# Symbiotic Adversarial Learning for Attribute-based Person Search

Yu-Tong Cao*, Jingya Wang*, and Dacheng Tao

UBTECH Sydney AI Centre, School of Computer Science, Faculty of Engineering,
The University of Sydney, Darlington, NSW 2008, Australia
{ycao5602@uni.,jingya.wang@,dacheng.tao@}sydney.edu.au

**Abstract.** Attribute-based person search is in significant demand for applications where no detected query images are available, such as identifying a criminal from witness. However, the task itself is quite challenging because there is a huge modality gap between images and physical descriptions of attributes. Often, there may also be a large number of unseen categories (attribute combinations). The current state-of-the-art methods either focus on learning better cross-modal embeddings by mining only seen data, or they explicitly use generative adversarial networks (GANs) to synthesize unseen features. The former tends to produce poor embeddings due to insufficient data, while the latter does not preserve intra-class compactness during generation. In this paper, we present a symbiotic adversarial learning framework, called SAL. Two GANs sit at the base of the framework in a symbiotic learning scheme: one synthesizes features of unseen classes/categories, while the other optimizes the embedding and performs the cross-modal alignment on the common embedding space. Specifically, two different types of generative adversarial networks learn collaboratively throughout the training process and the interactions between the two mutually benefit each other. Extensive evaluations show SALs superiority over nine state-of-the-art methods with two challenging pedestrian benchmarks, PETA and Market-1501. The code is publicly available at: https://github.com/ycao5602/SAL.

**Keywords:** Person search, Cross-modal retrieval, Adversarial learning

## 1 Introduction

The goal with person search is to find the same person in non-overlapping camera views at different locations. In surveillance analysis, it is a crucial tool for public safety. To date, the most common approach to person search has been to take one detected image captured from a surveillance camera and use it as a query [25, 45, 54, 14, 58, 40, 26, 3]. However, this is not realistic in many real-world applications  for example, where a human witness has identified the criminal, but no image is available.

---

*  Equal contribution.

Attributes, such as gender, age, clothing or accessories, are more natural to us as searchable descriptions, and these can be used as soft biometric traits to search for in surveillance data [28, 22, 16, 33]. Compared to the queries used in image-based person search, these attributes are also much easier to obtain. Further, semantic descriptions are more robust than low-level visual representations, where changes in viewpoint or diverse camera conditions can be problematic. Recently, a few studies have explored sentence-based person search [26, 24]. Although this approach provides rich descriptions, unstructured text tends to introduce noise and redundancy during modeling. Attributes descriptions, on the other hand, are much cheaper to collect, and they inherently have a robust and relatively independent ability to discriminate between persons. As such, attribute descriptions have the advantage of efficiency over sentence descriptions in person search tasks.

Unfortunately, applying a cross-modal, attribute-based person search to real-world surveillance images is very challenging for several reasons. (1) There is a huge semantic gap between visual and textual modalities (i.e., attribute descriptions). Attribute descriptions are of lower dimensionality than visual data, e.g., tens versus thousands, and they are very sparse. Hence, in terms of data capacity, they are largely inferior to visual data. From performance comparisons between single-modal retrieval (image-based person search) and cross-modal retrieval (attribute-based person search) using current state-of-the-art methods, there is still a significant gap between the two, e.g., mAP 84.0% [34] vs. 24.3% [7] on Market-1501 [56, 27]. (2) Most of the training and testing classes (attribute combinations) are non-overlapping, which leads to zero-shot retrieval problems a very challenging scenario to deal with [7]. Given an attribute-style query, model aims to have more capacity to search an unseen class given a query of attributes. (3) Compared to the general zero-shot learning settings one might see in classification tasks, surveillance data typically has large intra-class variations and inter-class similarities with only a small number of available samples per class (Fig. 4), e.g., $\sim$ 7 samples per class in PETA [5] and $\sim$ 26 samples per class in Market-1501 compared with $\sim$ 609 samples per class in AWA2 [51]. One category (i.e., general attribute combinations that are not linked to a specific person) may include huge variations in appearance  just consider how many dark-haired, brown-eyed people you know and how different each looks. Also, the inter-class distance between visual representations from different categories can be quite small considering fine-grained attributes typically become ambiguous in low resolution with motion blur.

One general idea for addressing this problem is to use cross-modal matching to discover a common embedding space for the two modalities. Most existing methods focus on representation learning, where the goal is to find projections of data from different modalities and combine them in a common feature space. This way, the similarity between the two can be calculated directly [12, 48, 10] (Fig. 1(a)). However, these approaches typically have a weaker modality-invariant learning ability and, therefore, weaker performance in cross-modal retrieval. More recently, some progress has been made with the develop-
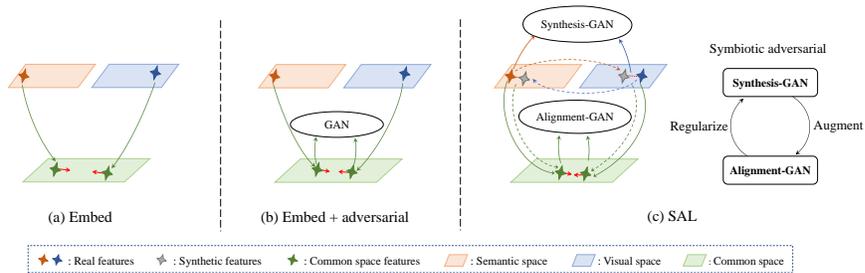
**Fig. 1.** Three categories of cross-modal retrieval with common space learning: (a) The embedding model that projects data from different modalities into common feature space. (b) The embedding model with common space adversarial alignment. (c) The proposed SAL that jointly considers feature-level data augmentation and cross-modal alignment by two GANs.

ment of GANs [11]. GANs can better align the distributions of representations across modalities in a common space [44], plus they have been successfully applied to attribute-based person search [53]. (Fig. 1(b)). There are still some bridges to cross, however. With only a few samples, only applying cross-modal alignment to a high-level common space may mean the model fails to capture variances in the feature space. Plus, the cross-modal alignment probably will not work for unseen classes. Surveillance data has all these characteristics, so all of these problems must be overcome.

To deal with unseen classes, some recent studies on zero-shot learning have explicitly used GAN-based models to synthesize those classes [30, 19, 63, 9]. Compared to zero-shot classification problems, our task is cross-modal retrieval, which requires learning a more complex search space with a finer granularity. As our experiments taught us, direct generation without any common space condition may reduce the intra-class compactness.

