

# DINO-Mix: Enhancing Visual Place Recognition with Foundational Vision Model and Feature Mixing

Gaoshuang Huang, Yang Zhou\*, Xiaofei Hu, Chenglong Zhang, Luying Zhao, Wenjian Gan, and Mingbo Hou

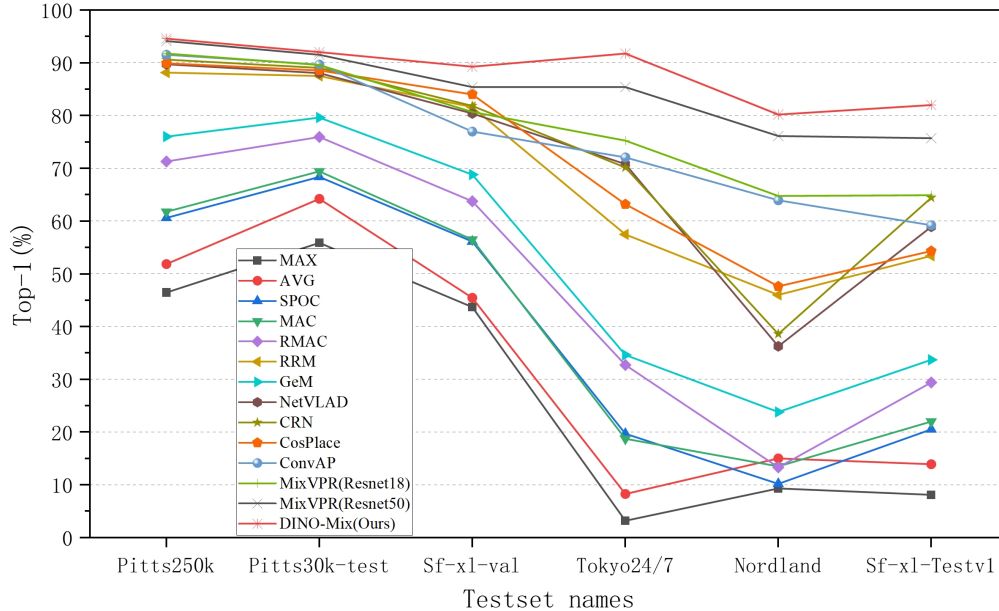


Fig. 1: Comparison of DINO-Mix and other VPR methods in terms of Top-1 accuracy for different test sets

**Abstract**—In the vast expanse of cyberspace, a plethora of publicly available images exist. Utilizing visual place recognition (VPR) technology to ascertain the geographical location of publicly available images is a pressing issue for real-world VPR applications. Although most current VPR methods achieve favorable results under ideal conditions, their performance in complex environments, characterized by lighting variations, seasonal changes, and occlusions caused by moving objects, is generally unsatisfactory. Therefore, obtaining efficient and robust image feature descriptors even in complex environments is a pressing issue in VPR applications. In this study, we utilize the DINOv2 model as the backbone network for trimming and fine-tuning to extract robust image features. We propose a novel VPR architecture called DINO-Mix, which combines a foundational vision model with feature aggregation. This architecture relies on the powerful image feature extraction capabilities of foundational vision models. We employ an MLP-Mixer-based mix module to aggregate image features, resulting in globally robust and generalizable descriptors that enable high-precision VPR. We experimentally demonstrate that the proposed DINO-Mix architecture significantly outperforms current state-of-the-art (SOTA) methods. In test sets having lighting variations, seasonal changes, and occlusions (Tokyo24/7, Nordland, SF-XL-Testv1), our proposed DINO-Mix architecture achieved Top-

1 accuracy rates of 91.75%, 80.18%, and 82%, respectively. Compared with SOTA methods, our architecture exhibited an average accuracy improvement of 5.14%. To further evaluate the performance of DINO-Mix, we compared it with other SOTA methods using representative image retrieval case studies. Our analysis revealed that DINO-Mix outperforms its competitors in terms of VPR performance. Furthermore, we visualized the attention maps of DINO-Mix and other methods to provide a more intuitive understanding of their respective strengths. These visualizations serve as compelling evidence of the superiority of the DINO-Mix framework in this domain. Code is available at <https://github.com/GaoShuang98/DINO-Mix>.

