# *Supplementary Materials*:
# Person Search by Text Attribute Query as Zero-Shot Learning

Qi Dong
Queen Mary University of London
q.dong@qmul.ac.uk

Shaogang Gong
Queen Mary University of London
s.gong@qmul.ac.uk

Xiatian Zhu
Vision Semantics Ltd.
eddy.zhuxt@gmail.com

As described in the paper, the proposed *Attribute-Image Hierarchical Matching* (AIHM) model consists of 4 components: (1) Hierarchical visual embedding (Sec. 1), (2) Hierarchical textual embedding (Sec. 2), (3) Cross-modality cross-level embedding (Sec. 3), and (4) Matching module (Sec. 4). We detail the network design of all the components below. The notations are consistent with the main paper. For facilitating the presentation, we start by giving the embedding dimensions as summarised in Table 1.

Table 1: Embedding dimensions.

| Definition | Notation | Value |
|---|---|---|
| Local embedding dimension | $\text{Dim}_{emb}^{\text{loc}}$ | 512 |
| Global embedding dimension | $\text{Dim}_{emb}^{\text{glo}}$ | 1024 |
| Cross-modal embedding dimension | $\text{Dim}_{emb}^{\text{s}}$ | 512 |

## 1. Hypercritical Visual Embedding Network

We describe the details of the 2-layers and 4-layers hierarchical visual embedding.

### 1.1. 2-Layers Hierarchical Visual Embedding

In the main experiments, we employ a 2-layers multi-task learning design for hierarchical visual embedding. The architecture details are shown in Figure 1 with the layer configurations listed in Table 2.

Table 2: Configuration of 2-layers visual embedding.

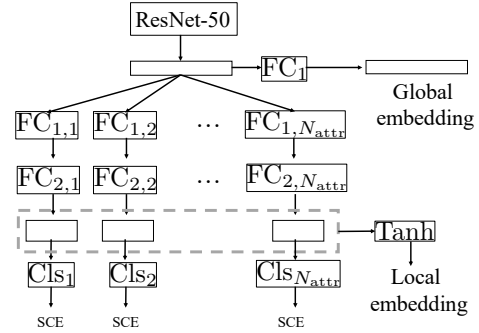| Structure | Size |
|---|---|
| ResNet-50 | Output size is 2048 |
| FC$_1$ | $2048 \times \text{Dim}_{emb}^{\text{glo}}$<br>Tanh |
| FC$_{1,i}$ (i=1, 2, ..., $N_{\text{Attr}}$) | $2048 \times 1024$<br>ReLU |
| FC$_{2,i}$ (i=1, 2, ..., $N_{\text{Attr}}$) | $1024 \times \text{Dim}_{emb}^{\text{loc}}$ |
| Classification$_i$ (i=1, 2, ..., $N_{\text{Attr}}$) | $\text{Dim}_{emb}^{\text{loc}} \times N_{\text{Attr}_i}$ |



Figure 1: Design of 2-layers visual embedding. Cls: Classification. SCE: Softmax Cross Entropy loss function.

### 1.2. 4-Layers Hierarchical Visual Embedding

The 4-layers hierarchical visual embedding is by a tree-structured multi-task learning design. The architecture design is shown in Figure 2 with the layer configuration listed in Table 3.

Table 3: Configuration of 4-layers visual embedding.

| Structure | Size |
|---|---|
| ResNet-50 | Output size is 2048 |
| FC$_1$ | $2048 \times \text{Dim}_{emb}^{\text{glo}}$<br>Tanh |
| FC$_{1,i}$ (i=1, 2) | $2048 \times 1024$<br>ReLU |
| FC$_{2,i}$ (i=1, 2, 3, 4) | $1024 \times 512$<br>ReLU |
| FC$_{3,i}$ (i=1, 2, ..., $N_{\text{Attr}}$) | $512 \times \text{Dim}_{emb}^{\text{loc}}$ |
| Classification$_i$ (i=1, 2, ..., $N_{\text{Attr}}$) | $\text{Dim}_{emb}^{\text{loc}} \times N_{\text{Attr}_i}$ |

