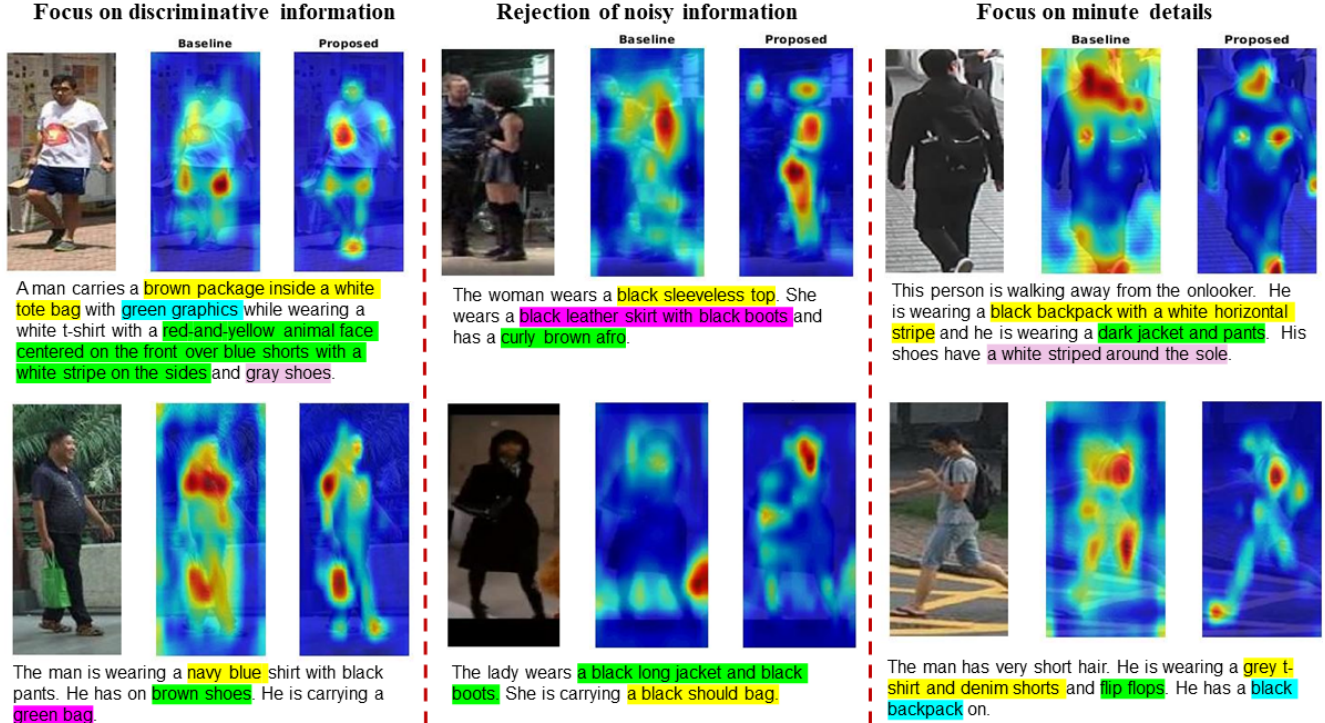


Figure 6. Attention maps visualisation for AXM-Net. The proposed network focuses on discriminative information and rejects background.

sider several parameters for designing the local visual feature learning branch as presented in Table 6. First of all, we test with separate attention networks for each strip of feature maps. We find that having a shared attention network helps in learning as well as reducing the number of parameters. It supports our idea of semi-global context sharing and signifies the connectedness of various semantic concepts across the strips, for example, a long coat covers both the upper and lower parts of the body. Next, we note that the global max pool helps in capturing local cues by focusing on the highest responses in each region. For each branch of AXM-Net, we also optimise the number of FC layers as it is critical to avoid any performance degradation and network overfitting. We observe from the table that having identical linear layer structure not always implies the best performing solution. Being a richer modality and focusing on information from the whole image, two FC layers help in the global branch to learn better representations. We also consider feature dropping technique [10] to obtain robust features. However, we observed a decrease in Rank@1 using both batchwise random location and horizontal strip(part) drops.