Hence, in this work, we present a fully-integrated symbiotic adversarial learning framework for cross-modal person search, called SAL. Inspired by symbiotic evolution, where two different organisms cooperate so both can specialize, we jointly explore the feature space synthesis and common space alignment with separate GANs in an integrated framework at different scales (Fig. 1(c)). The proposed SAL mainly consists of two GANs with interaction: (1) A synthesis-GAN that generates synthetic features from semantic representations to visual representations, and vice versa. The features are conditioned on common embedding information so as to preserve very fine levels of granularity. (2) An alignment-GAN that optimizes the embeddings and performs cross-modal alignment on the common embedding space. In this way, one GAN augments the data with synthetic middle-level features, while the other uses those features to optimize the embedding framework. Meanwhile, a new regularization term, called common space granularity-consistency loss, forces the cross-modal generated representations to be consistent with their original representations in the high-level
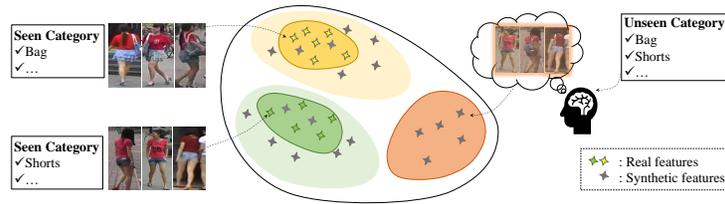
**Fig. 2.** An illustration of feature space data augmentation with a synthesized unseen class.

common space. These more reliable embeddings further boost the quality of the synthetic representations. To address zero-shot learning problems when there is no visual training data for an unseen category, we use new categories of combined attributes to synthesize visual representations for augmentation. An illustration of feature space data augmentation with a synthesized unseen class is shown in Fig. 2.

In summary, this paper makes the following main **contributions** to the literature:

- We propose a fully-integrated framework called symbiotic adversarial learning (SAL) for attribute-based person search. The framework jointly considers feature-level data augmentation and cross-modal alignment by two GANs at different scales. Plus, it handles cross-model matching, few-shot learning, and zero-shot learning problems in one unified approach.
- We introduce a symbiotic learning scheme where two different types of GANs learn collaboratively and specially throughout the training process. By generating qualified data to optimise the embeddings, SAL can better align the distributions in the common embedding space, which, in turn, yields superior cross-modal search results.
- Extensive evaluations on the PETA [5] and Market-1501 [56] datasets demonstrate the superiority of SAL at attribute-based person search over nine state-of-the-art attribute-based models.

## 2   Related Works

**Attribute-Based Person Search.** In recent years, semantic attribute recognition of pedestrians has drawn increasing attention [23, 21, 6, 46, 55] and it has been extensively exploited for image-based person search as a middle-level feature [20, 38, 39, 27, 47]. Some studies on exploit attribute-based person search for cross-modal retrieval [43, 37, 35, 53]. Early attribute-based retrieval methods intuitively rely on attribute prediction. For example, Vaquero et al. [43] proposed a human parts based attribute recognition method for person retrieval. Siddiquie et al. [37] utilized a structured prediction framework to integrate ranking and

retrieval with complex queries, while Scheirer et al. [35] constructed normalized multi-attribute spaces from raw classifier outputs. However, the pedestrian attribute recognition problem is far from being solved [23, 46, 29]. Consistently predicting all the attributes is a difficult task with sparsely-labeled training data: the images of people from surveillance cameras are low resolution and often contain motion blur. Also the same pedestrian would rarely be captured in the same pose. These imperfect attributes are significantly reducing the reliability of existing models. Shi et al. [36] suggested transfer learning as a way to overcome the limited availability of costly labeled surveillance data. They chose fashion images with richly labeled captions as the source domain to bridge the gap between unlabeled pedestrian surveillance images. They were able to produce semantic representations for a person search without annotated surveillance data, but the retrieval results were not as good as the supervised/semi-supervised methods, which raises a question of cost versus performance. [53] posed attribute-based person search as a cross-modality matching problem. They applied an adversarial leaning method to align the cross-modal distributions in a common space. However, their model design did not extend to the unseen class problem or few-shot learning challenges. Dong et al. [7] formulated a hierarchical matching model that fuses the similarity between global category-level embeddings and local attribute-level embeddings. Although they consider a zero-shot learning paradigm in the approach, the main disadvantage of this method is that it lacks the ability to synthesize visual representations of an unseen category. Therefore, it does not successfully handle unseen class problem for cross-modal matching.

**Cross-Modal Retrieval.** Our task of searching for people using attribute descriptions is closely related to studies on cross-modal retrieval as both problems require the attribute descriptions to be aligned with image data. This is a particularly relevant task to text-image embedding [52, 49, 26, 44, 57, 41], where canonical correlation analysis (CCA) [12] is a common solution for aligning two modalities by maximizing their correlation. With the rapid development of deep neural networks (DNN), a deep canonical correlation analysis (DCCA) model has since been proposed based on the same insight [1]. Yan et al. [52] have subsequently extended the idea to an image-text matching method based on DCCA. Beyond these core works, a variety of cross-modal retrieval methods have been proposed for different ways of learning a common space. Most use category-level (categories) labels to learn common representations [53, 44, 63, 57, 57]. However, what might be fine-grained attribute representations often lose granularity, and semantic categories are all treated as being the same distance apart. Several more recent works are based on using a GAN [11] for cross-modal alignment in the common subspace [44, 53, 63, 42]. The idea is intuitive since GANs have been proved to be powerful tools for distribution alignment [42, 15]. Further, they have produced some impressive results in image translation [62], domain adaptation [42, 15] and so on. However, when only applied to a common space, conventional adversarial learning may fail to model data variance where there are only a few samples and, again, the model design ignores the zero-shot learning challenge.

## 3    Symbiotic Adversarial Learning (SAL)

**Problem Definition**  Our deep model for attribute-based person search is based on the following specifications. There are $N$ labeled training images $\{\boldsymbol{x}_i, \boldsymbol{a}_i, y_i\}_{i=1}^N$ available in the training set. Each image-level attribute description $\boldsymbol{a}_i = [a_{(i,1)}, \ldots, a_{(i,n_{\text{attr}})}]$ is a binary vector with $n_{\text{attr}}$ number of attributes, where 0 and 1 indicate the absence and presence of the corresponding attribute. Images of people with the same attribute description are assigned to a unique category so as to derive a category-level label, specifically $y_i \in \{1, ..., M\}$ for $M$ categories. The aim is to find matching pedestrian images from the gallery set $\mathcal{X}_{\text{gallery}} = \{\boldsymbol{x}_j\}_{j=1}^G$ with $G$ images given an attribute description $\boldsymbol{a}_q$ from the query set.