## 2. Hypercritical Textual Embedding Network

The textual embedding consists of two parts: (1) local textual embedding and (2) global textual embedding. Simi-
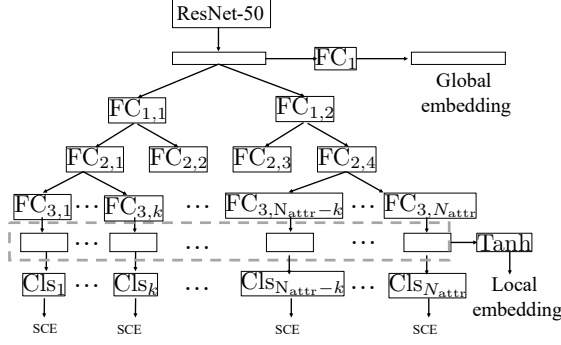
Figure 2: Design of 4-layers visual embedding. Cls: Classification. SCE: Softmax Cross Entropy loss function. The hierarchy is built as a balanced binary tree structure as possible, i.e. $k \simeq N_{\mathrm{attr}}/4$.

larly, we describe 2-layers and 4-layers hierarchical textual embedding, respectively.

## 2.1. 2-Layers Hierarchical Textual Embedding

In textual embedding, the input is a set of text attributes. Each text attribute is firstly passed into a word2vector model trained on wikipedia [1] and then into three FC layers, The resulting local embeddings are then utilised to form the global embedding. See the architecture in Figure 3 and layer configuration in Table 4.

Table 4: Configuration of 2-layers textual embedding. The setting of Conv layers: the number of input channels, the number of output channels, kernel size, stride, and padding. Cls: Classification.

| Structure | Size |
|---|---|
| $FC_1$ | $300 \times 512$ <br> Tanh |
| $FC_2$ | $512 \times 1024$ <br> Tanh |
| $FC_3$ | $1024 \times \mathrm{Dim}_{emb}^{loc}$ <br> Tanh |
| $Cls_i$ (i=1, 2, ..., $N_{\mathrm{Attr}}$) | $\mathrm{Dim}_{emb}^{loc} \times N_{\mathrm{Attr}_i}$ |
| $Fusion_1$ | $\mathrm{Conv}[\mathrm{Dim}_{emb}^{loc}, \mathrm{Dim}_{emb}^{glo}, 1, 1, 0]$ <br> $\mathrm{Conv}[N_{\mathrm{Attr}}, 1, 1, 1, 0]$ <br> Tanh |

## 2.2. 4-Layers Hierarchical Textual Embedding

The 4-layers textual embedding is in a similar structure as the 2-layers counterpart. See the architecture in Figure 4 and layer configuration in Table 5.
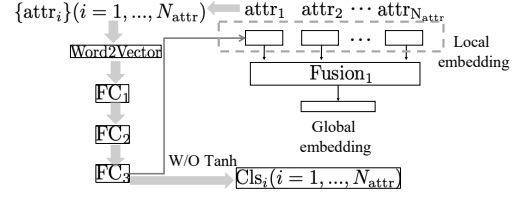


Figure 3: Design of 2-layers textual embedding. Cls: Classification.

Table 5: Configuration of 4-layers textual embedding. Cls: Classification.

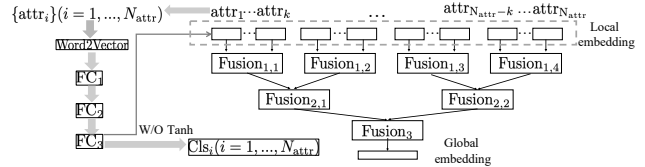| Structure | Size |
|---|---|
| $FC_1$ | $300 \times 512$ <br> Tanh |
| $FC_2$ | $512 \times 1024$ <br> Tanh |
| $FC_3$ | $1024 \times \mathrm{Dim}_{emb}^{loc}$ <br> Tanh |
| $Cls_i$ (i=1, 2, ..., $N_{\mathrm{Attr}}$) | $\mathrm{Dim}_{emb}^{loc} \times N_{\mathrm{Attr}_i}$ |
| $Fusion_{1,i}$ (i=1,2,3,4) | $\mathrm{Conv}[\mathrm{Dim}_{emb}^{loc}, 512, 1, 1, 0]$ <br> $\mathrm{Conv}[k,1,1,1,0]$ <br> Tanh |
| $Fusion_{2,i}$ (i=1,2) | $\mathrm{Conv}[512, 512, 1, 1, 0]$ <br> $\mathrm{Conv}[2,1,1,1,0]$ <br> Tanh |
| $Fusion_3$ | $\mathrm{Conv}[512, \mathrm{Dim}_{emb}^{glo}, 1, 1, 0]$ <br> $\mathrm{Conv}[2,1,1,1,0]$ <br> Tanh |



Figure 4: Design of 4-layers textual embedding. Cls: Classification. The hierarchy is built as a balanced binary tree structure as visual embedding, i.e. $k \simeq N_{\mathrm{attr}}/4$.

