# HIERARCHICAL ATTENTION FUSION FOR GEO-LOCALIZATION

*Liqi Yan[1], Yiming Cui[2], Yingjie Chen[3], Dongfang Liu[3*]*

[1]Department of Computer Science and Technology, Fudan University
[2]Department of Electrical and Computer Engineering, University of Florida
[3]Department of Computer Graphics Technology, Purdue University

## ABSTRACT

Geo-localization is a critical task in computer vision. In this work, we cast the geo-localization as a 2D image retrieval task. Current state-of-the-art methods for 2D geo-localization are not robust to locate a scene with drastic scale variations because they only exploit features from one semantic level for image representations. To address this limitation, we introduce a hierarchical attention fusion network using multi-scale features for geo-localization. We extract the hierarchical feature maps from a convolutional neural network (CNN) and organically fuse the extracted features for image representations. Our training is self-supervised using adaptive weights to control the attention of feature emphasis from each hierarchical level. Evaluation results on the image retrieval and the large-scale geo-localization benchmarks indicate that our method outperforms the existing state-of-the-art methods. Code is available here: https://github.com/YanLiqi/HAF.

***Index Terms***— Geo-localization, hierarchical attention, multi-scale feature extraction, image retrieval.
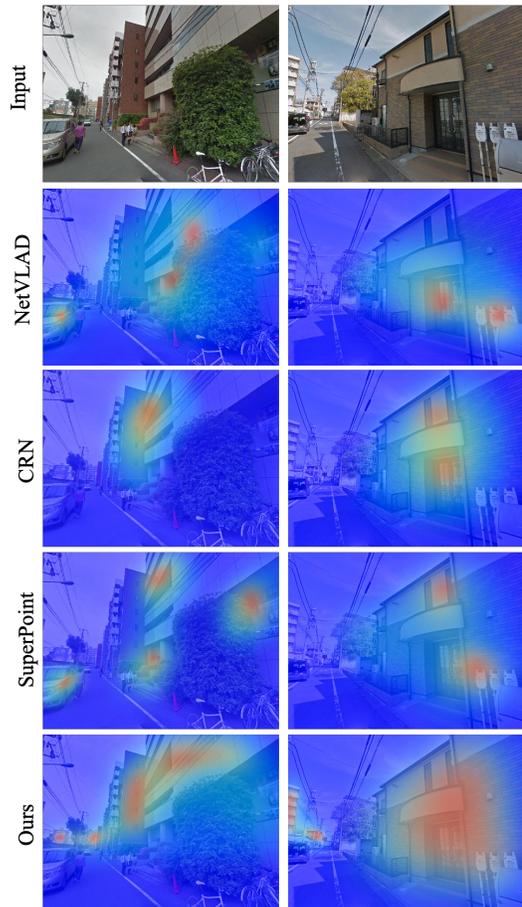
## 1. INTRODUCTION

Geo-localization is an important task in computer vision as it holds valuable potentials for applications such as autonomous driving [1] and robot navigation [2]. When working under the region with poor GPS signals, mobile agents require a supplementary localization for operation and geo-localization is a helpful addition to GPS [3].

In our work, we cast the geo-localization problem as a task of image retrieval [4], which searches over a pre-stored GPS-tagged image database to determine the current location [5]. The GPS-tagged image from the database with the closest distance to the query image in feature space is approximated as the location of the query [6].

### 1.1. Related Work

Nowadays, convolutional neural networks (CNNs) have become a powerful technique to explore image representations

---

*Dongfang Liu is the corresponding author.



**Fig. 1**: Comparison of feature emphasis. Compared to conventional methods [4, 7, 8], our method exploits the multi-scale features for hierarchical attention to depict image representation of landmarks with different scales and distance.

[9, 10]. The primary challenge for geo-localization is to produce discriminative image representation to identify two different places [3, 5]. A geo-localization database generally contains images having landmarks with different scales. For conventional methods [4, 7, 8], landmarks, with medium or small sizes, are difficult to be recognized because CNNs intend to downsample the spatial resolution of the input image by a significant margin [11,12]. However, many medium- and

small-size landmarks include valuable distinctiveness in image representation for geo-localization [7]. Thus, our method wants to exploit multi-scale features to delineate landmarks with different distances and scales (Fig. 1).

