

# Synthetic18K: Learning Better Representations for Person Re-ID and Attribute Recognition from 1.4 Million Synthetic Images

Onur Can Uner<sup>a</sup>, Cem Aslan<sup>b</sup>, Burak Ercan<sup>b</sup>, Tayfun Ates<sup>b</sup>, Ufuk Celikcan<sup>b</sup>,  
Aykut Erdem<sup>c</sup>, Erkut Erdem<sup>b,\*</sup>

<sup>a</sup>*Department of Computer Engineering, Bilkent University, Ankara, Turkey*

<sup>b</sup>*Department of Computer Engineering, Hacettepe University, Ankara, Turkey*

<sup>c</sup>*Department of Computer Engineering, Koç University, Istanbul, Turkey*

---

## Abstract

Learning robust representations is critical for the success of person re-identification and attribute recognition systems. However, to achieve this, we must use a large dataset of diverse person images as well as annotations of identity labels and/or a set of different attributes. Apart from the obvious concerns about privacy issues, the manual annotation process is both time consuming and too costly. In this paper, we instead propose to use synthetic person images for addressing these difficulties. Specifically, we first introduce Synthetic18K, a large-scale dataset of over 1 million computer generated person images of 18K unique identities with relevant attributes. Moreover, we demonstrate that pretraining of simple deep architectures on Synthetic18K for person re-identification and attribute recognition and then fine-tuning on real data leads to significant improvements in prediction performances, giving results better than or comparable to state-of-the-art models.

*Keywords:* person re-identification, attribute recognition, synthetic data

---

\*Corresponding author. Email address: erkut@cs.hacettepe.edu.tr

## 1. Introduction

In developed countries, video surveillance systems have become a vital component of public security, constantly monitoring cameras installed at various different key locations. Person re-identification (person re-ID) is one of the important tasks in video surveillance, which refers to the problem of automatically re-identifying across multiple different cameras with the help of computers. Different from other person-centric classification problems in computer vision such as face recognition, person re-ID systems use identity information only during training and assume the identities are unknown at test time. Person re-ID basically combines three different tasks, which are pedestrian detection, person tracking and person retrieval. Hence, it is extremely challenging due to large variations in lighting conditions, differences in pose and viewpoint changes.

The person re-ID models proposed in the literature can be divided into two main groups: image-based methods and video-based methods. All these approaches are usually evaluated on benchmark datasets annotated with either detected or ground truth human boxes so the problem mostly reduces to the person image retrieval where the aim is to retrieve images of a specific person identity from a large gallery of images involving many different identities. Therefore, success can be defined by comparing the identities of the retrieved images with the identity of the query image. In video-based approaches to person re-ID, the models use multiple bounding boxes for a video query and the videos in the gallery, and additionally try to integrate the temporal information between these boxes. Although research in both groups can benefit from each other, they are considered as different problems. In this work, we tackle the problem of image based re-ID.

As another critical task in video surveillance, person attribute recognition aims

26 at detecting various attributes of a person such as hair color, clothing type and  
27 color. Although person attribute recognition has been relatively less studied as  
28 compared to person re-ID, these tasks are, in fact, closely related since the mid-  
29 level semantic attributes, once identified, provide an intuitive way to describe a  
30 specific individual. Hence, a recent direction being explored in recent years is  
31 to consider these two challenging tasks in a joint manner in order to improve the  
32 performances of each other.

33 In the past few years, person re-ID research has reached a saturated point  
34 where researchers gently enhance the performances on popular benchmark datasets  
35 by designing more and more complex architectures or by applying complicated  
36 data augmentation schemes. That being said, most of the methods in the litera-  
37 ture fail to generalize well to in-the-wild settings because of the fact that existing  
38 datasets cover a limited range of samples that could be faced in real life. How-  
39 ever, obtaining a comprehensive dataset is very expensive to gather for which you  
40 have to use multiple camera sources located in very different environments. Even  
41 if you collect the right amount of visual data, the effort required for manual an-  
42 notation is quite costly. Hence, the existing datasets are not challenging enough  
43 in demonstrating different weather and lighting conditions, varieties in person at-  
44 tributes and/or body types. Finally, privacy of the individuals is an important  
45 issue for video surveillance. Although these datasets are initially collected for  
46 academic purposes, the intention of users may be different when the data become  
47 public, leading invasion of privacy of the individual. For example, DukeMTMC-  
48 reID dataset [1] has been recently shut down and cannot be downloaded publicly  
49 to conform to privacy regulations.

50 In this study, we deal with the problem of learning simple yet effective rep-

51 resentations for the person re-identification and attribute recognition tasks. In  
52 particular, in both of these two tasks, the main challenge lies in learning discrim-  
53 inative features which are not sensitive to the changes in the appearance of the  
54 person of interest due to illumination variations, scale and viewpoint changes.  
55 To overcome these difficulties, in this study, we first introduce a new synthetic  
56 dataset called Synthetic18K. The proposed dataset, compared to the existing syn-  
57 thetic datasets mentioned above, is much larger in scale in terms of the number  
58 of identities/virtual persons it contains and the number of images that each virtual  
59 person has. That is, it contains approximately 1.4 million images of 18K unique  
60 virtual persons captured in four synthetic environments (three outdoor and one  
61 indoor) as well as with various cubemaps taken in real-life. While generating  
62 these virtual persons and obtaining their images, we follow a procedural genera-  
63 tion method which gives us the ability to play with both the low and high-level  
64 attributes of these synthetic persons and the characteristics of the scenes (weather  
65 conditions, times of day, etc.). These aspects are of critical importance for feature  
66 learning as the existing real-world data are generally not diverse in various factors  
67 like illumination conditions, scenes, clothing, etc. are still considered as challeng-  
68 ing tasks. Moreover, covering each one of these factors in a dataset in a balanced  
69 manner could be very difficult to achieve, resulting in heavy-tailed data distribu-  
70 tions and poor performances for the rare cases. Tackling person re-identification  
71 and attribute recognition in a joint manner also introduces certain advantages as  
72 these tasks are considered complementary tasks. Yet, no other synthetic datasets  
73 handles these two in a combined manner.

74 Motivated with these, in our work, we also propose pretraining strategies and  
75 simple yet effective deep neural architectures for both person re-identification and

76 attribute recognition tasks. we show that our proposed Synthetic18K dataset can  
77 be used to learn more robust feature representations for these two tasks. In par-  
78 ticular, we proposed three different pretraining schemes, one for solely person re-  
79 identification, one for only attribute recognition and one final for a combination of  
80 these two tasks. We demonstrate that even simple neural architectures which are  
81 pretrained on our synthetically generated images using these strategies and later  
82 on fine-tuned on real data, perform competitively compared to complex state-of-  
83 the-art models. Our experiments also show handling person re-identification and  
84 attribute recognition together gives more accurate results than their single-task  
85 counterparts, indicating the importance of the proposed joint pretraining strategy.

86 Our dataset and models will be available at the project website<sup>1</sup>.

## 87 **2. Related Work**

88 In this section, we briefly review the literature on person re-ID and person  
89 attribute recognition, including the methods that perform these two tasks jointly,  
90 and mainly focusing on recent deep learning based techniques. Moreover, we look  
91 into existing re-ID datasets and their limitations.

### 92 *2.1. Person Re-Identification*

93 In recent years, person re-ID has emerged as a growing research topic, but its  
94 roots can be traced back to multi-camera tracking where the idea is to assign each  
95 person a unique latent label and try to re-identify them when they leave the scene  
96 and enter again to correctly keep tracking them. Hence, it is usually assumed  
97 that people wear the same clothes when they are captured by cameras. The early

---

<sup>1</sup><https://hucvl.github.io/synthetic18k>

98 approaches typically involve dividing the image into multiple image regions and  
99 represent each region with some local features such as color histograms and SIFT  
100 features and obtain a person representation by the concatenation of these features.  
101 One can apply metric learning to further increase the discriminative power of the  
102 features. A detailed analysis of these prior works can be found in [2, 3].

103 With the recent advancements in deep learning and the availability of large  
104 scale datasets, the aforementioned shallow models are replaced with their deep  
105 counterparts which combine feature learning and metric learning within a single  
106 framework. These approaches commonly use a pretrained convolutional neural  
107 network (CNN) such as ResNet-50 [4] trained on ImageNet [5] as a backbone  
108 and mostly differ from each other in their architectural details and the way their  
109 objectives are defined. As for the objectives, the most straightforward approach is  
110 to resort to classification loss where each person in the training set is treated as a  
111 distinct class [6]. More complicated objectives for Siamese architectures involve  
112 pair-wise contrastive loss [7], triplet ranking loss [8, 9] and quadruplet loss [10],  
113 or a combination of them [11, 12, 13]. Recent works even go beyond these ap-  
114 proaches and consider all the pairwise similarity relations within a batch com-  
115 posed of multiple identities and their images by using graph neural networks [14].

116 Solving person re-ID task requires features that have some invariance to scale  
117 changes [15], illumination and pose variations, viewpoint changes as well as mis-  
118 alignment errors in human detections [2, 3]. Apart from the training objectives  
119 mentioned above, some studies directly tackle these issues and design some spe-  
120 cialized architectures. Common trends include exploiting pose information by  
121 either using off-the-shelf pretrained pose detectors [16], or attaching a pose esti-  
122 mation network to the re-ID network and training them jointly [17]. Other alterna-

123 tives are learning to localize body parts [18, 19] or employing a human semantic  
124 parser network to additionally incorporate the body part features within the person  
125 re-ID network [20].