4.5. Visualisation

We analyse the attention maps for the proposed AXM-Net to deduce its capability to learn person specific semantic information from the input data. We show the gener-

ated visual feature maps for the baseline network and our AXM-Net in Figure 6. The visualisations are arranged in three columns, highlighting different aspects of a retrieval system. In all three columns, we observe that AXM-Net ignores person’s background with high confidence. Specifically, in the middle column (upper example), we note that the attention is focused on the lady whose description is provided, while the baseline model generates a spanned attention. In the third column, we observe that with the help of the textual description, the visual attention is refined to very minute details in the image. In all examples, we observe that the semantic information present in the textual description is highly emphasised across vision, which verifies the feature alignment induced by the inter-modal semantic alignment. Note that the baseline network also has access to multi-scale information, but the proposed dynamic cross-modal attention has intelligently fused this information.

5. Conclusion

In this work, we present a novel AXM-Net model to address the problem of cross-modal person re-identification and search. Our innovation involves a contextual alignment of features coming from the visual and textual modalities and local intra-modal attentions. In contrast to existing methods, the proposed AXM-Net is the first framework based on an integrated convolutional feature learning block, the AXM-block, and implicit semantic alignment of

the features across modalities. The contextual attention is learnt by a shared learnable network inside the AXM-block. The experimental results show that our network defines new SOTA performance on the CUHK-PEDES benchmark and also demonstrate the potential of the proposed network for more challenging cross-modal person Re-ID applications.