### 3.1    Multi-modal Common Space Embedding Base

One general solution for cross-modal retrieval problems is to learn a common joint subspace where samples of different modalities can be directly compared to each other. As mentioned in our literature review, most current approaches use category-level labels (categories) to generate common representations [53, 44, 63]. However, these representations never manage to preserve the full granularity of finely-nuanced attributes. Plus, all semantic categories are treated as being the same distance apart when, in reality, the distances between different semantic categories can vary considerably depending on the similarity of their attribute descriptions.

Thus, we propose a common space learning method that jointly considers the global category and the local fine-grained attributes. The common space embedding loss is defined as:

$$L_{\text{embed}} = L_{\text{cat}} + L_{\text{att}}. \tag{1}$$

The category loss utilizing the Softmax Cross-Entropy loss function for image embedding branch is defined as:

$$L_{\text{cat}} = -\sum_{i=1}^{N} \log \Big( p_{\text{cat}}(\boldsymbol{x}_i, y_i) \Big), \tag{2}$$

where $p_{\text{cat}}(\boldsymbol{x}_i, y_i)$ specifies the predicted probability on the ground truth category $y_i$ of the training image $\boldsymbol{x}_i$.

To make the common space more discriminative for fine-trained attributes, we use the Sigmoid Cross-Entropy loss function which considers all $m$ attribute classes:

$$L_{\text{att}} = -\sum_{i=1}^{N} \sum_{j=1}^{m} \Big( a_{(i,j)} \log \big( p_{\text{att}}^{(j)}(\boldsymbol{x}_i) \big) + (1 - a_{(i,j)}) \log \big( 1 - p_{\text{att}}^{(j)}(\boldsymbol{x}_i) \big) \Big), \tag{3}$$

where $a_{(i,j)}$ and $p_{\text{att}}^{(j)}(\boldsymbol{x}_i)$ define the ground truth label and the predicted classification probability on the $j$-th attribute class of the training image $\boldsymbol{x}_i$, i.e.
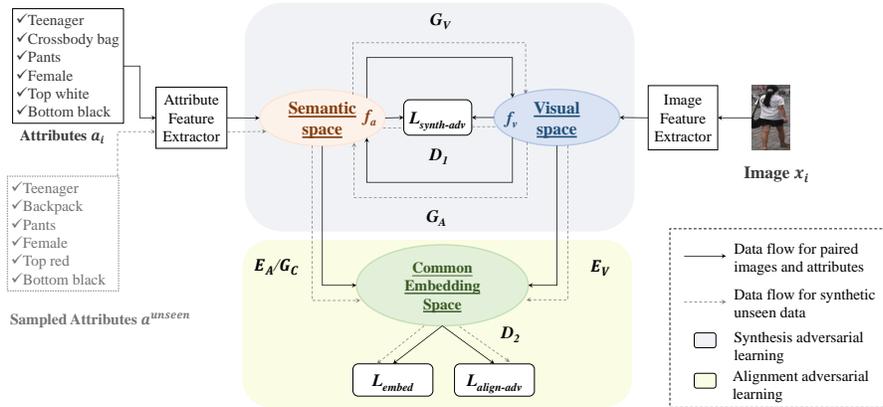
**Fig. 3.** An overview of the proposed SAL Framework.

$a_i = [a_{(i,1)}, \cdots, a_{(i,m)}]$ and $p_{\text{att}}^{(j)} = [p_{\text{att}}^{(1)}(x_i), \cdots, p_{\text{att}}^{(m)}(x_i)]$. Similar loss functions of Eq. (2) and Eq. (3) are applied to the attribute embedding branch as well.

The multi-modal common space embedding can be seen as our base, named as *Embed* (Fig. 1(a)). A detailed comparison is shown in Table 2.

### 3.2 Middle-Level Granularity-Consistent Cycle Generation

To address the challenge of insufficient data and bridge the modality gap between middle-level visual and semantic representations, we aim to generate synthetic features from the semantic representations to supplement the visual representations and vice versa. The synthetic features are conditioned on common embedding information so as to preserve very fine levels of granularity for retrieval. The middle-level representations are denoted as $f_v$ for visual and $f_a$ for attribute. More detail is provided in Fig. 3.

As paired supervision are available for seen categories during generation, conventional unsupervised cross-domain generation methods, e.g. CycleGAN [62], DiscoGAN [17] are not perfect suit. To better match joint distributions, we consider single discriminator distinguishes whether the paired data $(f_a, f_v)$ ) is from a real feature distribution $p(f_a, f_v)$ or not. This is inspired by Triple Generative Adversarial Networks [4]. The TripleGAN works on semi-supervised classification task with a generator, a discriminator and a classifier while our synthesis-GAN consists of two generators and a discriminator. As shown in Fig. 3, our synthesis-GAN consists of three main components: Generator $G_A$ synthesizes semantic representations from visual representations; Generator $G_V$ synthesizes visual representation from semantic representations, and a discriminator $D_1$ focused on identifying synthetic data pairs. Given that semantic representations are rather more sparse than visual representations, we add a noise vector $z \sim \mathcal{N}(0, 1)$ sampled from a Gaussian distribution for $G_V$ to get variation

in visual feature generation. Thus, $G_V$ can be regarded as a one-to-many matching from a sparse semantic space to a rich visual space. This accords with the one-to-many relationship between attributes and images. The training process is formulated as a minmax game between three players - $G_V$, $G_A$ and $D_1$ - where, from a game theory perspective, both generators can achieve their optima. The adversarial training scheme for middle-level cross-modal feature generation can, therefore, be formulated as:

$$
\begin{aligned}
L_{\text{gan1}}(G_A, G_V, D_1) &= \mathbb{E}_{(f_a, f_v) \sim p(f_a, f_v))}[log(D_1(f_a, f_v))] \\
&+ \frac{1}{2}\mathbb{E}_{f_a \sim p(f_a)}[log(1 - D_1(f_a, \widetilde{f}_v))] + \frac{1}{2}\mathbb{E}_{f_v \sim p(f_v)}[log(1 - D_1(\widetilde{f}_a, f_v))].
\end{aligned}
\tag{4}
$$

The discriminator $D_1$ is fed with three types of input pairs: (1) The fake input pairs $(\widetilde{f}_a, f_v)$, where $\widetilde{f}_a = G_A(f_v)$. (2) The fake input pairs $(f_a, \widetilde{f}_v)$, where $\widetilde{f}_v = G_V(f_a, z)$. (3) The real input pairs $(f_a, f_v) \sim p(f_a, f_v)$.