## 126 2.2. *Person Attribute Recognition*

127 Compared to person re-ID, person attribute recognition is a relatively less  
128 studied topic, but it has also gained attraction due to widespread interest in au-  
129 tomated visual surveillance systems. Traditional approaches typically employ  
130 hand-crafted features such as color histograms and involve classifiers trained inde-  
131 pendently for each attribute [21, 22]. However, attribute recognition is inherently  
132 a multi-label classification problem. Moreover, there are dependencies between  
133 some attributes. For example, there is a high probability for a woman wearing  
134 high heels to carry her bag in her left or right arm. A way to incorporate this  
135 knowledge and improve the prediction performance is to use a graphical model,  
136 *e.g.* a conditional Markov random field [23].

137 In one of the early deep learning based approaches to attribute recognition [24],  
138 the authors trained a single CNN model which considers the dependencies be-  
139 tween the attributes during training. In particular, the network shares most of its  
140 parameters among each attribute classifier and trained based on a KL-divergence  
141 based loss. Similarly, in another work, Zhu et al. [25] employ a multi-label loss  
142 function but their architecture is a multi-stream CNN model which takes mul-  
143 tiple overlapping image regions extracted from the original input person image.  
144 In [26], Wang et al. follow a different strategy and employ a recurrent CNN model  
145 which sequentially outputs the attribute predictions. Moreover, they utilize simi-  
146 lar images in the dataset during training in order to alleviate the issues regarding  
147 background clutter and uncontrolled viewing conditions.

### 148 2.3. *Joint Person Re-ID and Attribute Recognition*

149 In literature, some researchers have also explored how person attribute recog-  
150 nition and person re-ID tasks can promote each other via a joint learning strategy.  
151 The idea dates back to [27], where the authors use extracted attributes as addi-  
152 tional, mid-level semantic features to improve re-ID performance. The approach  
153 is based on training SVM-based attribute classifiers and then applying a greedy  
154 strategy to decide the optimum weights of the attributes for re-ID. The follow-up  
155 studies, however, approach the problem from a multi-task learning perspective that  
156 leverage attribute and identity information to train a single unified model [28, 29].

157 As for the examples of deep learning based approaches, in [30], Su et al. pro-  
158 posed a semi-supervised strategy, which involves training attribute classifiers on  
159 an attribute dataset and exploiting them to extend the annotations of a person re-  
160 ID dataset with person-specific attributes. Then, the person re-ID model is trained  
161 with a triplet loss defined on top of these attributes, assuming that predicted at-  
162 tribute labels should be similar for the same person. In [31], Lin et al. propose a  
163 multi-task learning framework which includes a shared CNN-based encoder and  
164 two task-specific branches, one for attribute prediction and another for re-ID. The  
165 re-weighted attribute predictions are concatenated to CNN features to incorporate  
166 semantic knowledge into person re-ID. In [32], Sun et al. present a deep person  
167 re-ID model which incorporates body parts and pose information as well as a sec-  
168 ondary attribute classification to improve the discriminative power of the learned  
169 features. Liu et al. [33] employ connectionist temporal classification (CTC) loss  
170 and self-attention to jointly learn attribute recognition and re-ID. Tay et al. [34]  
171 propose a unified architecture that combines attribute features and attribute atten-  
172 tion maps with identity and body part classification. In another recent study, Wang



173 et al. [35] suggest to learn and use hidden attributes other than the provided ones  
174 in an unsupervised manner to boost the re-identification performance.

#### 175 *2.4. Synthetic Data for Person Re-ID*

176 Several researchers have recently explored the idea of using synthetic data to  
177 address various issues in person re-ID. In [36], Barbosa et al. used MakeHu-  
178 man, a 3D character creation software, to generate 25 male and 25 female bodies  
179 which have 8 different outfit types to obtain a synthetic dataset, which they refer  
180 to SOMAset. The authors show that this dataset can be used to alleviate a main  
181 drawback of real-world datasets that they heavily rely on appearances of clothes,  
182 but not much to structural aspects of the human body. In another study [37], Sun  
183 et al. focused on issues related to viewpoint changes in re-ID datasets and devel-  
184 oped PersonX, a data generation framework which renders hand-crafted clothed  
185 human meshes onto several backgrounds in different lighting levels to better un-  
186 derstand the role of viewpoint. Lastly, in [38], Bak et al. aimed to address the  
187 lack of illumination variances in real re-ID datasets. They rendered 100 different  
188 virtual humans in multiple HDR environment maps to simulate different lighting  
189 conditions to create the SyRI dataset. They also proposed a domain adaptation  
190 technique which makes use of this synthetic data. Xiang et al. [39] proposed an-  
191 other synthetic dataset called GPR for person re-identification which consists of  
192 754 identities and around 440K bounding boxes by using the GTA5 computer  
193 game. The identities span a diverse set of persons with different gender, appear-  
194 ance, nationalities, etc. Then, they suggested a domain adaptation technique for  
195 unsupervised person re-identification that depends on these synthetically gener-  
196 ated person images. Concurrent to our work, Wang et al. [40] proposed to use  
197 Unity3D engine to generate a large-scale synthetic dataset called RandPerson that

198 is composed of images from 8K different virtual persons with different races and  
199 attributes. In particular, they automatize the clothing of these synthetic persons  
200 by generating and using a large number of random UV texture maps. In another  
201 recent study, Zeng et al. [41] addressed changes in the illumination conditions  
202 by constructing two synthetic datasets containing images with a wide range of il-  
203 lumination variations. These simulated images, however, were not generated by  
204 rendering synthetic images but obtained by applying random gamma adjustments  
205 to real images.

206 Our proposed Synthetic18K dataset, while sharing some features of earlier  
207 works, departs from them in the manner that it uses a framework which procedu-  
208 rally generates synthetic persons; resulting in a significantly higher count of unique  
209 persons (18K), covering a much more diverse range of looks in terms of body  
210 types, clothes, accessories, skin tones and facial features. And since the persons  
211 are synthetically generated, we can also provide semantic annotations about them,  
212 as well, which can be used to create similarity metrics in the given set of persons.  
213 Moreover, while the aforementioned datasets contribute mostly simple illumina-  
214 tion changes, Synthetic18K features images of each person at different environ-  
215 ments from numerous viewpoints in varying real-life -like environment conditions  
216 using a completely procedural atmosphere and weather rendering system.

### 217 **3. Synthetic18K Dataset**

218 The Synthetic18K dataset is a collection of 1,408,600 synthetically generated  
219 images of 18,306 unique persons in varying environmental conditions. The dataset  
220 is built for person re-ID and attribute recognition purposes; hence with each im-  
221 age, several annotations are also provided (Table 1). The dataset was generated

222 by re-purposing our procedural generation framework, which is built on Unity  
 223 graphics engine, specifically for the tasks at hand. Generation process is com-  
 224 pletely automatized, i.e., it does not need any supervision. Generating around  
 225 60 images per second, the whole process took about 8 hours on a system with  
 226 mid-range specifications (i7-6820HK, NVidia GTX-1070, 16GB DDR4, SSD).

### 227 3.1. Synthetic Persons

228 The synthetic persons in the dataset were procedurally generated at run-time  
 229 by making use of several content creation layers which consist of predefined set  
 230 of categorizable, annotatable randomizations as well as procedural, low-level ran-  
 231 domizations in order to yield a distinct look in each generated person (Fig. 1).

Table 1: A total of 84 ID-level attributes were used in annotating the Synthetic18K images.

Type	#	Type	#
Body Type	18*	Shoe Color	8
Color of Upper-Body Clothing	10	Hair Color	4
Length of Upper-Body Clothing Sleeve	3	Hair Type	6
Type of Lower-Body Clothing	4	Beard Type	2
Color of Lower-Body Clothing	10	Carrying Bag	Binary
Has Outerwear	Binary	Bag Color	6
Color of Outerwear	11		

\*including gender information

Table 2: The numbers showing the attainable variations of facial, clothing and accessory items that can be used in the procedural generation of synthetic persons are given below. Extended variations by color changes are additionally provided inside parentheses.

Facial Items			Clothing and Accessory Items		
Item	Male	Female	Item	Male	Female
Hair	4 (48)	3 (32)	Upper-Body Clothing	7 (28)	7 (28)
Eyebrows	2 (24)	2 (24)	Lower-Body Clothing	6 (240)	13 (520)
Beard	8 (96)	- / -	Outerwear	2 (80)	3 (120)
			Shoes	5 (40)	10 (80)
			Bags	3 (12)	3 (12)
			Other	2 (4)	3 (18)



Figure 1: An arbitrarily chosen sample of 24 synthetic persons from the Synthetic18K dataset indicating the distinct array of looks prevalent throughout the dataset.

232 Each person in the dataset has a unique body shape which can be categorized  
233 into one of 9 pre-defined major body types per gender. Uniqueness of a body  
234 shape is realized by applying a rather small white noise with uniform distribution  
235 to pre-defined body blend shapes. Facial attributes of the persons are also affected  
236 by these randomizations. The clothing and hair attributes for the persons are gen-  
237 erated from a set of several content sets and can be colored at run-time. A shared  
238 color system ensures a wide variety of distinct looking persons with randomized  
239 colors for their clothing, skin and hair which are then categorized into main groups  
240 of colors accordingly for annotation (Table 2).