References

- [1] Surbhi Aggarwal, Venkatesh Babu RADHAKRISHNAN, and Anirban Chakraborty. Text-based person search via attribute-aided matching. In *The IEEE Winter Conference on Applications of Computer Vision*, pages 2617–2625, 2020. 3, 6
- [2] Honglong Cai, Zhiguan Wang, and Jinxing Cheng. Multi-scale body-part mask guided attention for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, pages 0–0, 2019. 2
- [3] Gal Chechik, Varun Sharma, Uri Shalit, and Samy Bengio. Large scale online learning of image similarity through ranking. In *Iberian Conference on Pattern Recognition and Image Analysis*, pages 11–14. Springer, 2009. 4
- [4] Binghui Chen, Weihong Deng, and Jiani Hu. Mixed high-order attention network for person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 1, 2
- [5] Dapeng Chen, Hongsheng Li, Xihui Liu, Yantao Shen, Jing Shao, Zejian Yuan, and Xiaogang Wang. 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)*, pages 54–70, 2018. 3, 6
- [6] Tianlong Chen, Shaojin Ding, Jingyi Xie, Ye Yuan, Wuyang Chen, Yang Yang, Zhou Ren, and Zhangyang Wang. Abdnnet: Attentive but diverse person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 2
- [7] Tianlang Chen, Chenliang Xu, and Jiebo Luo. Improving text-based person search by spatial matching and adaptive threshold. In *2018 IEEE Winter Conference on Applications of Computer Vision (WACV)*, pages 1879–1887. IEEE, 2018. 3, 6
- [8] Yanbei Chen, Xiatian Zhu, and Shaogang Gong. Person re-identification by deep learning multi-scale representations. In *Proceedings of the IEEE International Conference on Computer Vision Workshops*, pages 2590–2600, 2017. 3
- [9] Yanbei Chen, Xiatian Zhu, and Shaogang Gong. Instance-guided context rendering for cross-domain person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 1
- [10] Zuozhuo Dai, Mingqiang Chen, Xiaodong Gu, Siyu Zhu, and Ping Tan. Batch dropblock network for person re-identification and beyond. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3691–3701, 2019. 1, 2, 8
- [11] Ammarah Farooq, Muhammad Awais, Fei Yan, Josef Kittler, Ali Akbari, and Syed Safwan Khalid. A convolutional baseline for person re-identification using vision and language descriptions. *arXiv preprint arXiv:2003.00808*, 2020. 1, 2, 3, 5
- [12] Ammarah Farooq, Muhammad Awais, Fei Yan, Josef Kittler, Ali Akbari, and Syed Safwan Khalid. Cross modal person re-identification with visual-textual queries. In *2020 IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 2020. 1, 2, 3, 5, 6, 7
- [13] Yang Fu, Yunchao Wei, Yuqian Zhou, Honghui Shi, Gao Huang, Xinchao Wang, Zhiqiang Yao, and Thomas Huang. Horizontal pyramid matching for person re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8295–8302, 2019. 2
- [14] Yixiao Ge, Zhuowan Li, Haiyu Zhao, Guojun Yin, Shuai Yi, Xiaogang Wang, et al. Fd-gan: Pose-guided feature distilling gan for robust person re-identification. In *Advances in neural information processing systems*, pages 1222–1233, 2018. 2
- [15] Zilong Huang, Xinggang Wang, Lichao Huang, Chang Huang, Yunchao Wei, and Wenyu Liu. Cnet: Criss-cross attention for semantic segmentation. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 4
- [16] Ya Jing, Chenyang Si, Junbo Wang, Wei Wang, Liang Wang, and Tieniu Tan. Pose-guided multi-granularity attention network for text-based person search. *arXiv preprint arXiv:1809.08440*, 2018. 2, 3, 6
- [17] Ya Jing, Wei Wang, Liang Wang, and Tieniu Tan. Cross-modal cross-domain moment alignment network for person search. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 1
- [18] Dangwei Li, Xiaotang Chen, Zhang Zhang, and Kaiqi Huang. Learning deep context-aware features over body and latent parts for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 384–393, 2017. 2
- [19] Jianing Li, Jingdong Wang, Qi Tian, Wen Gao, and Shiliang Zhang. Global-local temporal representations for video person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3958–3967, 2019. 2
- [20] Jianing Li, Shiliang Zhang, and Tiejun Huang. Multi-scale 3d convolution network for video based person re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8618–8625, 2019. 3
- [21] Shuang Li, Tong Xiao, Hongsheng Li, Wei Yang, and Xiaogang Wang. Identity-aware textual-visual matching with latent co-attention. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 1890–1899, 2017. 3, 6
- [22] Shuang Li, Tong Xiao, Hongsheng Li, Bolei Zhou, Dayu Yue, and Xiaogang Wang. Person search with natural language description. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1970–1979, 2017. 1, 2, 3, 5, 6