To constrain the generation process to only produce representations that are close to the real distribution, we assume the synthetic visual representation can be generated back to the original semantic representation, which is inspired by the CycleGAN structure [62]. Given this is a one-to-many matching problem as mentioned above, we flex the ordinary two-way cycle consistency loss to a one-way cycle consistency loss (from the semantic space to the visual space and back again):

$$
L_{\text{cyc}}(G_A, G_V) = \mathbb{E}_{f_a \sim p(f_a)}[||G_A(G_V(f_a, z)) - f_a||_2].
\tag{5}
$$

Furthermore, feature generation is conditioned on embeddings in the high-level common space, with the aim of generating more meaningful representations that preserve as much granularity as possible. This constraint is a new regularization term that we call the common space granularity-consistency loss, formulated as follows:

$$
\begin{aligned}
L_{\text{consis}}(G_A, G_V) &= \mathbb{E}_{f_v \sim p(f_v)}[||E_A(\widetilde{f}_a) - E_V(f_v)||_2] + \mathbb{E}_{f_a \sim p(f_a)}[||E_V(\widetilde{f}_v) - E_A(f_a)||_2] \\
&+ \mathbb{E}_{(f_a, f_v) \sim p(f_a, f_v)}[||E_A(\widetilde{f}_a) - E_A(f_a)||_2] + \mathbb{E}_{(f_a, f_v) \sim p(f_a, f_v)}[||E_V(\widetilde{f}_v) - E_V(f_v)||_2],
\end{aligned}
\tag{6}
$$

where $E_A$ and $E_V$ stand for the common space encoders for attribute features and visual features respectively.

Lastly, the full objective of synthesis-GAN is:

$$
L_{\text{synth-adv}} = L_{\text{gan1}} + L_{\text{cyc}} + L_{\text{consis}}.
\tag{7}
$$

### 3.3   High-Level Common Space Alignment with Augmented Adversarial Learning

On the high-level common space which is optimized by Eq. (1), we introduce the second adversarial loss for cross-modal alignment, making the common space

more modality-invariant. Here, $E_A$ works as a synchronous common space encoder and generator. Thus, we can define $E_A$ as $G_C$ in this adversarial loss:

$$L_{\text{gan2}}(G_C, D_2) = \mathbb{E}_{f_v \sim p(f_v)}[logD_2(E_V(f_v))] \\ + \mathbb{E}_{f_a \sim p(f_a)}[log(1 - D_2(E_A(f_a)))]. \tag{8}$$

Benefit from the middle-level GAN that bridges the semantic and visual spaces, we can access the augmented data, $\widetilde{f}_a$ and $\widetilde{f}_v$, to optimise the common space, where $\widetilde{f}_a = G_A(f_v)$ and $\widetilde{f}_v = G_V(f_a, z)$. Thus the augmented adversarial loss is applied for cross-modal alignment by:

$$L_{\text{aug1}}(G_C, D_2) = \mathbb{E}_{f_a \sim p(f_a)}[logD_2(E_V(\widetilde{f}_v))] \\ + \mathbb{E}_{f_v \sim p(f_v)}[log(1 - D_2(E_A(\widetilde{f}_a)))]. \tag{9}$$

And can be used to optimize the common space by:

$$L_{\text{aug2}}(E_A, E_V) = L_{\text{embed}}(\widetilde{f}_a) + L_{\text{embed}}(\widetilde{f}_v). \tag{10}$$

The total augmented loss from the synthesis-GAN interaction is calculated by:

$$L_{\text{aug}} = L_{\text{aug1}} + L_{\text{aug2}}. \tag{11}$$

The final augmented adversarial loss for alignment-GAN is defined as:

$$L_{\text{align-adv}} = L_{\text{gan2}} + L_{\text{aug}}. \tag{12}$$

### 3.4   Symbiotic Training Scheme for SAL

Algorithm 1 summarizes the training process. As stated above, there are two GANs learn collaboratively and specially throughout the training process. The synthesis-GAN (including $G_A$, $G_V$ and $D_1$) aims to build a bridge of semantic and visual representations in the middle-level feature space, and to synthesize features from both to each other. The granularity-consistency loss was further proposed to constrain the generation that conditions on high-level embedding information. Thus, the high-level alignment-GAN (including $E_A$, $E_V$ and $D_2$) that optimizes the common embedding space benefits the synthesis-GAN with the better common space constrained by Eq. (6). Similarly, while the alignment-GAN attempts to optimize the embedding and perform the cross-modal alignment, the synthesis-GAN supports the alignment-GAN with more realistic and reliable augmented data pairs via Eqs. (9) and (10). Thus, the two GANs update iteratively with SALs full objective becomes:

---

**Algorithm 1** Learning the SAL model.

---

   **Input:** $N$ labeled data pairs $(\boldsymbol{x}_i, \boldsymbol{a}_i)$ from the training set;

   **Output:** SAL attribute-based person search model;

   **for** $t = 1$ **to** *max-iteration* **do**

      Sampling a batch of paired training data $(\boldsymbol{x}_i, \boldsymbol{a}_i)$;

      **Step 1: Multi-modal common space Embedding**:

      Updating both image branch and attribute branch by common space embedding loss $L_{\text{embed}}$ (Eq. (1));

      **Step 2: Middle-Level Granularity-Consistent Cycle Generation** : Updating $G_A$, $G_V$ and $D_1$ by $L_{\text{synth-adv}}$ (Eq. (7));

      **Step 3: High-Level Common Space Alignment**: Updating $E_A$, $E_V$ and $D_2$ by $L_{\text{align-adv}}$ (Eq. (12));

   **end for**

---

$$L_{\text{SAL}} = L_{\text{embed}} + L_{\text{synth-gan}} + L_{\text{align-gan}}. \tag{13}$$

**Semantic augmentation for unseen classes.** To address the zero-shot problem, our idea is to synthesize visual representations of unseen classes from semantic representations. In contrast to conventional zero-shot settings with pre-defined category names, we rely on a myriad of attribute combinations instead. Further, this inspired us to also sample new attributes in the model design. Hence, during training, some new attribute combinations $\boldsymbol{a}^{unseen}$ are dynamically sampled in each iteration. The sampling for the binary attributes follows a Bernoulli distribution, and a 1-trial multinomial distribution for the multi-valued attributes (i.e., mutually-exclusive attributes) according to the training data probability. Once the attribute feature extraction is complete, $f_a$ in the objective function is replaced with $[f_a, f_{a^{unseen}}]$, $f_{a^{unseen}}$ is used to generate synthetic visual feature via $G_V$. The synthetic feature pairs are then used for the final common space cross-modal learning.