241 An arbitrarily chosen sample of 24 generated synthetic persons from the Syn-  
242 thetic18K dataset (Fig. 1) demonstrates that the generated persons are easily dis-  
243 tinguishable from one another. Figure 2b presents a comparison of the images of  
244 three different synthetic persons from the Synthetic18K dataset to the ones from  
245 the real person re-ID datasets Market1501 [42] and DukeMTMC-reID [1].

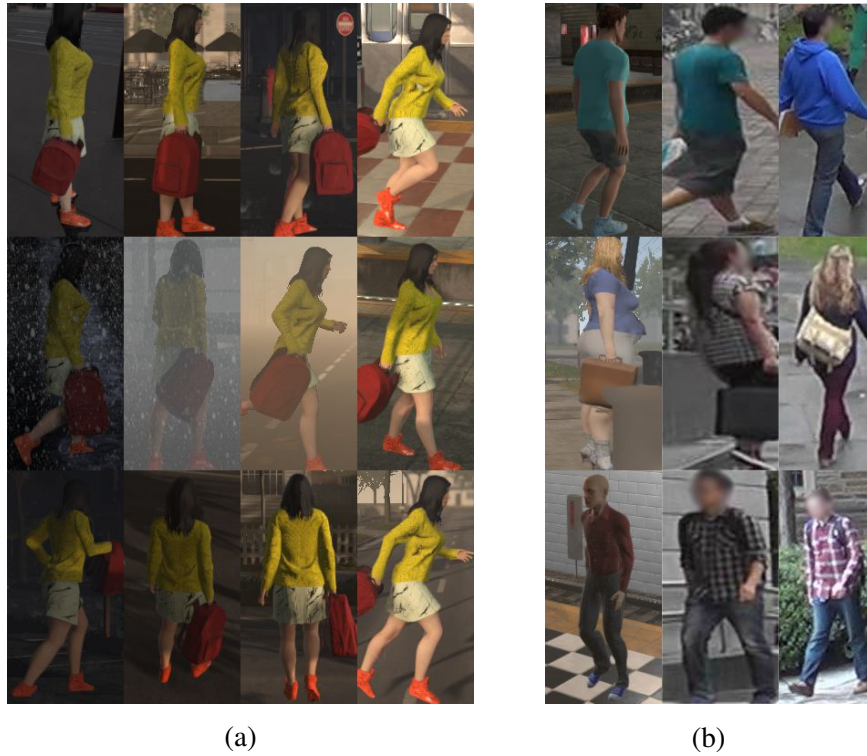


Figure 2: (a) Sample images of a synthetic person captured from various viewpoints at different locations, time of day and weather conditions. (b) A comparison of the images of three different synthetic persons from the Synthetic18K dataset (left) to the ones from the real person re-ID datasets Market1501 [42] (middle) and DukeMTMC-reID [1] (right). Real person faces are blurred for privacy concerns.

246 *3.2. Environments*

247 The Synthetic18K dataset contains images that are captured from different 3D  
 248 environments, of which three are outdoors (a town square, a suburban street and a  
 249 metropolitan urban district) and one is indoors (a subway station) (Fig. 3a). In ad-  
 250 dition, Synthetic18K also contains images that use HDR cubemaps captured from  
 251 real-life as background environments (Fig. 3b). These make up approximately  
 252 24% of the entire dataset. However, as the cubemaps are static, these images do

253 not convey background variations in terms of illumination and weather, which are  
254 present in the images captured inside the 3D environments.

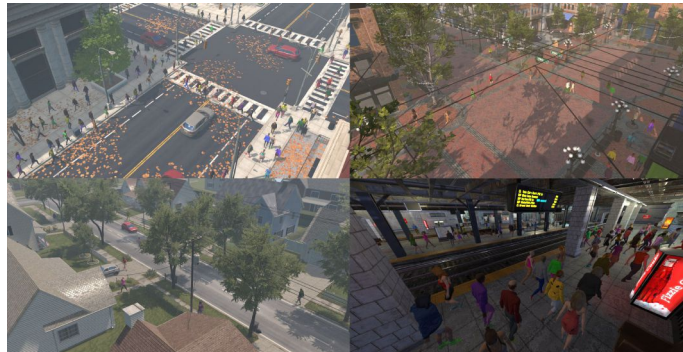
255 Each synthetic person has images that are taken at 3 different locations in each  
256 scene and cubemap (Fig. 2a). The locations for the scenes are chosen randomly  
257 from a set of pre-determined points distributed throughout each scene. At each  
258 location, the person’s images are captured for each of the time-of-day and weather  
259 variations that the framework can simulate (Fig. 3c).

## 260 4. Approach

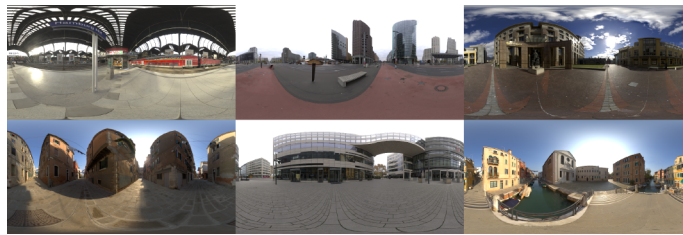
261 In our work, we consider the following three different tasks to demonstrate  
262 the importance of our pretraining strategy with synthetic data: (1) person re-ID,  
263 (2) attribute recognition, and (3) joint person re-ID and attribute recognition. The  
264 general overview of our approach is given in Fig. 4. In particular, for each afore-  
265 mentioned task, we utilize the same backbone network for feature extraction, and  
266 add additional modules to address the specifics of the task. In our experiments,  
267 we firstly pretrain the related model parameters on Synthetic18K at first and then  
268 fine-tune them on relevant real-world datasets. To validate the effectiveness of our  
269 approach, we also evaluate against the widely used strategy of using ImageNet-  
270 pretrained models. In the following, we give formal definitions of the tasks and  
271 describe the our model architectures together with some training and implemen-  
272 tation details.

### 273 4.1. Person Re-ID

274 For person re-ID task, the training data consists of a set of person images  
275  $\mathcal{S} = \{(x_1, y_1), (x_2, y_2), \dots, (x_N, y_N)\}$  where  $N$  is the total number of images, and  
276  $x_i$  and  $y_i \in [1, K]$  refer to  $i$ -th person image and its identity label , respectively,



(a)



(b)



(c)

Figure 3: Illustrating the diversity of the environments used to generate the Synthetic18K dataset. (a) Sample images of the 3D environments clockwise from top-left: a metropolitan urban district, a town square, a subway station and a suburban street. (b) Sample HDR cubemaps captured from real-world [43]. (c) Simulation of different times of day and weather conditions at the same environment setting.

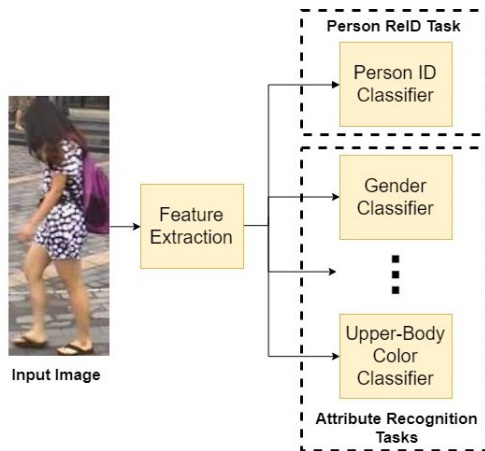


Figure 4: General overview of the proposed framework. In our work, we address person re-ID and attribute recognition by using a common backbone network as the feature extractor and extend this architecture according to the requirements of the task.

277 with  $K$  indicating the number of identities. In our analysis, we use two different  
 278 network architectures for person re-ID. The first one is a simple network consist-  
 279 ing of a basic backbone network and a classifier module. The second one, which  
 280 we refer to as HUCVReidNet, is a network that employs attention mechanisms to  
 281 better exploit contextual information and learn more discriminative features. In  
 282 the test phase, each query and gallery image are represented in terms of feature  
 283 responses at the last fully-connected layer, and we use the Euclidean distance as  
 284 the metric to retrieve the nearest neighbors.

285 **Basic Re-ID Network.** This network consists of a backbone network and a clas-  
 286 sifier module attached to the end. Features extracted from the backbone network  
 287 are first passed through a global average pooling layer, and then passed to clas-  
 288 sifier module to determine the identity. Classifier module consists of two fully  
 289 connected (FC) layers, where the second fully connected layer has units as much  
 290 as the number of identities in the training set. Any standard commonly-used



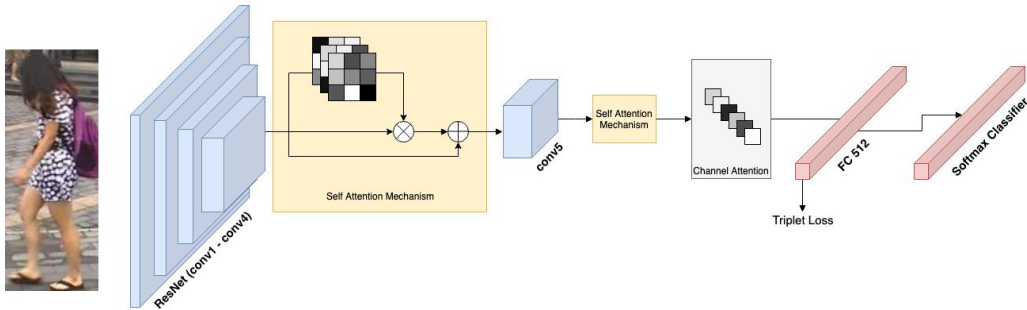


Figure 5: Our proposed HUCVReidNet model for person re-identification.