- [23] Wei Li, Xiatian Zhu, and Shaogang Gong. Harmonious attention network for person re-identification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2285–2294, 2018. 2
- [24] Jinxian Liu, Bingbing Ni, Yichao Yan, Peng Zhou, Shuo Cheng, and Jianguo Hu. Pose transferrable person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 4099–4108, 2018. 2
- [25] Yan Lu, Yue Wu, Bin Liu, Tianzhu Zhang, Baopu Li, Qi Chu, and Nenghai Yu. Cross-modality person re-identification with shared-specific feature transfer. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, June 2020. 1, 3
- [26] Jiaxu Miao, Yu Wu, Ping Liu, Yuhang Ding, and Yi Yang. Pose-guided feature alignment for occluded person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 2
- [27] Tomas Mikolov, Ilya Sutskever, Kai Chen, Greg S Corrado, and Jeff Dean. Distributed representations of words and phrases and their compositionality. In *Advances in neural information processing systems*, pages 3111–3119, 2013. 5
- [28] Xuelin Qian, Yanwei Fu, Yu-Gang Jiang, Tao Xiang, and Xiangyang Xue. Multi-scale deep learning architectures for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 5399–5408, 2017. 3
- [29] Ruijie Quan, Xuanyi Dong, Yu Wu, Linchao Zhu, and Yi Yang. Auto-reid: Searching for a part-aware convnet for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3750–3759, 2019. 2
- [30] Chunfeng Song, Yan Huang, Wanli Ouyang, and Liang Wang. Mask-guided contrastive attention model for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 1179–1188, 2018. 2
- [31] Yifan Sun, Liang Zheng, Yi Yang, Qi Tian, and Shengjin Wang. Beyond part models: Person retrieval with refined part pooling (and a strong convolutional baseline). In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 480–496, 2018. 2, 4
- [32] Guangcong Wang, Jianhuang Lai, Peigen Huang, and Xiaohua Xie. Spatial-temporal person re-identification. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 33, pages 8933–8940, 2019. 2
- [33] Guan'an Wang, Tianzhu Zhang, Jian Cheng, Si Liu, Yang Yang, and Zengguang Hou. Rgb-infrared cross-modality person re-identification via joint pixel and feature alignment. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 3
- [34] Xiaolong Wang, Ross Girshick, Abhinav Gupta, and Kaiming He. Non-local neural networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 7794–7803, 2018. 4
- [35] Zhe Wang, Zhiyuan Fang, Jun Wang, and Yezhou Yang. Vi-taa: Visual-textual attributes alignment in person search by natural language. *arXiv preprint arXiv:2005.07327*, 2020. 3, 6
- [36] Zijie Wang, Aichun Zhu, Zhe Zheng, Jing Jin, Zhouxin Xue, and Gang Hua. Img-net: inner-cross-modal attentional multigranular network for description-based person re-identification. *Journal of Electronic Imaging*, 29(4):043028, 2020. 3, 6
- [37] Longhui Wei, Shiliang Zhang, Hantao Yao, Wen Gao, and Qi Tian. Glad: Global-local-alignment descriptor for scalable person re-identification. *IEEE Transactions on Multimedia*, 21(4):986–999, 2018. 2
- [38] Bryan (Ning) Xia, Yuan Gong, Yizhe Zhang, and Christian Poellabauer. Second-order non-local attention networks for person re-identification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision (ICCV)*, October 2019. 1, 2
- [39] Jing Xu, Rui Zhao, Feng Zhu, Huaming Wang, and Wanli Ouyang. Attention-aware compositional network for person re-identification. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, pages 2119–2128, 2018. 2
- [40] Fei Yan, Josef Kittler, and Krystian Mikolajczyk. Person re-identification with vision and language. In *2018 24th International Conference on Pattern Recognition (ICPR)*, pages 2136–2141. IEEE, 2018. 1, 3
- [41] M. Ye, X. Lan, Q. Leng, and J. Shen. Cross-modality person re-identification via modality-aware collaborative ensemble learning. *IEEE Transactions on Image Processing*, 29:9387–9399, 2020. 3
- [42] Ying Zhang and Huchuan Lu. Deep cross-modal projection learning for image-text matching. In *Proceedings of the European Conference on Computer Vision (ECCV)*, pages 686–701, 2018. 3, 6
- [43] Liang Zheng, Liyue Shen, Lu Tian, Shengjin Wang, Jingdong Wang, and Qi Tian. Scalable person re-identification: A benchmark. In *Proceedings of the IEEE international conference on computer vision*, pages 1116–1124, 2015. 2
- [44] Zhedong Zheng, Liang Zheng, Michael Garrett, Yi Yang, Mingliang Xu, and Yi-Dong Shen. Dual-path convolutional image-text embeddings with instance loss. *ACM Transactions on Multimedia Computing, Communications, and Applications (TOMM)*, 16(2):1–23, 2020. 3, 6
- [45] Kaiyang Zhou, Yongxin Yang, Andrea Cavallaro, and Tao Xiang. Omni-scale feature learning for person re-identification. In *Proceedings of the IEEE International Conference on Computer Vision*, pages 3702–3712, 2019. 1, 2, 3, 6, 7