## 4   Experiments

**Datasets.** To evaluate SAL with an attribute-based person search task, we used two widely adopted pedestrian attribute datasets[1]: (1) The ***Market-1501 attribute*** dataset [56] which the Market-1501 dataset annotated with 27 attributes for each identity [27] - see Fig. 4. The training set contains an average of $\sim 26$ labeled images in each category. During testing, 367 out of 529 ( **69.4%**) were from **unseen** categories. (2) The ***PETA*** dataset [5] consists of $19,000$ pedestrian images collected from 10 small-scale datasets of people. Each image is annotated with 65 attributes. Following [53], we labeled the images with categories according to their attributes, then split the dataset into a training set and a gallery set by category. In the training set, there was on average $\sim 7$ labeled images in each category, and all 200 categories in the gallery set (**100%**) were **unseen**.

---

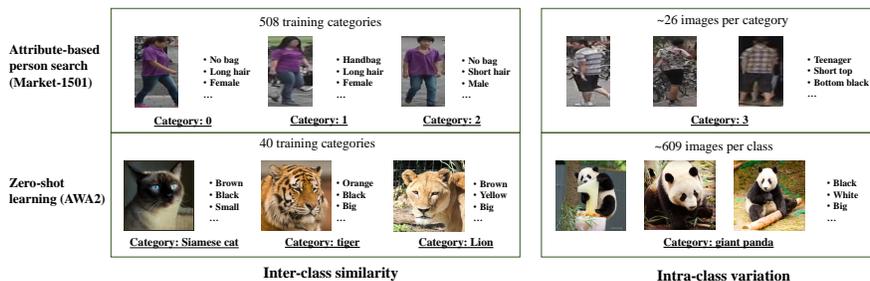[1] The DukeMTMC dataset is not publicly available.

**Fig. 4.** Example images from the person search Market-1501 [56] dataset and compared with the general zero-shot AWA2 [51] dataset.

**Table 1.** Attribute-based person search performance evaluation. Best results are shown in **bold**. The second-best results are <u>underlined</u>.

| Metric (%) | | Market-1501 Attributes | | | | PETA | | | |
|---|---|---|---|---|---|---|---|---|---|
| Model | Reference | mAP | rank1 | rank5 | rank10 | mAP | rank1 | rank5 | rank10 |
| DeepCCA [1] | ICML'13 | 17.5 | 30.0 | 50.7 | 58.1 | 11.5 | 14.4 | 20.8 | 26.3 |
| DeepMAR [23] | ACPR'15 | 8.9 | 13.1 | 24.9 | 32.9 | 12.7 | 17.8 | 25.6 | 31.1 |
| DeepCCAE [50] | ICML'15 | 9.7 | 8.1 | 24.0 | 34.6 | 14.5 | 14.2 | 22.1 | 30.0 |
| 2WayNet [8] | CVPR'17 | 7.8 | 11.3 | 24.4 | 31.5 | 15.4 | 23.7 | 38.5 | 41.9 |
| CMCE [24] | ICCV'17 | 22.8 | 35.0 | 51.0 | 56.5 | 26.2 | 31.7 | 39.2 | 48.4 |
| ReViSE [41] | ICCV'17 | 17.7 | 24.2 | 45.2 | 57.6 | 31.1 | 30.5 | <u>57.0</u> | 61.5 |
| MMCC [9] | ECCV'18 | 22.2 | 34.9 | <u>58.7</u> | <u>70.2</u> | <u>33.9</u> | 33.5 | <u>57.0</u> | <u>69.0</u> |
| AAIPR [53] | IJCAI'18 | 20.7 | 40.3 | 49.2 | 58.6 | 27.9 | <u>39.0</u> | 53.6 | 62.2 |
| AIHM [7] | ICCV'19 | <u>24.3</u> | <u>43.3</u> | 56.7 | 64.5 | - | - | - | - |
| **SAL** (Ours) | | **29.8** | **49.0** | **68.6** | **77.5** | **41.2** | **47.0** | **66.5** | **74.0** |

**Performance Metric.** The two metrics used to evaluate performance were the cumulative matching characteristic (CMC) and mean Average Precision (mAP). We followed [31] and computed the CMC on each rank k as the probe cumulative percentage of truth matches appearing at ranks $\leq k$. mAP measures the recall of multiple truth matches and was computed by first computing the area under the Precision-Recall curve for each probe, then calculating the mAP over all probes. **Implementation Details.** SAL was implemented based on Torchreid [59, 61, 60] in the Pytorch framework [32] with ResNet-50 [13] as the image feature extractor. We used fully connected (FC) layers to form the attribute feature extractor (64, 128, 256, 512), the middle-level generators (256, 128, 256, 512), and the encoders (512, 256, 128). Batch normalization and a ReLU nonlinear activation followed each of the first three layers with a Tanh (activation) before the output. The discriminators were also FC layers (512, 256, 1) with batch normalization and a leaky ReLU activation following each of the first two layers, and a Sigmoid activation prior to the output. After the common space embedding, the classifiers output with 1-FC layer. **Training process:** We first pre-trained the image branch and *Embed* model as a person search baseline (Fig. 3.1) until

it converges. Then we trained the full SAL model for 60 epochs with a learning rate of 0.001 for the image branch and 0.01 for the attribute branch. We chose Adam as the training optimizer [18] with a batch size of 128. During testing, we calculated the cosine similarity between the query attributes and the gallery images in the common embedding space with 128-D deep feature representations for matching. **Training time:** It took 77 minutes for SAL to converge (26.3M parameters) compared to 70 minutes for a single adversarial model Embed+adv to converge (25.0M parameters). Both training processes were run on the same platform with 2 NVIDIA 1080Ti GPUs.