291 CNN architecture can be adapted as the backbone network. In our work, we con-  
 292 ducted experiments with three different architectures (Section 5.2) and selected  
 293 DenseNet-121 [44] as it gave the best performance.

294 **HUCVReidNet.** Our HUCVReidNet has a novel but simple architecture which  
 295 employs attentional mechanisms, as shown in Fig. 5. The model uses ResNet-  
 296 50 [4] as its backbone network. After the last convolutional layer of ResNet-50  
 297 (conv5 layer), we apply self attention [45] mechanisms both before and after this  
 298 layer to refine the feature maps extracted by ResNet-50. Moreover, we apply  
 299 a channel attention right before the global average pooling layer. Since pool-  
 300 ing takes average along spatial dimensions for each channel separately, applying  
 301 channel attention scales the channel features according to their importance before  
 302 vectorizing them. These features are then passed to FC block which has linear  
 303 layer with 1024 units, batch normalization layer and ReLU activation function.  
 304 After this FC block, a final classification layer whose dimension is equal to the  
 305 number of identities in the training set.

306 **Loss function.** We train the aforementioned models by using a composite loss

307 function that includes an identification loss and a triplet ranking loss functions:

$$\mathcal{L}_{re-ID} = \mathcal{L}_{id} + \mathcal{L}_{tri} \quad (1)$$

308 Here, the first term  $\mathcal{L}_{id}$  treats re-ID as a multi-class classification problem with  
 309 each person identity representing a distinct class label. In particular, we use cross  
 310 entropy loss on the final fully connected layer of the classifier module to enforce  
 311 identity consistency, given as below:

$$\mathcal{L}_{id} = \mathbb{E} [-\log p(y_i|x_i)] \quad (2)$$

312 where  $p(y_i|x_i)$  denotes the predicted probability of  $x_i$  belonging to the identity  
 313 label  $y_i$  based on its extracted deep features. In our implementation, we also apply  
 314 label smoothing to regularize the trained classifier by adding small constants to  
 315 ground truth values instead of using 1 and 0s.

316 The second term in Eqn. (1) is the triplet loss  $\mathcal{L}_{tri}$  which casts person re-ID  
 317 task as a metric learning problem. Specifically, we represent each person image  
 318 with the deep features extracted by the first fully connected layer of the classi-  
 319 fier module. During training, we consider a triplet  $(x_i, x_j, x_k)$  consisting of two  
 320 distinct images of the same person  $x_i$  and  $x_j$  and an image of a different person  
 321  $x_k$ . Triplet loss defined as follows enforces the model to learn a feature space in  
 322 which images of the same person are mapped closer to each other where images  
 323 of different persons are separated from each other by a large margin:

$$\mathcal{L}_{tri} = \mathbb{E} [[d(f(x_i), f(x_j)) + m - d(f(x_i), f(x_k))]_{+}] \quad (3)$$

324 where  $[z]_{+} = \max(0, z)$ ,  $m$  is a scalar representing the margin,  $f(x_i)$  is the deep  
 325 feature representation of image  $x_i$ , and  $d$  denotes the Euclidean distance. In our

326 implementation, to increase the robustness of the learned feature space, we use  
 327 hard positive and hard negative mining proposed by Hermans et al. [9].

328 **Training details.** Our person re-ID networks take 3-channel RGB images of per-  
 329 sons that are resized to  $128 \times 256$  pixels. We train our basic re-ID and HUCVRei-  
 330 dNet models for 80 epochs, by using Adam with a batch size of 32. For both of  
 331 these models, we set the initial learning rate to 0.001 for the newly added layers  
 332 and 0.0001 for the layers of the backbone network. We do not use any data aug-  
 333 mentation during training on our Synthetic18K dataset. Translation and horizontal  
 334 flip are applied randomly during fine-tuning of real-world datasets.

#### 335 4.2. Attribute Recognition

336 The training data for attribute recognition task contains a set of pairs  $\mathcal{S} =$   
 337  $\{(x_i, \mathbf{a}_i)\}$  where each pair consists of a person image  $x$  and a set of attributes  
 338  $\mathbf{a} = (a^1, a^2, \dots, a^M)$ . In our work, we use a multi-task learning approach for  
 339 attribute recognition, where classification of each attribute is considered as a sep-  
 340 arate classification task. Like our basic architecture in person re-ID, we use  
 341 DenseNet-121 as our backbone network and define separate classification mod-  
 342 ules for each attribute. Similarly, each classifier module consist of two FC layers  
 343 and a classification layer whose dimension of the final is equal to the number of  
 344 different labels for that attribute. As the backbone network is shared between all  
 345 classifier modules, it learns to extract features that is useful for all of the attributes.

346 **Loss function.** To train our network, we use a weighted cross entropy loss for  
 347 each attribute classifier module, which results in the following joint loss function:

348

$$\mathcal{L}_{attr} = \mathbb{E} \left[ - \sum_{j=1}^M \lambda_j \log p(a_i^j | x_i) \right] \quad (4)$$

349 Here,  $\lambda_j$  is a scalar denoting the importance of  $j$ -th attribute and  $p(a_i^j|x_i)$  denotes  
350 the predicted probability of  $x_i$  having the attribute  $a_i^j$  based on the extracted deep  
351 features. In our implementation, we set  $\lambda_j$  in accordance with the total number of  
352 class samples to avoid class imbalance problem.

353 **Training details.** Our attribute recognition network use  $128 \times 256$  pixels RGB  
354 images of persons as input. We train our model for 80 epochs by using Adam  
355 algorithm with batches of size 32. We set the initial learning rate to 0.001 for the  
356 newly added layers and 0.0001 for the layers of the backbone network. No data  
357 augmentation is used during training on synthetic person images. For real data,  
358 we apply a similar data augmentation scheme we used as in person re-ID.

### 359 4.3. Joint Person Re-ID and Attribute Recognition

360 In this task, we jointly train for person re-ID and attribute recognition tasks.  
361 The intuition is to utilize a shared backbone network for these tasks where the  
362 training of this network involves supervision signals from both person id and at-  
363 tribute labels. Since person re-ID and attribute recognition are two closely related  
364 tasks, we expect that a joint training scheme would result in better performances  
365 for both of these two tasks. The overall system architecture can be seen in Fig. 6.  
366 It is similar to the attribute recognition architecture but with an additional clas-  
367 sifier module for person re-ID, containing  $N + 1$  separate classifier modules,  $N$   
368 modules for classifying  $N$  distinct attributes and one for the identification. These  
369 classifier modules are same as the ones that are used in basic person re-ID and  
370 attribute recognition networks.

371 **Loss function.** Person re-ID classifier is trained with both cross entropy loss and  
372 triplet ranking loss, and attribute classification classifiers are trained with only  
373 cross entropy losses. The common backbone network is trained via supervision

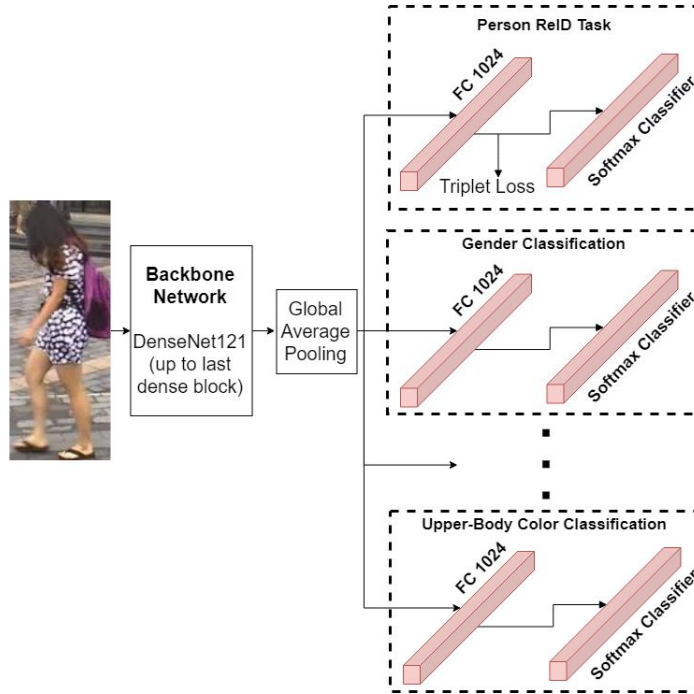


Figure 6: Person Re-ID and Attribute Recognition Multi-Task Network

374 signals from the combination of these losses as defined below:

$$\mathcal{L}_{joint} = \mathcal{L}_{attr} + \beta \mathcal{L}_{re-ID} \quad (5)$$

375 where  $\beta$  is a weight factor. In our experiments, we have observed that setting  $\beta$  to  
 376 2 gives a good trade-off between to the two tasks.

377 **Training details.** Training strategy and hyperparameters for our joint network for  
 378 the person-reid and attribute recognition tasks is completely same with those for  
 379 attribute recognition task.