## I. INTRODUCTION

Visual place recognition (VPR), also known as image geo-localization (IG) or visual geo-localization (VG), has been extensively applied in various fields, such as cyberspace mapping, intelligence gathering, image target localization, autonomous driving, and outdoor user localization. Currently, most VPR studies focus on image retrieval in urban scenarios. However, images captured in urban environments may have varying shooting angles, temporal lighting changes, seasonal variations, occlusions, and similar repetitive textures. These conditions can pose significant difficulty for achieving high-precision image retrieval. Therefore, extracting robust and generalizable image feature descriptors for accurate image retrieval is a critical issue.

In previous approaches to VPR, handcrafted SIFT [3],

This paper is supported in part by the National Natural Science Foundation of China (NSFC) under Grant No.42001338.

\* Corresponding author

All the authors are with the Institute of Geospatial Information, PLA Strategic Support Force Information Engineering University, Zheng Zhou, China. E-mail: huanggaoshuang123@163.com, zhouyang3d@163.com, huxiaofei@163.com, z1204149693@163.com, 2014202130068@whu.edu.cn, 14737117985@163.com, hobefrank@163.com.

HOG [13], SURF [17], ORB [18], and other techniques were used to extract features from images. These features were then aggregated using methods such as bag-of-words (BoW) [30], Fisher Vector (FV) [31], and Vector of Locally Aggregated Descriptors (VLAD) [32] to obtain image descriptors for image retrieval, enabling image geo-localization. More recently, deep-learning techniques have emerged as mainstream approaches for extracting image features. Compared with handcrafted features, these methods significantly improve VPR accuracy. Examples include NetVLAD [14], which combines convolutional neural networks (CNNs) with VLAD, and other variants that incorporate attention, semantics, context, and multiscale features. Other methods based on the generalized mean (GeM) [11], such as CosPlace [24], and those utilizing fully-connected multilayer perceptron (MLP)-based feature aggregation, such as MixVPR [21], have been proposed. However, in practical testing, the accuracy of these methods for image geo-localization has been suboptimal under challenging conditions, such as varying shooting angles, temporal lighting changes, seasonal variations, and occlusions.

The rapid development of foundational visual models has enabled the generation of universal visual features from images [26]. By training on billions of data points, foundational visual models can extract image features that are more generalizable and robust than those extracted by conventional models. They can effectively handle the challenging conditions encountered in practice. Therefore, incorporating foundational visual models into VPR is a promising approach.

Considering the aforementioned challenges, this study proposes a method based on the DINOv2 [26] model, called DINO-Mix, which combines foundational visual models with feature aggregation. This architecture possesses exceptional discriminative power. Efficient and robust image features suitable for image geo-localization are extracted by fine-tuning and trimming. Furthermore, it utilizes a feature mixer [27] module to aggregate image features, resulting in a global feature descriptor vector. DINO-Mix is experimentally demonstrated to achieve superior test accuracy on multiple benchmarks, surpassing state-of-the-art (SOTA) methods.

The remainder of this paper is organized as follows. In **II**, we summarize previous relevant research on image geo-localization. In **III**, the DINO-Mix method is introduced. In **IV**, we provide details on the training set, testing set, training, and evaluation parameters used in our experiments. The proposed method is compared with existing methods in terms of accuracy, and ablation experiments are conducted. In **V**, We qualitatively demonstrate the state-of-the-art of DINO-Mix architecture by comparing DINO-Mix with other VPR methods through a typical VPR example and visualizing the corresponding attention map. Finally, our conclusions are presented in **VI**.