A critical reason for the aforementioned problem is that concurrent methods [4, 6–8] only use features from one semantic level for the geo-localization task. The feature maps from a single semantic level fail to fully explore rich visual clues from landmarks of different scales. This observation motivates us to exploit hierarchical features with different semantics to improve the geo-localization task.

## 1.2. Principal Contributions

Our work brings the following three contributions. First, we introduce a hierarchical attention fusion network, a novel algorithm for geo-localization (Fig. 2.). We find inspiration from the feature pyramid networks [11] which uses hierarchical feature maps to predict multi-scale image representations. Different from the conventional methods [4, 8], we extract hierarchical features at the lower-level, mid-level, and higher-level layer of the backbone network. Thus, the obtained features have rich multi-scale information. We then perform a feature fusion over the obtained features, so that our method can simultaneously pay attention to landmarks with various scales to identify discriminative keypoints for geo-localization. Second, we propose a self-supervised loss function to captures pairwise image relationships in training. Our training only needs GPS-tagged image pairs instead of expensive pixel-wise correspondences between images. Our training strategy can encourage the trained algorithm to learn which imagery context should be focused or suppressed to achieve a better image representation for localization. Finally, we use extensive experiments to assess our method. Results demonstrate that the proposed method sets a new state-of-the-art on several geo-localization benchmarks.

## 2. METHOD

This section details our method. The architecture of our method is shown in Fig. 2, which includes two principal modules: (i). Hierarchical feature extraction, and (ii). Feature fusion decoder. Using CNN, we first extract hierarchical feature maps to close the semantic gap in feature learning. We then perform attention fusion over the obtained multi-scale features to predict strong image representation.

## 2.1. Hierarchical Feature Extraction

We use VGG16 [9] as the backbone network for feature extraction. We extract hierarchical features from `Con3_2`, `Con4_3`, and `Con5_3` respectively (Fig. 2.). The corresponding feature maps are $\{C3, C4, C5\}$ and each has strides of $\{4, 8, 16\}$ pixels with respect to the input image. The obtained

hierarchical feature maps are then processed by a modified SuperPoint structure [8], which includes a non-linear $1 \times 1$ convolutional layer to control the channel dimensions and an upsampling layer to increase the feature map resolution. Specifically, the non-linear $1 \times 1$ convolutional layer uses `ReLU6` to determine the feature activation. The upsampling layer increases the resolution of the feature maps in a non-learned manner. In order to reduce the aliasing effect from upsampling, we perform the Adaptive Spatial Fusion on each output [13]. The output set of feature maps is $\{F_l, F_m, F_h\}$ corresponding to $\{C3, C4, C5\}$ as lower-level, mid-level, and higher-level features. Compared to $\{C3, C4, C5\}$, the channel dimensions for $\{F_l, F_m, F_h\}$ are multiplied $1.5\times$, $1\times$, $0.5\times$ times respectively and their spatial resolutions is brought back to half of the input resolution. $\{F_l, F_m, F_h\}$ together has rich features of 1,152 dimensions.

## 2.2. Attention Fusion Decoder

**Feature attention mask.** We denote the channel dimensions for $F_l, F_m, F_h$ as $x$, $y$, and $z$ respectively. Thus, $F_l, F_m, F_h$ can be expressed as a set of feature representations as $\{f_1, ..., f_x\}$, $\{f_1, ..., f_y\}$, and $\{f_1, ..., f_z\}$ accordingly, where $f$ is one feature map from a single channel.