#### 380 4.4. Pretraining for Feature Learning

381 The proposed Synthetic18K dataset differs from the existing synthetic datasets  
 382 proposed for person re-identification in certain aspects as mentioned before. It

383 contains images of larger number of synthetic identities captured under various  
384 illumination conditions and backgrounds. But, more importantly, our procedural  
385 generation framework allows us to play with the attributes of the generated per-  
386 sons as well. We used this capability to collect a large-scale dataset which can  
387 be used for both person re-identification and attribute recognition tasks. Our Syn-  
388 thetic18K dataset can be used for learning feature representations robust for these  
389 two tasks. As mentioned in the previous subsections, this is achieved by carefully  
390 designing pretraining strategies for the proposed deep models.

391 In our work, we considered three different pretraining schemes. While our  
392 first pretraining strategy involves only the person re-identification task, our sec-  
393 ond strategy considers the attribution recognition task. For all these settings, we  
394 develop simple neural deep models. Finally, our third pretraining scheme em-  
395 ploys a combination of these two and involves a multi-task learning setting for  
396 joint person re-identification and attribution recognition. In that respect, in our  
397 third strategy, we combine our proposed deep models by taking into account their  
398 common backbone network architecture and introduce separate heads for each one  
399 of these tasks. For all of our pretraining strategies, we follow a similar training  
400 scheme. That is, as the first step, we train our proposed model either by using the  
401 individual tasks or by utilizing both in a multi-task learning setting. This initial  
402 training step lets the deep models learn distinctive features or filter weights for  
403 the task(s) under consideration. We then use these weights to initialize the model  
404 parameters for the experiments done on the real datasets and perform finetuning  
405 the actual real data. While doing so, we set the learning rate to ...

## 406 **5. Experimental Results**

407 In the following, we first summarize the evaluation metrics used in our exper-  
408 iments. We then provide an analysis on the test set of the Synthetic18K dataset  
409 for our deep models for person re-ID and attribute recognition. Next, we test the  
410 performances of these pretrained models on several real-life datasets. Finally, we  
411 compare our results with the state-of-the-art models proposed for person re-ID,  
412 attribute recognition as well as jointly trained ones.

### 413 *5.1. Evaluation Metrics*

414 To evaluate performance on person re-ID task, we use cumulative matching  
415 characteristics (CMC) and mean average precision (mAP). To compute these,  
416 gallery images are sorted by their similarity to the query image for each query.  
417 CMC curve represents the expectation to include true person identity in the first  $k$   
418 images of the sorted gallery images. mAP is the mean value of the precision  
419 scores for all queries, where average precision for a single query is the area under  
420 the precision-recall curve. Since CMC curve considers only the first match of the  
421 true identity within  $k$  images, mAP metric is also used for person re-ID, which re-  
422 wards retrieving multiple true identities. For attribute recognition task, we report  
423 classification accuracy for each attribute as well as their averages (mA).

### 424 *5.2. Validation on Synthetic Images*

425 We first analyze how does the choice of backbone network and loss functions  
426 affect the performances of our basic re-ID model. We split our Synthetic18K  
427 dataset into training and test sets according to person identities, each contain-  
428 ing 12K and 6K different persons, respectively. Table 3 shows the results of our

Table 3: Person Re-ID Performance on Synthetic 18K.

Backbone	$\mathcal{L}_{id}$				$\mathcal{L}_{id} + \mathcal{L}_{tri}$			
	mAP	Rank-1	Rank-5	Rank-10	mAP	Rank-1	Rank-5	Rank-10
DenseNet-121	98.9	99.2	100.0	100.0	99.4	99.9	100.0	100.0
ResNet-50	96.6	98.2	99.7	100.0	97.4	99.1	100.0	100.0
MobileNetV2	94.5	96.7	98.9	99.6	96.1	97.1	99.0	99.7

429 analysis. In particular, we consider DenseNet-121 [44], ResNet-50 [4], and Mo-  
 430 bileNetV2 [46] models as our backbone network and train them by using solely  
 431 the identity loss ( $\mathcal{L}_{id}$ ) and the joint loss function containing both the identity and  
 432 the triplet loss ( $\mathcal{L}_{id} + \mathcal{L}_{tri}$ ). We observe that our model with DenseNet-121 as  
 433 its backbone achieves slightly higher scores than the other two models. More-  
 434 over, training the models with the joint loss function improves the performances.  
 435 Hence, as we mentioned before, for the rest of the experiments we use DenseNet-  
 436 121 model in our basic network models for person re-ID and attribute recognition.

437 In Table 4, we provide the performance of our attribute recognition model on  
 438 our Synthetic18K dataset. Synthetic person images introduce certain challenges  
 439 as compared to the aforementioned analysis regarding person re-ID that the in-  
 440 dividual attribute prediction scores are not very high, especially predicting color  
 441 attributes seems more difficult. This demonstrates that Synthetic18K dataset could  
 442 be used as a test-bed for attribute recognition approaches.

### 443 5.3. Using Synthetic Data for Pretraining

444 In this section, we provide several experiments regarding how pretraining on  
 445 our Synthetic18K dataset can help improving performances on real datasets, espe-



Table 4: Attribute Recognition Performance on Synthetic18K.

Gender	Age	Hair	Beard	Weight	Sleeve Len.	L.Body Clth	
99.5	94.8	97.4	95.0	93.0	82.2	97.3	
L.Body Clth Col	U.Body Clth Col	Over Clth	Hand Bag	Hand Bag Col	Hair Col	Shoe Col	mA
58.8	84.1	76.1	85.8	85.2	77.1	93.4	87.1

446 cially compared to commonly used strategy of using ImageNet pretrained models.  
 447 **Person Re-ID.** In our experiments, we use Market1501 [42] and DukeMTMC-  
 448 reID [1] datasets containing 32,668 and 34,183 real person images, respectively.  
 449 In Market1501, 751 identities are allocated for training and the rest 750 iden-  
 450 tities for testing. In DukeMTMC-reID contains 1404 identities, of which 702  
 451 identities are selected for training and the rest for testing. In our analysis, we fine-  
 452 tune our basic re-ID and HUCVReidNet models, which were pretrained on our  
 453 Synthetic18K dataset, on the training sets of these datasets and report their per-  
 454 formances on the corresponding test sets accordingly. As a comparison, we also  
 455 provide the results of our models which instead utilize ImageNet pretrained back-  
 456 bones. Table 5 reports these comparisons. As can be seen, pretraining on Syn-  
 457 thetic18K improves re-ID performances on both Market1501 and DukeMTMC-  
 458 reID datasets. For our basic re-ID model, pretraining on Synthetic18K results in  
 459 2.9 and 0.6 increase in mAP, and 1.7 and 0.8 increase on Rank-1 scores on Mar-  
 460 ket1501 and DukeMTMC-reID datasets, respectively. Again, for HUCVReidNet,  
 461 pretraining gives much better results in terms of mAP, Rank-1 and Rank-5 scores.  
 462 Moreover, with its inherent attention mechanisms, our proposed HUCVReidNet  
 463 model gives better results than our basic re-ID model.  
 464 **Attribute Recognition.** We conduct our pretraining analysis on Market1501-  
 465 Attributes [31] dataset, an extended version of Market1501 [42] where each per-

Table 5: Analysis of pretraining on Synthetic18K for person ReID.

Model	Dataset	Synthetic18K			ImageNet		
		mAP	Rank-1	Rank-5	mAP	Rank-1	Rank-5
Basic ReidNet	Market1501	77.3	91.6	96.6	74.4	89.9	96.1
	DukeMTMC-reID	63.3	81.3	90.6	62.7	80.5	89.4
HUCVReidNet	Market1501	78.8	92.4	97.9	78.3	91.9	97.1
	DukeMTMC-reID	63.7	83.0	91.9	63.7	81.7	91.0

Table 6: Analysis of pretraining on Synthetic18K for attribute recognition

Pretrain	gender	age	hair	L.slv	L.low	S.clth	B.pack	H.bag	bag	hat	C.up	C.low	mA
ImageNet	88.7	84.8	85.7	92.5	92.7	93.1	86.2	87.6	73.8	95.2	77.0	70.1	85.6
Synthetic18K	<b>91.4</b>	<b>85.6</b>	<b>86.6</b>	<b>93.5</b>	<b>93.6</b>	<b>94.0</b>	<b>87.8</b>	<b>88.1</b>	<b>76.9</b>	<b>97.5</b>	<b>78.1</b>	<b>71.0</b>	<b>87.0</b>

466 son image is annotated with 27 different attributes. We used the same training  
 467 and testing splits as in person re-ID. Table 6 shows accuracy scores for each at-  
 468 tribute as well as the average accuracy (mA). We observe that the model that has  
 469 pretrained on Synthetic18K outperforms the ImageNet pretrained model for all at-  
 470 tributes, resulting in 1.4 increase in the mean accuracy. This demonstrates that our  
 471 synthetic data pretraining approach is also effective for attribute recognition task.