### 4.1   Comparisons to the State-Of-The-Arts

The results of the search task appear in Table 1, showing SALs performance in comparison to 9 state-of-the-art models across 4 types of methods. These are: (I) *attribute recognition* based method: (1) DeepMAR [23]. (II) *correlation* based methods: (2) Deep Canonical Correlation Analysis (DCCA) [1]; (3) Deep Canonically Correlated Autoencoders (DCCAE) [50]; (4) 2WayNet [8] (III) *common space embedding* : (5) Cross-modality Cross-entropy (CMCE) [24]; (6) ReViSE [41]; (7) Attribute-Image Hierarchical Matching (AIHM) [7]; (IV) *adversarial learning*: (8) Multi-modal Cycle-consistent (MMCC) model [9]; (9) Adversarial Attribute-Image Person Re-ID (AAIPR) [53].

From the results, we made the following observations: (1) SAL outperformed all 9 state-of-the-art models on both datasets in terms of mAP. On the Market-1501 dataset, the improvement over the second best scores were 5.5% (29.8 v 24.3) and 5.7% (49.0 v 43.3) for rank1. On the PETA dataset, the improvement was 7.3% (41.2 v 33.9) and 8.0% (47.0 v 39.0) for rank1. This illustrates SALs overall performance advantages in cross-modal matching for attribute-based person search. (2) Directly predicting the attributes using the existing recognition models and matching the predictions with the queries is not efficient (e.g., DeepMAR). This may be due to the relatively low prediction dimensions and the sparsity problems with attributes in semantic space. (3) Compared to the common space learning-based method (CMCE), the conventional correlation methods (e.g., DCCA, DCCAE, 2WayNet) witnessed relatively poor results. This demonstrates the power of common space learning with a common embedding space. It is worth mentioning that CMCE is specifically designed to search for people using natural language descriptions, yet SAL outperformed CMCE by 7.0% mAP/14.0% rank1 with Market-1501, and by 15.0% mAP/15.3% rank1 with PETA. (4) The adversarial model comparisons, AAIPR and MMCC, did not fare well against SAL, especially AAIPR, which utilizes single adversarial learning to align the common space. This directly demonstrates the advantages of our approach with symbiotic adversarial design compared to traditional adversarial learning methods. (5) Among the compared state-of-the-arts, MMCC, ReViSE and AIHM addressed *zero-shot problems*. Against MMCC, SAL outperformed by 7.6% mAP/14.1% rank1 on Market-1501 and 7.3% mAP/13.5% rank1 on PETA. AIHM is the most recent state-of-the-art in this category of methods

**Table 2.** Component analysis of SAL on PETA dataset.

| Metric (%) | mAP | rank1 | rank5 | rank10 |
|---|---|---|---|---|
| *Embed* | 31.3 | 34.0 | 57.0 | 64.5 |
| *Embed + adv* | 35.0 | 37.5 | 60.5 | 66.5 |
| *Embed + symb-adv* | 40.6 | 44.0 | 64.0 | 70.5 |
| *Embed + symb-adv + unseen*(SAL) | 41.2 | 47.0 | 66.5 | 74.0 |

**Table 3.** Effect of interactions between two GANs on PETA dataset.

| Metric (%) | mAP | rank1 | rank5 | rank10 |
|---|---|---|---|---|
| SAL - $L_{aug}$ | 35.4 | 38.0 | 60.0 | 69.0 |
| SAL - $L_{consis}$ | 35.2 | 39.5 | 56.5 | 66.0 |
| SAL (Full interaction) | 41.2 | 47.0 | 66.5 | 74.0 |

and SALs performance improvement was 5.5% mAP/5.7% rank1 on Market-1501. This demonstrates the advantages of SALs new regularization term, the common space granularity-consistency loss, for generating middle-level features.

### 4.2 Further Analysis and Discussions

To further evaluate the different components in the model, we conduct studies on the PETA dataset.

**Component analysis of SAL** Here, we compared: (1) the base embedding model (*Embed*, Fig. 1(a)), which comprises the multi-modal common space embedding base (Sec. 3.1) and is optimized by embedding loss (Eq. (1)) only; (2) the base embedding model plus single adversarial learning (*Embed+adv*, Fig. 1(b)), (3) the base embedding model plus our symbiotic adversarial learning (*Embed+symb-adv*, Fig. 1(c)), and (4) our full SAL model, which includes the attribute sampling for unseen classes. The results are shown in Table 2.

Compared to the *Embed* model, *Embed+adv* saw an improvement of 3.7% mAP/3.5% rank1, whereas *Embed+symb-adv* achieved an improvement of 9.3% mAP/10% rank1, and finally SAL (*Embed+symb-adv+unseen*) witnessed a significant improvement of 9.9% mAP/13.0% rank1. This is a clear demonstration of the benefits of jointly considering middle-level feature augmentation and high-level cross-modal alignment instead of only having common space cross-modal alignment (41.2% vs 35.0% mAP and 47.0% vs 37.5% rank1). We visualize the retrieved results in the supplementary material.

**Effect of interactions between two GANs** Our next test was designed to assess the influence of the symbiotic interaction, i.e., where the two GANs iteratively regularize and augment each other. Hence, we removed the interaction loss between the two GANs, which is the total augmented loss from Eq. (11) and the common space granularity-consistency loss from Eq. (6). Removing the augmented loss (SAL - $L_{aug}$) means the model no longer uses the augmented

**Table 4.** Comparing stage-wise training vs. symbiotic training scheme.

| Metric (%) | mAP | rank1 | rank5 | rank10 |
|---|---|---|---|---|
| SAL w/ stage-wise training | 35.0 | 41.0 | 58.0 | 65.0 |
| SAL w/ symbiotic training | 41.2 | 47.0 | 66.5 | 74.0 |

data. Removing the common space granularity-consistency loss (SAL - $L_{consis}$) means the middle-level generators are not conditioned on high-level common embedding information, which should reduce the augmented data quality. The results in Table 3 show that, without augmented data, SALs mAP decreased by 5.8% and by 6.0% without the granularity-consistency loss as a regularizer.