## II. RELATED WORK

By retrieving the most similar image from the image database, the geographical location of the retrieved image can be used as the location of the target image [29]. In

recent years, numerous researchers have made significant contributions to the field of image retrieval for VPR. The features used in image retrieval can be broadly categorized into handcrafted and deep features. Zhang and Kosecka [5] first extracted scale-invariant feature transform (SIFT) features from images to establish an image feature database. They performed a brute-force global search of the database and validated and ranked the top five candidate images using the random sample consensus (RANSAC) [6] algorithm. Finally, the geographical location of the target image was obtained by triangulating the top three images. Zamir and Shah [4] extracted SIFT feature vectors from images to build a database and employed a nearest-neighbor tree search to improve the retrieval efficiency. Zamir and Shah [2] further improved the nearest-neighbor matching technique by pruning outliers, and applied the generalized minimum clique problem (GMCP) in conjunction with approximate feature matching. This resulted in a 5% improvement in the localization accuracy compared to their previous work [4]. The advantages of using handcrafted features for VPR are their simplicity and strong interpretability. However, these methods tend to have high redundancy, require dimensionality reduction, are susceptible to environmental changes, and generally have low accuracy.

Deep features are extracted by neural networks with modules such as convolutional layers and attention mechanisms. These features often outperform handcrafted features owing to their strong expressive power, ability to freely define feature dimensions, and flexibility in designing neural network frameworks. Noh et al. [10] proposed a deep local feature (DELf) descriptor and an attention mechanism for keypoint selection to identify semantic local features. Ng et al. [22] introduced a global descriptor called Second-Order Loss and Attention for image Retrieval (SOLAR) that utilizes spatial attention and descriptor similarity to perform large-scale image retrieval using second-order information. Chu et al. [8] constructed a CNN to extract dense features, embedded an attention module within the network to score features, and proposed a grid feature point selection (GFS) method to reduce the number of image features. Chu et al. [9] combined deep features with handcrafted features, extracted average pooling features from the intermediate layers of a CNN for retrieval on street-view datasets, and used SIFT to re-rank them. Yan [7] extracted hierarchical feature maps from CNNs and organically fused them for image feature representation. Chu et al. [15] employed a CNN with a HOW module [23] to extract local image features, aggregated them into a feature vector using VLAD, used the aggregated selective match kernel (ASMk), and estimated the geographical location of the query image using kernel density prediction (KDP).

To address environmental factors, Mishkin et al. [12] employed a BoW method with multiple detectors, descriptors, and adaptive thresholds. Relja et al. [14] designed a trainable NetVLAD layer inspired by VLAD, which provides a pooling mechanism that can be integrated into other CNN structures. In addition, variants of NetVLAD have been proposed, such as CRN [33], SPE-VLAD [34], MultiRes-NetVLAD [35], SARE [36], and SFRS [37]. Ali-bey et al.

proposed ConvAP [20], which combines 1x1 convolutions with adaptive mean pooling to encode local features.

### III. METHODOLOGY

#### A. Proposed architecture

We propose an image-geolocation framework that integrates the foundational Vision model with feature aggregation. The proposed framework utilizes the truncated DINOv2 [26] model as the backbone network. Owing to its exceptional image understanding capability, DINOv2 is well-suited for various downstream tasks. Therefore, we pre-trained the DINOv2 model as the primary network for image feature extraction and employed an efficient and lightweight Mixer module to aggregate the obtained image features. The DINO-Mix visual geolocation architecture is illustrated in Fig. 2.

We modified the DINOv2 model by removing its layer norms and head modules, which were subsequently used as the backbone network. Furthermore, to maximize the pre-trained parameter benefits of the DINOv2 model for image understanding, we used the output from the last layer of the Vision Transformer (ViT) blocks [1] as the input to the Mixer module. Given that the output of the modified ViT block module is a feature matrix of size  $C \times D$  (channels  $\times$  feature vector length), we transform it into  $s$  feature maps of size  $h \times w$ , as expressed by Equation 1. These transformed feature maps serve as inputs for the mix module.

$$\begin{cases} D = s \\ C = hw \end{cases} \quad (1)$$

where  $D$  represents the length of the feature vector output by the backbone network,  $C$  denotes the number of channels in the output of the backbone network,  $h$  and  $w$  are the height and width of the feature map, respectively, and  $s$  is the number of feature maps.

#### B. Foundational vision model: DINOv2

Foundational vision models (FVMs) are typically constructed using structures such as CNNs or Transformers. These models often have parameters on the order of tens to hundreds of millions, giving them a greater representational capacity than smaller models. In addition, because of the use of larger and more diverse datasets during training, FVMs can learn more features and have better generalization capabilities.