472 **Joint Person Re-ID and Attribute Recognition.** We use Market1501 and Market1501-  
 473 Attributes datasets and follow a similar strategy and fine-tune our joint model pre-  
 474 trained on our Synthetic18K dataset on the training set of these datasets using both  
 475 person attributes and identities, and subsequently evaluate it on the corresponding  
 476 test set. Table 7 reports prediction accuracies of the person attributes together with  
 477 mAP and Rank-1 for re-ID. We again observe that our joint model pretrained on

Table 7: Analysis of pretraining on Synthetic18K for Joint Re-ID and Attribution Recognition

Pretrain	gender	age	hair	L.slv	L.low	S.clth	B.pack	H.bag	bag	hat	C.up	C.low	mA	mAP	Rank-1
ImageNet	89.9	85.7	86.1	93.7	92.9	92.9	86.1	88.2	76.3	97.1	76.7	70.3	86.3	76.4	88.9
Synthetic18K	<b>91.5</b>	<b>86.2</b>	<b>88.2</b>	<b>93.9</b>	<b>94.0</b>	<b>94.0</b>	<b>88.0</b>	<b>88.8</b>	<b>78.9</b>	<b>97.2</b>	<b>78.4</b>	<b>71.4</b>	<b>87.5</b>	<b>78.4</b>	<b>90.3</b>

478 Synthetic18K outperforms the ImageNet pretrained model by a large margin. Pre-  
 479 training on Synthetic18K gives 1.2 increase in mA, 2.0 increase in mAP and 1.4  
 480 increase in Rank-1 score. Moreover, these results demonstrate that jointly train-  
 481 ing a model for person re-ID and attribute recognition model improves the model  
 482 performances for the individual tasks (cf. Table 5 and 6).

#### 483 5.4. Comparison with the state-of-the-art

484 In this section, we compare the results of our person re-ID, attribute recogni-  
 485 tion models and their joint version against those of the state-of-the-art approaches.  
 486 **Person Re-ID.** Table 8 and 9 show the comparison between our basic re-ID and  
 487 HUCVReidNet models that we pretrained on our Synthetic18K dataset and the  
 488 state-of-the-art algorithms on Market-1501 and DukeMTMC-reID datasets, re-  
 489 spectively. We observe that our models outperforms most of the recently proposed  
 490 re-ID models. There are a few methods that give relatively better results than ours  
 491 but these models have highly complex network architectures. For instance, the  
 492 models in [20, 47, 19, 48] consider part structures of the human to extract local  
 493 features from images, the method proposed in [14] employs graph neural network  
 494 instead of the commonly use Siamese networks to compare query-gallery pairs,  
 495 or the work in [49] uses self-distillation learning along with a novel code pyramid  
 496 based coarse-to-fine (CtF) hashing code search strategy. In addition to these, the  
 497 model by Chen et al. [50] utilizes text descriptions of the person images to guide

Table 8: Person Re-ID performances on Market1501.

Method	mAP	Rank-1	Rank-5	Method	mAP	Rank-1	Rank-5
OIM Loss [52]	60.9	82.1	-	HA-CNN [53]	75.7	91.2	-
SpindleNet [17]	-	76.9	91.5	Pose-transfer [54]	58.0	79.8	-
MSCAN [18]	57.5	80.3	-	MLFN [55]	74.3	90.0	-
SSM [56]	68.8	82.2	-	DML [57]	70.5	89.3	-
k-reciprocal [58]	63.6	77.1	-	Suh [47]	79.6	91.7	96.9
Point 2 Set [12]	44.3	70.7	-	SPReID [20]	83.4	93.7	97.6
CADL [59]	47.1	73.8	-	SGGNN [14]	82.8	92.3	96.1
DPFL [60]	73.1	88.9	-	GLILA [50]	81.8	93.3	-
VI+LSRO [61]	66.1	84.0	-	Mancs [51]	82.3	93.1	-
SVDNet [62]	62.1	82.3	92.3	PCB+RPP [19]	81.6	93.8	97.5
OL-MANS [63]	-	60.7	-	IID [41]	71.5	88.5	-
Pose Driven [16]	63.4	84.1	92.7	Wang et al. [40]	70.9	87.2	-
Part Aligned [64]	63.4	81.0	92.0	Xiang et al [39]	50.8	76.2	89.2
HydraPlus-Net [65]	-	76.9	91.3	Generalizing-Reid [66]	71.5	88.1	94.4
TriNet [9]	69.1	84.9	94.2	RGA-SC [67]	88.4	96.1	-
DarkRank [68]	74.3	89.8	-	ISP [48]	88.6	95.3	98.6
PN-GAN [69]	72.6	89.4	-	CtF [49]	84.9	93.7	-
DuATM [70]	76.6	91.4	97.1	Zhuang et al. [71]	77.3	91.3	-
HAP2S_E [72]	69.8	84.2	-	Ours (Basic ReidNet)	77.3	91.6	96.6
HAP2S_P [72]	69.4	84.6	-	Ours (HUCVReid)	78.8	92.4	97.9

498 learning of the local and global visual features. The model suggested by Wang et  
499 al. [51] has many attention blocks at different layers of the backbone network and  
500 during its training separate classification loss functions are attached to output of  
501 each of these blocks.

Table 9: Person Re-ID performances on DukeMTMC-reID.

Method	mAP	Rank-1	Rank-5	Method	mAP	Rank-1	Rank-5
BoW+KISSME [42]	12.2	25.1	-	SPReID [20]	73.3	86.0	93.0
LOMO+XQDA [73]	17.0	30.8	-	SGGNN [14]	68.2	81.1	88.4
APR [31]	51.9	70.7	-	Suh [47]	69.3	84.4	92.2
ACRN [74]	52.0	72.6	84.8	Mancs [51]	71.8	84.9	-
DPFL [60]	60.6	79.2	-	AANet-152 [34]	74.3	87.7	-
OIM Loss [52]	47.4	68.1	-	PCB+RPP [19]	69.2	83.3	90.5
Basel.+LSRO [61]	47.1	67.7	-	IID [41]	60.6	78.1	-
SVDNet [62]	56.8	76.7	86.4	Wang et al. [40]	60.6	79.4	-
CamStyle [75]	57.6	78.3	-	Xiang et al [39]	51.9	71.2	82.7
Pose-transfer [54]	48.1	68.6	-	Generalizing-Reid [66]	65.2	79.5	88.3
MLFN [55]	62.8	81.0	-	ISP [48]	80.0	89.6	95.5
DuATM [70]	64.6	81.8	90.2	CtF [49]	74.8	87.6	-
PN-GAN [69]	53.2	73.6	-	Zhuang et al. [71]	67.3	82.5	-
HAP2S_P [72]	60.6	75.9	-	Ours (Basic ReidNet)	63.3	81.3	90.6
HAP2S_E [72]	59.6	76.1	-	Ours (HUCVReid)	63.7	83.0	91.9

502 **Attribute Recognition.** Table 10 provides a comparison between our model  
503 against the recent attribute recognition models on the Market1501-Attributes dataset.  
504 Similar to person re-ID, our proposed networks that have been pretrained on our  
505 Synthetic18K dataset give better or comparable accuracies as compared to the  
506 state-of-the-art models. Only the average accuracies of the methods proposed  
507 in [76, 33] are a bit higher than ours, but these models achieve these results either  
508 by considering training attribute recognition jointly with person re-ID [33] or by  
509 explicitly learning the importance of each attribute on a validation set [76].

Table 10: Attribute Recognition Performances on Market1501-Attributes

Method	gender	age	hair	L.slv	L.low	S.clth	B.pack	H.bag	bag	hat	C.up	C.low	mA
ARN [31]	87.5	85.8	84.2	93.5	93.6	93.6	86.6	88.1	78.6	97.0	72.4	71.7	86.0
APR [31]	88.9	88.6	84.4	93.6	93.7	92.8	84.9	90.4	76/4	97.1	74.0	73.8	86.6
Sun et al. [32]	88.9	84.8	78.3	93.5	92.1	84.8	85.5	88.4	67.3	97.1	87.5	87.2	86.3
AWMDN [76]	-	-	-	-	-	-	-	-	-	-	-	-	88.5
MLFN [55]	-	-	-	-	-	-	-	-	-	-	-	-	85.3
PANDA [77]	-	-	-	-	-	-	-	-	-	-	-	-	86.8
JCM [33]	89.7	87.4	82.5	93.7	93.3	89.2	85.2	86.2	86.9	97.2	92.4	93.1	89.7
AANet-152 [34]	92.3	88.2	86.6	94.5	94.2	94.8	87.8	89.6	79.7	98.0	77.0	70.8	87.8
Ours (Basic)	91.4	85.6	86.6	93.5	93.6	94.0	87.8	88.1	76.9	97.5	78.1	71.0	87.0
Ours (Joint)	91.5	86.2	88.2	93.9	94.0	94.0	88.0	88.8	78.9	97.2	78.4	71.4	87.5

510 **Joint Attribute Prediction and Person Re-ID.** In Table 11, we show the results  
511 of our joint model along with those of the recent approaches which also consider  
512 joint training of a model on both person re-ID and attribute recognition tasks.  
513 We find that our model achieves much better re-ID and recognition performances  
514 than most of the state-of-the-art approaches. The model JCM-57344 in [33] out-  
515 performs our model but it achieves this by using a 57344 dimensional feature  
516 embedding. In fact, the re-ID performance of its second version with a 1024 di-  
517 mensional representation is much lower than ours. Moreover, the AANet models  
518 in [34] give a bit better predictions than ours. However, it is important to mention  
519 that AANet employs body part locations to extract local features while we only  
520 consider a global representation of person images.