## 3. Cross-Modality Cross-Level Embedding

Given the hierarchical visual and textual embedding, we then conduct global-level and local-level cross-modality embedding followed with cross-level cross-modality embedding. The configuration of layers are listed in Table 6.

### 3.1. Cross-Modality Global-Level Embedding

The global-level fusion module takes as input the global visual embedding $x^{\mathrm{glo}}$ and the global textual embedding

Table 6: Configuration of cross-level (CL) cross-modality (CM) embedding.

| Structure | Size |
|---|---|
| $FC_{T/V}^i$, $i \in \{glo, loc\}$ | $Dim_{emb}^j \times 512$ |
| $FC_1$ | $512 \times 512$<br>Tanh |
| $FC_2$ | $512 \times 512$<br>Tanh |
| $Fusion_{CM}$ | $Conv[512, Dim_{emb}^s, 1, 1, 0]$<br>$Conv[N_{attr}, 1, 1, 1, 0]$<br>Tanh |
| $Fusion_{CL}$ | $Conv[512, Dim_{emb}^s, 1, 1, 0]$<br>$Conv[2, 1, 1, 1, 0]$<br>Tanh |

$z^{glo}$, outputting the global cross-modality embedding $s^{glo}$. The architecture is shown in Figure 5 (a).

### 3.2. Cross-Modality Local-level Embedding

The local-level fusion module takes as input the local visual embedding $\{x_i^{loc}\}_{i=1}^{N_{att}}$ and the local textual embedding $\{z_i^{loc}\}_{i=1}^{N_{att}}$, outputting the local cross-modality embedding $s^{loc}$. The architecture is given in Figure 5 (b).
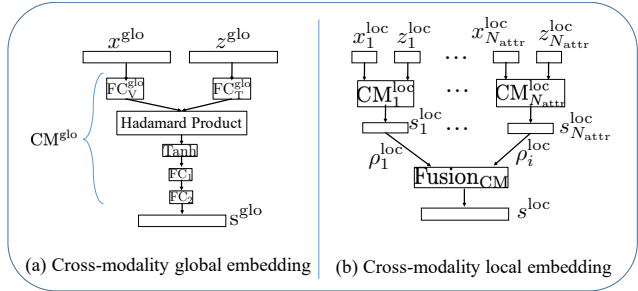


(a) Cross-modality global embedding    (b) Cross-modality local embedding

Figure 5: Design of cross-modality (a) global-level and (b) local-level embedding.

### 3.3. Cross-Modality Cross-Level Embedding

Given the global $s^{glo}$ and local $s^{loc}$ cross-modality embedding, we obtain the cross-modality cross-level embedding $s$ as shown in Figure 6.
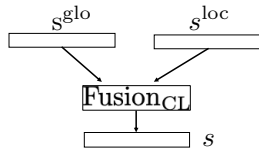


Figure 6: Design of cross-level (CL) embedding.

## 4. Matching Module

The matching module takes as input the cross-modality cross-level embedding $s$, and outputs the similarity score $\hat{y} \in [0, 1]$ of the input image and attribute set. In training, we set the ground-truth similarity score 1 for the *matching* attribute-image pairs and 0 for the *unmatched* attribute-image pairs. The details are shown in Table 7 and Figure 7.

Table 7: Configuration of matching module.

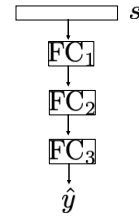| Structure | Size |
|---|---|
| $FC_1$ | $Dim_{emb}^s \times 256$<br>ReLU |
| $FC_2$ | $256 \times 128$<br>ReLU |
| $FC_3$ | $128 \times 1$<br>Sigmoid |



Figure 7: Design of matching module.

## References

[1] T. Mikolov, I. Sutskever, K. Chen, G. S. Corrado, and J. Dean. Distributed representations of words and phrases and their compositionality. In *NIPS*, 2013. 2