Inspired by [7], we implement three learnable feature attention masks $\{m_1, m_2, m_3\}$ which are appended to $F_l, F_m, F_h$ separately. The attention mask is to indicate which spatial regions from the feature maps are discriminatively representative for localization. We define the attention-weighted features as:

$$F_l' = \sum_{n=1}^{x} \sum_{r \in R} m^r \cdot f_n^r, F_m' = \sum_{n=1}^{y} \sum_{r \in R} m^r \cdot f_n^r, F_h' = \sum_{n=1}^{z} \sum_{r \in R} m^r \cdot f_n^r \tag{1}$$
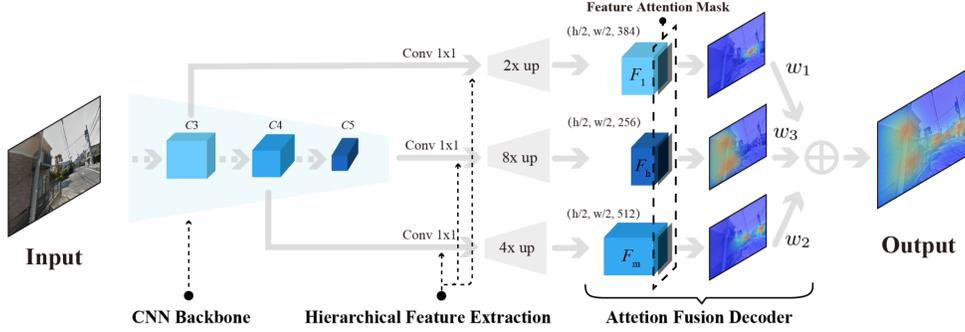
where $R$ denotes a set of spatial regions on the feature map. During backpropagation, each $m$ is learned to emphasize or suppress certain features from different spatial regions to encourage discriminative representation.

**Coupled descriptor and detector.** Similar to [8], we implement a coupled detector and descriptor based on the weighted features to prevent information loss. Using the attention-weighted features $F'$, we define descriptor as a set of vectors $K$:

$$K = \sum_{i=1}^{h} \sum_{j=1}^{w} F'^{ij:}, K^{ij} \in \mathbb{R}^x. \tag{2}$$

The obtained descriptors $K$ are $L_2$-normalized to be a unit length. At a pixel point $(i, j)$ on the image, we calculate the Euclidean distance of each descriptor $K^{ij}$ between images to establish feature correspondences. For detectors, we exploit the feature maps $F'$ in the same manner. Thus the detectors $D$ can be denoted as:

$$D = \sum_{n=1}^{x} F'^{::n}, \quad D^n \in \mathbb{R}^{h \times w} \tag{3}$$

**Fig. 2**: The architecture of the proposed method. Our method uses hierarchical features to close the semantic gap in feature learning. We perform the attention fusion over the obtained features to produce strong image representation for landmarks with different scales.

If the pixel point $(i, j)$ is detected, we denote the most strong detection of all channels as $D^{(ij)n'}$ ($n' \in \mathbb{R}^x$) on the response maps. We then perform an image-wise normalization of the detection to obtain the detection score at a pixel $(i, j)$:

$$s_{ij} = \frac{D_l^{(ij)n'}}{\sum_{i'=1}^{h} \sum_{j'=1}^{w} D_l^{(ij)n'}}. \tag{4}$$

We perform attention fusion decoding by using all the detectors and descriptors from different hierarchical levels to predict the image representation (Fig. 2).

### 2.3. Training Objective

All our modules are trained in an end-to-end fashion which facilitates the task-relevant feature learning. Given a query image $I_q$, our goal is to approximate its location by finding the reference images $\{I_r\}$ which are the nearest neighbors in feature space. To achieve hierarchical attention fusion, we propose a novel triplet ranking loss to jointly optimize the detectors and descriptors based on the different hierarchical features. Our training is weakly supervised. Thus, instead of using expensive feature correspondences at the pixel level to learn, our training only needs image-level annotations (the positive references $\{I_r^+\}$ and the negative references $\{I_r^-\}$). Namely, our method is trained to match the positive references $\{I_r^+\}$ and discriminate the negative ones $\{I_r^-\}$.