**Comparing stage-wise training vs. symbiotic training scheme**    We wanted to gain a better understanding of the power of the symbiotic training scheme (Sec. 3.4). So, we replaced the iterative symbiotic training scheme (SAL w/ symbiotic training) with a stage-wise training (SAL w/ stage-wise training). The stage-wise training [2] breaks down the learning process into sub-tasks that are completed stage-by-stage. So during implementation, we first train the synthesis-GAN conditioned on the common embedding. Then for the next stage, we optimise the alignment-GAN on the common space with synthetic data. As shown in Table 4, there was a 6.2% drop in mAP using the stage-wise training, which endorses the merit of the symbiotic training scheme. During the symbiotic training, the synthesized augmentation data and the common space alignment iteratively boost each others learning. A better common space constrains the data synthesis to generate better synthetic data. And, with better synthetic data for augmentation, a better common space can be learned.

## 5    Conclusion

In this work, we presented a novel symbiotic adversarial learning model, called SAL. SAL is an end-to-end framework for attribute-based person search. Two GANs sit at the base of the framework: one synthesis-GAN uses semantic representations to generate synthetic features for visual representations, and vice versa, while the other alignment-GAN optimizes the embeddings and performs cross-modal alignment on the common embedding space. Benefiting from the symbiotic adversarial structure, SAL is able to preserve finely-grained discriminative information and generate more reliable synthetic features to optimize the common embedding space. With the ability to synthesize visual representations of unseen classes, SAL is more robust to zero-shot retrieval scenarios, which are relatively common in real-world person search with diverse attribute descriptions. Extensive ablation studies illustrate the insights of our model designs. Further, we demonstrate the performance advantages of SAL over a wide range of state-of-the-art methods on two challenging benchmarks.

# References

1. Andrew, G., Arora, R., Bilmes, J., Livescu, K.: Deep canonical correlation analysis. In: ICML (2013)
2. Barshan, E., Fieguth, P.: Stage-wise training: An improved feature learning strategy for deep models. In: Feature Extraction: Modern Questions and Challenges. pp. 49–59 (2015)
3. Chen, Y.C., Zhu, X., Zheng, W.S., Lai, J.H.: Person re-identification by camera correlation aware feature augmentation. IEEE TPAMI **40**(2) (2018)
4. Chongxuan, L., Xu, T., Zhu, J., Zhang, B.: Triple generative adversarial nets. In: NIPS (2017)
5. Deng, Y., Luo, P., Loy, C.C., Tang, X.: Pedestrian attribute recognition at far distance. In: ACM MM. ACM (2014)
6. Deng, Y., Luo, P., Loy, C.C., Tang, X.: Learning to recognize pedestrian attribute. arXiv:1501.00901 (2015)
7. Dong, Q., Gong, S., Zhu, X.: Person search by text attribute query as zero-shot learning. In: Proceedings of the IEEE International Conference on Computer Vision. pp. 3652–3661 (2019)
8. Eisenschtat, A., Wolf, L.: Linking image and text with 2-way nets. In: CVPR (2017)
9. Felix, R., Kumar, V.B., Reid, I., Carneiro, G.: Multi-modal cycle-consistent generalized zero-shot learning. In: ECCV (2018)
10. Feng, F., Wang, X., Li, R.: Cross-modal retrieval with correspondence autoencoder. In: ACM MM. ACM (2014)
11. Goodfellow, I., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A., Bengio, Y.: Generative adversarial nets. In: NIPS (2014)
12. Hardoon, D.R., Szedmak, S., Shawe-Taylor, J.: Canonical correlation analysis: An overview with application to learning methods. Neural computation **16**(12) (2004)
13. He, K., Zhang, X., Ren, S., Sun, J.: Deep residual learning for image recognition. In: CVPR (2016)
14. Hermans, A., Beyer, L., Leibe, B.: In defense of the triplet loss for person re-identification. arXiv:1703.07737 (2017)
15. Hoffman, J., Tzeng, E., Park, T., Zhu, J.Y., Isola, P., Saenko, K., Efros, A.A., Darrell, T.: Cycada: Cycle-consistent adversarial domain adaptation. arXiv:1711.03213 (2017)
16. Jaha, E.S., Nixon, M.S.: Soft biometrics for subject identification using clothing attributes. In: IJCB. IEEE (2014)
17. Kim, T., Cha, M., Kim, H., Lee, J.K., Kim, J.: Learning to discover cross-domain relations with generative adversarial networks. In: Proceedings of the 34th International Conference on Machine Learning-Volume 70. pp. 1857–1865. JMLR. org (2017)
18. Kingma, D.P., Ba, J.: Adam: A method for stochastic optimization. arXiv:1412.6980 (2014)
19. Kumar Verma, V., Arora, G., Mishra, A., Rai, P.: Generalized zero-shot learning via synthesized examples. In: Proceedings of the IEEE conference on computer vision and pattern recognition. pp. 4281–4289 (2018)
20. Layne, R., Hospedales, T.M., Gong, S.: Towards person identification and re-identification with attributes. In: ECCV. Springer (2012)
21. Layne, R., Hospedales, T.M., Gong, S.: Attributes-based re-identification. In: Person Re-Identification. Springer (2014)