## 521 6. Conclusion

522 In this work, we have introduced Synthetic18K dataset that consists of synthet-  
523 ically generated photo-realistic person images. Each image in our dataset is anno-  
524 tated with both the person identity label and the relevant person attributes. Particu-

Table 11: Joint Attribute Prediction and Person Re-ID Performances on Market1501

Method	mA	mAP	Rank-1	Method	mA	mAP	Rank-1
ACRN [74]	-	62.6	83.6	AANet-50 [34]	-	82.4	93.9
JCM-1024 [33]	-	75.7	84.9	AANet-152 [34]	87.8	83.4	93.9
JCM-57344 [33]	89.7	81.2	91.3	Wang et al. [35]	-	76.0	91.3
Sun et al. [32]	87.0	70.1	87.0	Ours (Joint)	87.5	78.4	90.3
APR [31]	86.6	66.9	87.0				

525 larly, we have addressed person re-ID and attribute recognition tasks and demon-  
 526 strated that large-scale pretraining of simple deep models on our Synthetic18K  
 527 dataset greatly improves the model performances on the real-life datasets. More-  
 528 over, we have demonstrated that joint training of a basic deep model for person  
 529 re-ID and attribute recognition on Synthetic18K outperforms the individual model  
 530 performances and gives better or comparable results than the state-of-the-art meth-  
 531 ods. As a future work, we plan to investigate the use of synthetic data to boost the  
 532 performance of video re-ID using computer generated video sequences.

### 533 Acknowledgments

534 This work was supported in part by TUBA GEBIP fellowship awarded to E.  
 535 Erdem and by TUBITAK-1001 Program Award No. 217E029.

### 536 References

- 537 [1] E. Ristani, F. Solera, R. Zou, R. Cucchiara, C. Tomasi, Performance mea-  
 538 sures and a data set for multi-target, multi-camera tracking, in: Proc. ECCV  
 539 Workshop on Benchmarking Multi-Target Tracking, 2016.

- 540 [2] L. Zheng, Y. Yang, A. G. Hauptmann, Person re-identification: Past, present  
541 and future, arXiv preprint arXiv:1610.02984 (2016).
- 542 [3] S. Karanam, M. Gou, Z. Wu, A. Rates-Borras, O. Camps, R. J. Radke, A  
543 systematic evaluation and benchmark for person re-identification: Features,  
544 metrics, and datasets, *IEEE Trans. Pattern Anal. Mach. Intell.* 41 (3) (2019)  
545 523–536.
- 546 [4] K. He, X. Zhang, S. Ren, J. Sun, Deep residual learning for image recogni-  
547 tion, in: *Proc. CVPR, 2016*, pp. 770–778.
- 548 [5] J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, L. Fei-Fei, ImageNet: A large-  
549 scale hierarchical image database, in: *Proc. CVPR, 2009*.
- 550 [6] L. Zheng, Z. Bie, Y. Sun, J. Wang, C. Su, S. Wang, Q. Tian, Mars: A video  
551 benchmark for large-scale person re-identification, in: *Proc. ECCV, 2016*.
- 552 [7] D. Chung, K. Tahboub, E. J. Delp, A two stream siamese convolutional neu-  
553 ral network for person re-identification, in: *Proc. ICCV, 2017*.
- 554 [8] S. Ding, L. Lin, G. Wang, H. Chao, Deep feature learning with relative dis-  
555 tance comparison for person re-identification, *Pattern Recognition* 48 (10)  
556 (2015) 2993–3003.
- 557 [9] A. Hermans, L. Beyer, B. Leibe, In defense of the triplet loss for person  
558 re-identification, arXiv preprint arXiv:1703.07737 (2017).
- 559 [10] W. Chen, X. Chen, J. Zhang, K. Huang, Beyond triplet loss: a deep quadru-  
560 plet network for person re-identification, in: *Proc. CVPR, 2017*.



- 561 [11] F. Wang, W. Zuo, L. Lin, D. Zhang, L. Zhang, Joint learning of single-  
562 image and cross-image representations for person re-identification, in: Proc.  
563 CVPR, 2016, pp. 1288–1296.
- 564 [12] S. Zhou, J. Wang, J. Wang, Y. Gong, N. Zheng, Point to set similarity based  
565 deep feature learning for person re-identification, in: Proc. CVPR, 2017, pp.  
566 3741–3750.
- 567 [13] W. Chen, X. Chen, J. Zhang, K. Huang, A multi-task deep network for per-  
568 son re-identification, in: Proc. AAAI, 2017, pp. 3988–3994.
- 569 [14] Y. Shen, H. Li, S. Yi, D. Chen, X. Wang, Person re-identification with deep  
570 similarity-guided graph neural network, in: Proc. ECCV, 2018, pp. 486–504.
- 571 [15] Z. Wang, R. Hu, Y. Yu, J. Jiang, C. Liang, J. Wang, Scale-adaptive low-  
572 resolution person re-identification via learning a discriminating surface, in:  
573 Proc. IJCAI, 2016, p. 26692675.
- 574 [16] C. Su, J. Li, S. Zhang, J. Xing, W. Gao, Q. Tian, Pose-driven deep convo-  
575 lutional model for person re-identification, in: Proc. ICCV, 2017, pp. 3960–  
576 3969.
- 577 [17] H. Zhao, M. Tian, S. Sun, J. Shao, J. Yan, S. Yi, X. Wang, X. Tang, Spindle  
578 net: Person re-identification with human body region guided feature decom-  
579 position and fusion, in: Proc. CVPR, 2017, pp. 1077–1085.
- 580 [18] D. Li, X. Chen, Z. Zhang, K. Huang, Learning deep context-aware features  
581 over body and latent parts for person re-identification, in: Proc. CVPR, 2017,  
582 pp. 384–393.

- 583 [19] Y. Sun, L. Zheng, Y. Yang, Q. Tian, S. Wang, Beyond part models: Person  
584 retrieval with refined part pooling (and a strong convolutional baseline), in:  
585 Proc. ECCV, 2018, pp. 480–496.
- 586 [20] M. M. Kalayeh, E. Basaran, M. Gökmen, M. E. Kamasak, M. Shah, Human  
587 semantic parsing for person re-identification, in: Proc. CVPR, 2018, pp.  
588 1062–1071.
- 589 [21] G. Sharma, F. Jurie, C. Schmid, Expanded parts model for human attribute  
590 and action recognition in still images, in: Proc. CVPR, 2013, pp. 652–659.
- 591 [22] Y. Deng, P. Luo, C. C. Loy, X. Tang, Pedestrian attribute recognition at far  
592 distance, in: Proc. ACM-MM, 2014, pp. 789–792.
- 593 [23] H. Chen, A. Gallagher, B. Girod, Describing clothing by semantic attributes,  
594 in: Proc. ECCV, 2012, pp. 609–623.
- 595 [24] P. Sudowe, H. Spitzer, B. Leibe, Person attribute recognition with a jointly-  
596 trained holistic CNN model, in: Proc. CVPRW, 2015, pp. 87–95.
- 597 [25] J. Zhu, S. Liao, Z. Lei, S. Z. Li, Multi-label convolutional neural network  
598 based pedestrian attributeclassification, Image and Vision Computing 58  
599 (2017) 224–229.
- 600 [26] J. Wang, X. Zhu, S. Gong, W. Li, Attribute recognition by joint recurrent  
601 learning of context and correlation, in: Proc. ICCV, 2017, pp. 531–540.
- 602 [27] R. Layne, T. Hospedales, S. Gong, Person re-identification by attributes, in:  
603 Proc. BMVC, 2012, pp. 24.1–24.11.

- 604 [28] C. Su, F. Yang, S. Zhang, Q. Tian, L. S. Davis, W. Gao, Multi-task learn-  
605 ing with low rank attribute embedding for person re-identification, in: Proc.  
606 ICCV, 2015, pp. 3739–3747.
- 607 [29] P. Peng, Y. Tian, T. Xiang, Y. Wang, T. Huang, Joint learning of semantic  
608 and latent attributes, in: Proc. ECCV, 2016.
- 609 [30] C. Su, S. Zhang, J. Xing, W. Gao, Q. Tian, Deep attributes driven multi-  
610 camera person re-identification, in: Proc. ECCV, 2016.
- 611 [31] Y. Lin, L. Zheng, Z. Zheng, Y. Wu, Z. Hu, C. Yan, Y. Yang, Improving person  
612 re-identification by attribute and identity learning, *Pattern Recognition* 95  
613 (2019) 151–161.
- 614 [32] C. Sun, N. Jiang, L. Zhang, Y. Wang, W. Wu, Z. Zhou, Unified framework for  
615 joint attribute classification and person re-identification, in: Proc. ICANN,  
616 2018, pp. 637–647.
- 617 [33] H. Liu, J. Wu, J. Jiang, M. Qi, R. Bo, Sequence-based person attribute  
618 recognition with joint ctc-attention model, arXiv preprint arXiv:1811.08115  
619 (2018).
- 620 [34] C.-P. Tay, S. Roy, K.-H. Yap, Aanet: Attribute attention network for person  
621 re-identifications, in: Proc. CVPR, 2019, pp. 7134–7143.
- 622 [35] Z. Wang, J. Jiang, Y. Wu, M. Ye, X. Bai, S. Satoh, Learning sparse and  
623 identity-preserved hidden attributes for person re-identification, *IEEE Trans.*  
624 *Image Process.* 29 (2020) 2013–2025.