For a pair of image $(I_q, I_r)$, we define their feature differences by calculating the the descriptor distance $\sum_{c \in \mathcal{C}} \| K_q^c - K_r^c \|_2$, where $\mathcal{C}$ indicates all the corresponding feature points between the two images. In training, we maximize the distance of the corresponding descriptors between the negative pairs while minimizing the distance between the positive ones. Additionally, in order to increase the detection repeatability [14], we include a detection term to compute differences in feature space between two images:

$$\Delta \mathcal{D}(I_q, I_r) = \sum_{c \in \mathcal{C}} \frac{s_q^c s_r^c}{\sum_{c' \in \mathcal{C}} s_q^{c'} s_r^{c'}} \| K_q^c - K_r^c \|_2 . \tag{5}$$

where $s$ is the detection scores in (4). Thus, the triple tranking loss is defined as:

$$\mathcal{L}(I_q, I_r^+, I_r^-) = \max\left(M + \Delta\mathcal{D}(I_q, I_r^+) - \Delta\mathcal{D}(I_q, I_r^-), 0\right). \tag{6}$$

Since we jointly optimize the detectors and descriptors based on the different hierarchical features, our overall loss is:

$$\mathcal{L}_{total} = w_1 \cdot \mathcal{L}_l + w_2 \cdot \mathcal{L}_m + w_3 \cdot \mathcal{L}_h, \tag{7}$$

where $\mathcal{L}_l$, $\mathcal{L}_m$, and $\mathcal{L}_h$ are individual loss for each hierarchical attention. $w$ is the adaptive weight ($\sum_{i=1}^{3} w_i = 1$) which determines the contribution of each hierarchical attention to the final prediction. Using the proposed loss function, our method effectively learns which features need to be suppressed or emphasized for image representation.

## 3. EXPERIMENTS AND RESULTS

This section details our experiments to evaluate the proposed method on several benchmark detasets.

### 3.1. Implementation Setup

In training, we use the margin m = 0.1, 30 epochs, learning rate 0.0001 which is halved in every 5 epochs, momentum 0.9, weight decay 0.001, and a batch size of 4 triplets. We use the Precision-Recall curve to evaluate the training performance [4]. The trained models which yield the best $recall@5$ on the validation set is used for testing. We utilize a grid search to find the best adaptive weights in training. Eventually, we have $w_1 = 0.1$, $w_2 = 0.4$, and $w_3 = 0.5$ respectively.

### 3.2. Evaluation Datasets and Metrics

We evaluate our method on two types of benchmarks: (i) Image retrieval datasets, which are Oxford 5k [15], Paris 6k [16], and Holidays [17]. We employ the mean-Average-Precision (mAP) in our evaluation; and (ii) Geo-localization datasets, which are Pitts250k-test [18], Tokyo 24/7 [19], TokyoTM-val [4], and Sf-0 [20]. We use the Precision-Recall curve to test the performance of geo-localization.

### 3.3. Empirical Results

To assess the benefits of the proposed method, we compare our method with the state-of-the-art methods, NetLAVD [4], CRN [7], and SuperPoint [8], on geo-localization and image retrieval benchmarks. In order to have a fair comparison, we retrain all the methods with the same setup.

**Image retrieval benchmarks.** To test the generalizability of our approach, our method is trained only on Pitts30k [4] without any fine-tuning on the image retrieval datasets. For Oxford 5k [15] and Paris 6k [16], we use both the full and cropped images; for Holidays [17], we use original and rotated images. The results are displayed in Table 1. Our results set the state-of-the-art for compact image representations (256-D) on all three benchmarks. On all metrics, our margins consistently exceed the mAP of other methods by 1 to 5%. For example, there are a 3.86% improvements on Oxford 5k (full) than the next best method; and there are a 4% improvements on Oxford 5k (crop) than the next best method. Our methods can be further improved by fine-tuning using the three image retrieval datasets.