22. Layne, R., Hospedales, T.M., Gong, S., Mary, Q.: Person re-identification by attributes. In: BMVC (2012)
23. Li, D., Chen, X., Huang, K.: Multi-attribute learning for pedestrian attribute recognition in surveillance scenarios. In: IAPR ACPR. IEEE (2015)
24. Li, S., Xiao, T., Li, H., Yang, W., Wang, X.: Identity-aware textual-visual matching with latent co-attention. In: ICCV (2017)
25. Li, W., Zhao, R., Xiao, T., Wang, X.: Deepreid: Deep filter pairing neural network for person re-identification. In: CVPR (2014)
26. Li, W., Zhu, X., Gong, S.: Person re-identification by deep joint learning of multi-loss classification. arXiv:1705.04724 (2017)
27. Lin, Y., Zheng, L., Zheng, Z., Wu, Y., Yang, Y.: Improving person re-identification by attribute and identity learning. arXiv:1703.07220 (2017)
28. Liu, C., Gong, S., Loy, C.C., Lin, X.: Person re-identification: What features are important? In: ECCV. Springer (2012)
29. 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: ICCV (2017)
30. Long, Y., Liu, L., Shao, L., Shen, F., Ding, G., Han, J.: From zero-shot learning to conventional supervised classification: Unseen visual data synthesis. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1627–1636 (2017)
31. Paisitkriangkrai, S., Shen, C., Van Den Hengel, A.: Learning to rank in person re-identification with metric ensembles. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. pp. 1846–1855 (2015)
32. Paszke, A., Gross, S., Chintala, S., Chanan, G., Yang, E., DeVito, Z., Lin, Z., Desmaison, A., Antiga, L., Lerer, A.: Automatic differentiation in pytorch. In: NIPS-W (2017)
33. Reid, D.A., Nixon, M.S., Stevenage, S.V.: Soft biometrics; human identification using comparative descriptions. IEEE TPAMI **36**(6) (2014)
34. Saquib Sarfraz, M., Schumann, A., Eberle, A., Stiefelhagen, R.: A pose-sensitive embedding for person re-identification with expanded cross neighborhood re-ranking. In: CVPR (2018)
35. Scheirer, W.J., Kumar, N., Belhumeur, P.N., Boult, T.E.: Multi-attribute spaces: Calibration for attribute fusion and similarity search. In: CVPR. IEEE (2012)
36. Shi, Z., Hospedales, T.M., Xiang, T.: Transferring a semantic representation for person re-identification and search. In: CVPR (2015)
37. Siddiquie, B., Feris, R.S., Davis, L.S.: Image ranking and retrieval based on multi-attribute queries. In: CVPR. IEEE (2011)
38. Su, C., Yang, F., Zhang, S., Tian, Q., Davis, L.S., Gao, W.: Multi-task learning with low rank attribute embedding for person re-identification. In: ICCV (2015)
39. Su, C., Zhang, S., Xing, J., Gao, W., Tian, Q.: Deep attributes driven multi-camera person re-identification. In: ECCV. Springer (2016)
40. Sun, Y., Zheng, L., Deng, W., Wang, S.: Svdnet for pedestrian retrieval. In: ICCV (2017)
41. Tsai, Y.H.H., Huang, L.K., Salakhutdinov, R.: Learning robust visual-semantic embeddings. In: ICCV. IEEE (2017)
42. Tzeng, E., Hoffman, J., Saenko, K., Darrell, T.: Adversarial discriminative domain adaptation. In: CVPR (2017)
43. Vaquero, D.A., Feris, R.S., Tran, D., Brown, L., Hampapur, A., Turk, M.: Attribute-based people search in surveillance environments. In: WACV. IEEE (2009)

44. Wang, B., Yang, Y., Xu, X., Hanjalic, A., Shen, H.T.: Adversarial cross-modal retrieval. In: ACM MM. ACM (2017)
45. Wang, F., Zuo, W., Lin, L., Zhang, D., Zhang, L.: Joint learning of single-image and cross-image representations for person re-identification. In: CVPR (2016)
46. Wang, J., Zhu, X., Gong, S., Li, W.: Attribute recognition by joint recurrent learning of context and correlation. In: ICCV (2017)
47. Wang, J., Zhu, X., Gong, S., Li, W.: Transferable joint attribute-identity deep learning for unsupervised person re-identification. In: CVPR (2018)
48. Wang, K., He, R., Wang, W., Wang, L., Tan, T.: Learning coupled feature spaces for cross-modal matching. In: ICCV (2013)
49. Wang, L., Li, Y., Lazebnik, S.: Learning deep structure-preserving image-text embeddings. In: CVPR (2016)
50. Wang, W., Arora, R., Livescu, K., Bilmes, J.: On deep multi-view representation learning. In: ICML (2015)
51. Xian, Y., Lampert, C.H., Schiele, B., Akata, Z.: Zero-shot learninga comprehensive evaluation of the good, the bad and the ugly. IEEE transactions on pattern analysis and machine intelligence $\mathbf{41}(9)$, 2251–2265 (2018)
52. Yan, F., Mikolajczyk, K.: Deep correlation for matching images and text. In: CVPR (2015)
53. 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: IJCAI (7 2018)
54. Zhang, L., Xiang, T., Gong, S.: Learning a discriminative null space for person re-identification. In: CVPR (2016)
55. Zhao, X., Sang, L., Ding, G., Guo, Y., Jin, X.: Grouping attribute recognition for pedestrian with joint recurrent learning. In: IJCAI (2018)
56. Zheng, L., Shen, L., Tian, L., Wang, S., Wang, J., Tian, Q.: Scalable person re-identification: A benchmark. In: ICCV (2015)
57. Zheng, Z., Zheng, L., Garrett, M., Yang, Y., Shen, Y.D.: Dual-path convolutional image-text embedding with instance loss. arXiv:1711.05535 (2017)
58. Zhong, Z., Zheng, L., Cao, D., Li, S.: Re-ranking person re-identification with k-reciprocal encoding. In: CVPR (2017)
59. Zhou, K., Xiang, T.: Torchreid: A library for deep learning person re-identification in pytorch. arXiv preprint arXiv:1910.10093 (2019)
60. Zhou, K., Yang, Y., Cavallaro, A., Xiang, T.: Learning generalisable omni-scale representations for person re-identification. arXiv preprint arXiv:1910.06827 (2019)
61. Zhou, K., Yang, Y., Cavallaro, A., Xiang, T.: Omni-scale feature learning for person re-identification. In: ICCV (2019)
62. Zhu, J.Y., Park, T., Isola, P., Efros, A.A.: Unpaired image-to-image translation using cycle-consistent adversarial networks. In: ICCV (2017)
63. Zhu, Y., Elhoseiny, M., Liu, B., Peng, X., Elgammal, A.: A generative adversarial approach for zero-shot learning from noisy texts. In: CVPR (2018)

# Supplementary material

Query: Teenager, **No bag**, Short top, Short bottom, Short hair, Male, Top gray, Bottom black

Query: Adult, **Handbag**, Short top, **Long bottom**, **Long hair**, **Female**, Top black, **Bottom blue**

Query: Teenager, **Backpack**, Short top, **Short bottom**, Short hair, Male, Top gray, Bottom white

Query: Adult, **Backpack**, Short top, Long bottom, **Long hair**, **Female**, **Top red**, Bottom black

**Fig. 5.** Ranked retreival results. The query attributes are shown above the retrieved images. The green/red border represents correct/wrong selections respectively. The attributes in **bold** correspond to false matches.

Fig. 5 visualizes the ranked results from *Embed*, *Embed+adv*, *Embed+symb-adv* and SAL to qualitatively illustrate the performance of the three models. Although the models are able to pick out some correct images from candidates in the top 10 ranks, SAL selected more with higher ranks. Checking the wrong attributes of the false ranks, we can see that SAL was better able to discern fine-grained attributes, such as different kinds of bags. Moreover, SAL picked correct images of diverse appearances, which may indicate it has some ability to overcome intra-class variation problems.