- 625 [36] I. B. Barbosa, M. Cristani, B. Caputo, A. Rognhaugen, T. Theoharis, Look-  
626 ing beyond appearances: Synthetic training data for deep CNNs in re-  
627 identification, *Comput. Vis. Image Und.* 167 (2018) 50–62.
- 628 [37] X. Sun, L. Zheng, Dissecting person re-identification from the viewpoint of  
629 viewpoint, in: *Proc. CVPR*, 2019.
- 630 [38] S. Bak, P. Carr, J.-F. Lalonde, Domain adaptation through synthesis for un-  
631 supervised person re-identification, in: *Proceedings of the European Confer-*  
632 *ence on Computer Vision (ECCV)*, 2018, pp. 189–205.
- 633 [39] S. Xiang, Y. Fu, G. You, T. Liu, Unsupervised domain adaptation through  
634 synthesis for person re-identification, in: *Proc. ICME*, 2020.
- 635 [40] Y. Wang, S. Liao, L. Shao, Surpassing real-world source training data: Ran-  
636 dom 3d characters for generalizable person re-identification, in: *Proc. ACM-*  
637 *MM*, 2020, p. 34223430.
- 638 [41] Z. Zeng, Z. Wang, Z. Wang, Y. Zheng, Y.-Y. Chuang, S. Satoh, Illumination-  
639 adaptive person re-identification, *IEEE Trans. Multimedia* 22 (12) (2020)  
640 3064–3074.
- 641 [42] L. Zheng, L. Shen, L. Tian, S. Wang, J. Wang, Q. Tian, Scalable person  
642 re-identification: A benchmark, in: *Proc. ICCV*, 2015, pp. 1116–1124.
- 643 [43] G. Zaal, HDRI Haven, <https://hdrihaven.com/hdriis>, online; ac-  
644 cessed: 2019-02-20 (2019).
- 645 [44] G. Huang, Z. Liu, L. Van Der Maaten, K. Q. Weinberger, Densely connected  
646 convolutional networks, in: *Proc. CVPR*, 2017, pp. 4700–4708.

- 647 [45] H. Zhang, I. Goodfellow, D. Metaxas, A. Odena, Self-attention generative  
648 adversarial networks, arXiv preprint arXiv:1805.08318 (2018).
- 649 [46] M. Sandler, A. Howard, M. Zhu, A. Zhmoginov, L.-C. Chen, Mobilenetv2:  
650 Inverted residuals and linear bottlenecks, in: Proc. CVPR, 2018, pp. 4510–  
651 4520.
- 652 [47] Y. Suh, J. Wang, S. Tang, T. Mei, K. Mu Lee, Part-aligned bilinear represen-  
653 tations for person re-identification, in: Proc. ECCV, 2018, pp. 402–419.
- 654 [48] K. Zhu, H. Guo, Z. Liu, M. Tang, J. Wang, Identity-guided human semantic  
655 parsing for person re-identification, in: Proc. ECCV, 2020.
- 656 [49] G. Wang, S. Gong, J. Cheng, Z. Hou, Faster person re-identification, in:  
657 Proc. ECCV, 2020.
- 658 [50] D. Chen, H. Li, X. Liu, Y. Shen, J. Shao, Z. Yuan, X. Wang, Improving  
659 deep visual representation for person re-identification by global and local  
660 image-language association, in: Proc. ECCV, 2018, pp. 54–70.
- 661 [51] C. Wang, Q. Zhang, C. Huang, W. Liu, X. Wang, Mancs: A multi-task at-  
662 tentional network with curriculum sampling for person re-identification, in:  
663 Proc. ECCV, 2018, pp. 365–381.
- 664 [52] T. Xiao, S. Li, B. Wang, L. Lin, X. Wang, Joint detection and identification  
665 feature learning for person search, in: Proc. CVPR, 2017, pp. 3415–3424.
- 666 [53] W. Li, X. Zhu, S. Gong, Harmonious attention network for person re-  
667 identification, in: Proc. CVPR, 2018, pp. 2285–2294.

- 668 [54] J. Liu, B. Ni, Y. Yan, P. Zhou, S. Cheng, J. Hu, Pose transferrable person  
669 re-identification, in: Proc. CVPR, 2018, pp. 4099–4108.
- 670 [55] X. Chang, T. M. Hospedales, T. Xiang, Multi-level factorisation net for per-  
671 son re-identification, in: Proc. CVPR, 2018, pp. 2109–2118.
- 672 [56] S. Bai, X. Bai, Q. Tian, Scalable person re-identification on supervised  
673 smoothed manifold, in: Proc. CVPR, 2017, pp. 2530–2539.
- 674 [57] Y. Zhang, T. Xiang, T. M. Hospedales, H. Lu, Deep mutual learning, in:  
675 Proc. CVPR, 2018, pp. 4320–4328.
- 676 [58] Z. Zhong, L. Zheng, D. Cao, S. Li, Re-ranking person re-identification with  
677 k-reciprocal encoding, in: Proc. CVPR, 2017, pp. 1318–1327.
- 678 [59] J. Lin, L. Ren, J. Lu, J. Feng, J. Zhou, Consistent-aware deep learning for  
679 person re-identification in a camera network, in: Proc. CVPR, 2017, pp.  
680 5771–5780.
- 681 [60] Y. Chen, X. Zhu, S. Gong, Person re-identification by deep learning multi-  
682 scale representations, in: Proc. ICCV, 2017, pp. 2590–2600.
- 683 [61] Z. Zheng, L. Zheng, Y. Yang, Unlabeled samples generated by gan improve  
684 the person re-identification baseline in vitro, in: Proc. ICCV, 2017, pp. 3754–  
685 3762.
- 686 [62] Y. Sun, L. Zheng, W. Deng, S. Wang, Svdnet for pedestrian retrieval, in:  
687 Proc. ICCV, 2017, pp. 3800–3808.

- 688 [63] J. Zhou, P. Yu, W. Tang, Y. Wu, Efficient online local metric adaptation  
689 via negative samples for person re-identification, in: Proc. ICCV, 2017, pp.  
690 2420–2428.
- 691 [64] L. Zhao, X. Li, Y. Zhuang, J. Wang, Deeply-learned part-aligned represen-  
692 tations for person re-identification, in: Proc. ICCV, 2017, pp. 3219–3228.
- 693 [65] X. Liu, H. Zhao, M. Tian, L. Sheng, J. Shao, S. Yi, J. Yan, X. Wang,  
694 Hydraplus-net: Attentive deep features for pedestrian analysis, in: Proc.  
695 ICCV, 2017, pp. 350–359.
- 696 [66] C. Luo, C. Song, Z. Zhang, Generalizing person re-identification by camera-  
697 aware invariance learning and cross-domain mixup, in: Proc. ECCV, 2020.
- 698 [67] Z. Zhang, C. Lan, W. Zeng, X. Jin, Z. Chen, Relation-aware global attention  
699 for person re-identification, in: Proc. CVPR, 2020.
- 700 [68] Y. Chen, N. Wang, Z. Zhang, Darkrank: Accelerating deep metric learning  
701 via cross sample similarities transfer, in: Proc. AAAI, 2018.
- 702 [69] X. Qian, Y. Fu, T. Xiang, W. Wang, J. Qiu, Y. Wu, Y.-G. Jiang, X. Xue, Pose-  
703 normalized image generation for person re-identification, in: Proc. ECCV,  
704 2018, pp. 650–667.
- 705 [70] J. Si, H. Zhang, C.-G. Li, J. Kuen, X. Kong, A. C. Kot, G. Wang, Dual  
706 attention matching network for context-aware feature sequence based person  
707 re-identification, in: Proc. CVPR, 2018, pp. 5363–5372.
- 708 [71] Z. Zhuang, L. Wei, L. Xie, T. Zhang, H. Zhang, H. Wu, H. Ai, Q. Tian, Re-

- 709 thinking the distribution gap of person re-identification with camera-based  
710 batch normalization, in: Proc. ECCV, 2020.
- 711 [72] R. Yu, Z. Dou, S. Bai, Z. Zhang, Y. Xu, X. Bai, Hard-aware point-to-set deep  
712 metric for person re-identification, in: Proc. ECCV, 2018, pp. 188–204.
- 713 [73] S. Liao, Y. Hu, X. Zhu, S. Z. Li, Person re-identification by local maximal  
714 occurrence representation and metric learning, in: Proc. CVPR, 2015, pp.  
715 2197–2206.
- 716 [74] A. Schumann, R. Stiefelwagen, Person re-identification by deep learning  
717 attribute-complementary information, in: Proc. CVPRW, 2017, pp. 20–28.
- 718 [75] Z. Zhong, L. Zheng, Z. Zheng, S. Li, Y. Yang, Camstyle: A novel data aug-  
719 mentation method for person re-identification, IEEE Trans. Image Process.  
720 28 (3) (2018) 1176–1190.
- 721 [76] K. He, Z. Wang, Y. Fu, R. Feng, Y.-G. Jiang, X. Xue, Adaptively weighted  
722 multi-task deep network for person attribute classification, in: Proc. ACM-  
723 MM, 2017, pp. 1636–1644.
- 724 [77] N. Zhang, M. Paluri, M. Ranzato, T. Darrell, L. Bourdev, Panda: Pose  
725 aligned networks for deep attribute modeling, in: Proc. CVPR, 2014, pp.  
726 1637–1644.