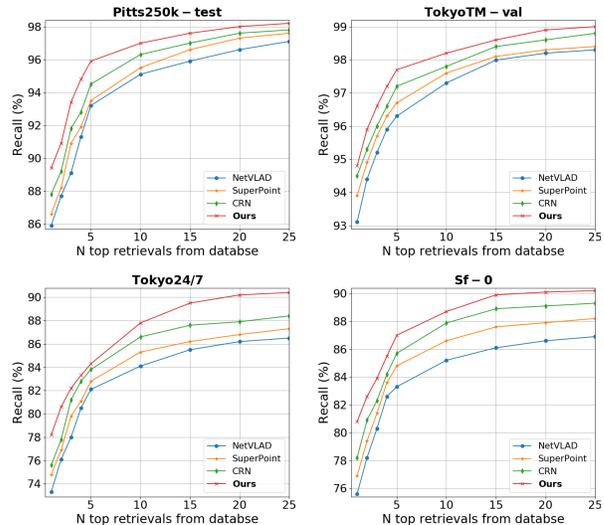| Method | Oxford 5K | | Paris 6k | | Holidays | |
|---|---|---|---|---|---|---|
| | full | crop | full | crop | orig | rot |
| Ours | **67.81** | **69.52** | **75.10** | **78.29** | **84.82** | **88.41** |
| CRN | 63.95 | 65.52 | 72.88 | 75.85 | 83.19 | 87.30 |
| NetVLAD | 63.09 | 65.33 | 72.53 | 75.67 | 82.67 | 86.83 |
| SuperPoint | 63.14 | 65.50 | 72.83 | 75.10 | 82.92 | 86.90 |

**Table 1**: Results for compact image representations (256-D).

**Geo-localization benchmarks.** We report the Precision-Recall plot for each method in Fig 3. Our method outperforms other methods under different $recall@n$ thresholds on all benchmarks. For qualitative analysis, we use the heatmap to visualize the feature emphasis for localization using the input image (Fig. 1). The qualitative examples illuminate that our method can effectively exploit multi-scale features and demonstrate the capacity of having hierarchical attentions on landmarks with different scales and distances for geo-localization. In contrast, other methods mainly focus on large-scale landmarks for discriminative visual clues. Our method also focuses on the distinctive details of buildings while avoiding confusing objects such as pedestrians, vegetation, or vehicles which are hard for feature repeatability.

Both quantitative and qualitative results advocate our hypothesis that we can leverage multi-scale features which are fused for hierarchical attention to produce discriminative yet compact image representations.

### 3.4. Adaptive Weight Analysis

In training, we train a set of models using different combinations of weights with a change of 0.1. The adaptive weight $w_1$, $w_2$, and $w_3$ controls lower-level features (small scale), mid-level features (middle scale), and higher-level features



**Fig. 3**: Comparison of recalls at $N$ top retrievals with the state-of-the-arts methods.

(large scale) respectively. We report the best adaptive weight which produces the best $recall@5$ results for each benchmark in Table 2. Pitts 250k-test dataset focuses on middle- and large-scale buildings. Thus, for Pitts 250k-test, $w_2$ and $w_3$ are much larger than $w_1$. TokyoTM dataset generally includes small-, middle-, and large-scale buildings which can be reflected from its best weights. For Tokyo 24/7, it has a similar adaptive weight as TokyoTM. Although Tokyo 24/7 has less small-scale buildings, it includes a lot of landmark details such as billboards, city lights, or traffic signs by the road. Sf-0 has a dominant $w_3$ as it mainly focuses on buildings with a large scale.

| Method | Pitts 250k-test | TokyoTM-val | Tokyo 24/7 | Sf-0 |
|---|---|---|---|---|
| $w_1$ | 0.1 | 0.3 | 0.2 | 0.1 |
| $w_2$ | 0.4 | 0.3 | 0.3 | 0.1 |
| $w_3$ | 0.5 | 0.4 | 0.5 | 0.8 |

**Table 2**: Best adaptive weights for each benchmarks.

### 4. CONCLUSION

In this work, we introduce a hierarchical attention fusion network for geo-localization. We extract the multi-scale feature maps from a convolutional neural network (CNN) to perform hierarchical attention fusion for image representations. Since the hierarchical features are scale-sensitive, our method is robust to landmarks with different scales and distances. We evaluate our method extensively on the image retrieval benchmarks and the large-scale geo-localization benchmarks. Results indicate that our method is competitive with the latest state-of-the-art approaches.

## 5. REFERENCES

[1] D. Liu, Y. Cui, Z. Cao, and Y. Chen, "A large-scale simulation dataset: Boost the detection accuracy for special weather conditions," in *2020 IJCNN*, 2020, pp. 1–8.

[2] D. Liu, Y. Cui, Z. Cao, and Y. Chen, "Indoor navigation for mobile agents: A multimodal vision fusion model," in *2020 IJCNN*, 2020, pp. 1–8.

[3] Dongfang Liu, Yiming Cui, Xiaolei Guo, Wei Ding, Baijian Yang, and Yingjie Chen, "Visual localization for autonomous driving: Mapping the accurate location in the city maze," *arXiv preprint arXiv:2008.05678*, 2020.

[4] Relja Arandjelovic, Petr Gronat, Akihiko Torii, Tomas Pajdla, and Josef Sivic, "Netvlad: Cnn architecture for weakly supervised place recognition," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2016, pp. 5297–5307.

[5] Mahdi Salarian, Nick Iliev, Ahmet Enis Cetin, and Rashid Ansari, "Improved image-based localization using sfm and modified coordinate system transfer," *IEEE Transactions on Multimedia*, vol. 20, no. 12, pp. 3298–3310, 2018.

[6] Paul-Edouard Sarlin, Cesar Cadena, Roland Siegwart, and Marcin Dymczyk, "From coarse to fine: Robust hierarchical localization at large scale," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 12716–12725.

[7] Hyo Jin Kim, Enrique Dunn, and Jan-Michael Frahm, "Learned contextual feature reweighting for image geo-localization," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 2136–2145.

[8] Daniel DeTone, Tomasz Malisiewicz, and Andrew Rabinovich, "Superpoint: Self-supervised interest point detection and description," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition Workshops*, 2018, pp. 224–236.

[9] Karen Simonyan and Andrew Zisserman, "Very deep convolutional networks for large-scale image recognition," *arXiv preprint arXiv:1409.1556*, 2014.

[10] Dongfang Liu, Yiming Cui, Yingjie Chen, Jiyong Zhang, and Bin Fan, "Video object detection for autonomous driving: Motion-aid feature calibration," *Neurocomputing*, vol. 409, pp. 1 – 11, 2020.

[11] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie, "Feature pyramid networks for object detection," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 2117–2125.

[12] Gao Huang, Zhuang Liu, Laurens Van Der Maaten, and Kilian Q Weinberger, "Densely connected convolutional networks," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2017, pp. 4700–4708.

[13] Chaoxu Guo, Bin Fan, Qian Zhang, Shiming Xiang, and Chunhong Pan, "Augfpn: Improving multi-scale feature learning for object detection," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2020, pp. 12595–12604.

[14] Mihai Dusmanu, Ignacio Rocco, Tomas Pajdla, Marc Pollefeys, Josef Sivic, Akihiko Torii, and Torsten Sattler, "D2-net: A trainable cnn for joint detection and description of local features," *arXiv preprint arXiv:1905.03561*, 2019.

[15] James Philbin, Ondrej Chum, Michael Isard, Josef Sivic, and Andrew Zisserman, "Object retrieval with large vocabularies and fast spatial matching," in *2007 IEEE conference on computer vision and pattern recognition*. IEEE, 2007, pp. 1–8.

[16] James Philbin, Ondrej Chum, Michael Isard, Josef Sivic, and Andrew Zisserman, "Lost in quantization: Improving particular object retrieval in large scale image databases," in *2008 IEEE conference on computer vision and pattern recognition*. IEEE, 2008, pp. 1–8.

[17] Herve Jegou, Matthijs Douze, and Cordelia Schmid, "Hamming embedding and weak geometric consistency for large scale image search," in *European conference on computer vision*. Springer, 2008, pp. 304–317.

[18] Akihiko Torii, Josef Sivic, Tomas Pajdla, and Masatoshi Okutomi, "Visual place recognition with repetitive structures," in *Proceedings of the IEEE conference on computer vision and pattern recognition*, 2013, pp. 883–890.

[19] Akihiko Torii, Relja Arandjelovic, Josef Sivic, Masatoshi Okutomi, and Tomas Pajdla, "24/7 place recognition by view synthesis," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 1808–1817.

[20] David M Chen, Georges Baatz, Kevin Köser, Sam S Tsai, Ramakrishna Vedantham, Timo Pylvänäinen, Kimmo Roimela, Xin Chen, Jeff Bach, Marc Pollefeys, et al., "City-scale landmark identification on mobile devices," in *CVPR 2011*. IEEE, 2011, pp. 737–744.