DINOv2 [26] is capable of extracting powerful image features and performs well across various tasks. Compared with Segment Anything [38], DINOv2 has a broader scope of application and areas of use. The architecture of the DINOv2 model is illustrated in Fig. 3. First, an input image is passed through a patch-embedded module consisting of a two-dimensional (2D) convolutional layer with a kernel size of  $14 \times 14$  and a stride of 14, followed by a normalization layer. This process uniformly outputs patches of size  $14 \times 14$ . These patches are then fed into ViT blocks, which vary in number according to the size of the model. The ViT blocks

TABLE 1: *Four ViT model parameters for DINOv2*

Name	Patch embed	Blocks	Feature dim	Size(MB)
ViTs14	$14 \times 14$	12	384	86.2
ViTb14	$14 \times 14$	24	768	338.2
ViTl14	$14 \times 14$	24	1024	1189.0
ViTg14	$14 \times 14$	40	1536	4439.5

output a feature matrix of size  $C$  (number of channels)  $\times$   $D$  (dimension of the feature vector), which is then normalized by a layer-norm module before being transformed into a feature vector of size  $1 \times n$ . Finally, the head module can be flexibly selected based on specific image task requirements.

DINOv2 is characterized by several key features: a) it presents a novel approach for training high-performance computer vision models; b) it offers superior performance without the need for fine-tuning; c) it can learn from any image dataset and capture certain features that existing methods struggle with; and d) it leverages knowledge distillation to transfer knowledge from more complex teacher models to smaller student models. Through knowledge distillation, three smaller models were obtained from the ViTg14 model: ViTl14 (large), ViTb14 (base), and ViTs14 (small) (see Tab. 1).

The primary advantage of DINOv2 is the ability to create a large dataset for model training. This dataset, called LVD-142M, comprises 142 million images and includes ImageNet-22k, ImageNet-1k, Google Landmarks, various fine-grained datasets, and image datasets crawled from the Internet. For model training, an Nvidia A100 40-GB GPU was utilized, with a total of 22k GPU hours dedicated to training the DINOv2-g model.

#### C. Feature Mixer

Currently, the most advanced techniques propose shallow aggregation layers that are inserted into very deep pre-trained backbones cropped to the last feature-rich layer. By contrast, Wang et al. proposed TransVPR [19], which achieved good results in local feature matching. However, its global representation performance did not surpass that of NetVLAD [14] or CosPlace [24]. Recent advancements in isotropic architectures have demonstrated that self-attention is not crucial for ViT. However, Mixer utilizes feature maps extracted from a pre-trained backbone and iteratively merges global relationships into each feature map. This is achieved through an isotropic block stack composed of MLPs, referred to as a feature mixer [27]. The effectiveness of Mixer has been demonstrated through several qualitative and quantitative results, demonstrating its high performance and lightweight nature [21]; the architecture is illustrated in Fig. 4.

Mixer treats the input feature map  $F \in R^{(s \times h \times w)}$  as a set of  $s$  2D features, each of size  $h \times w$ , as expressed by Equation 2:

$$F = \{X^i, i = \{1, \dots, s\}\} \quad (2)$$

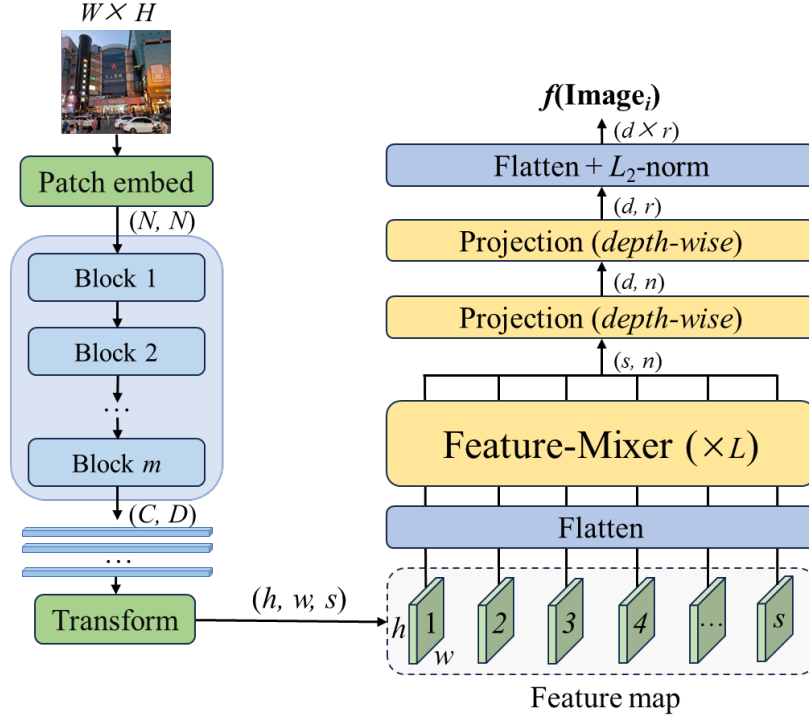


Fig. 2: The visual place recognition structure of DINO-Mix.

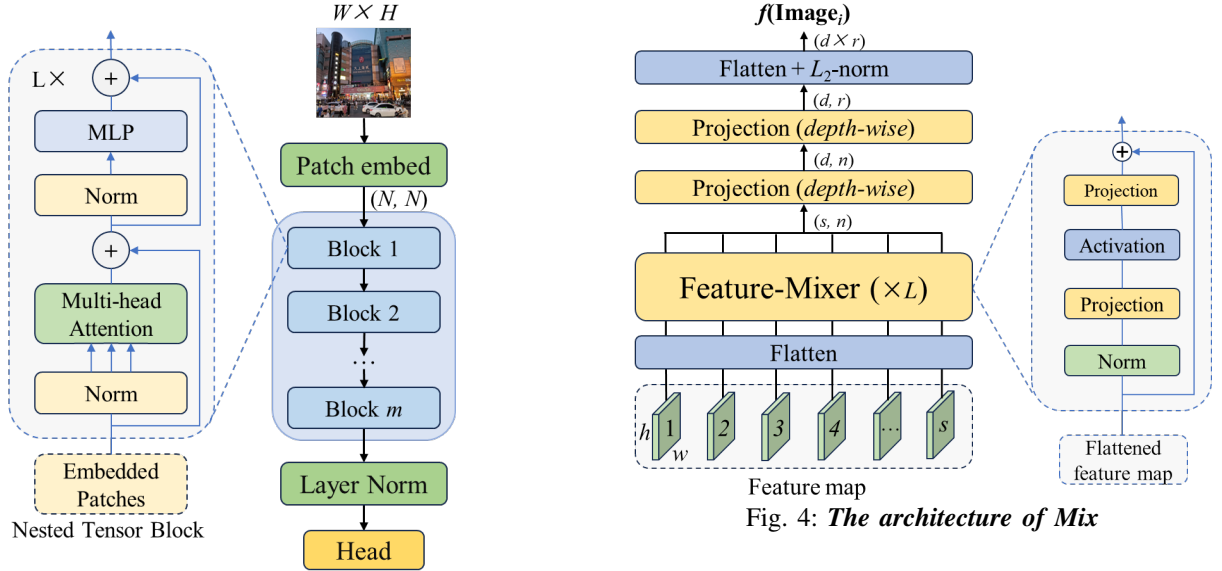


Fig. 4: The architecture of Mix

Fig. 3: The structural diagram of the DINOv2 model.

where  $X^i$  corresponds to the  $i$ th activation map in the feature map  $F$ . Secondly, each 2D feature map  $X^i$  is expanded into a 1D vector representation, resulting in a flattened feature map  $F \in R^{(s \times n)}$ , where  $n = h \times w$ .

The flattened feature maps are then fed into the feature mixer, which is composed of  $L$  MLPs with the same structure, as shown in Fig.4. The feature mixer takes the flattened feature map ensemble as input and successively incorporates spatial global relationships into each  $X^i \in F$  as per Equation

3:

$$X^i \leftarrow W_2(\sigma(W_1 X^i)) + X^i, i = \{1, \dots, s\} \quad (3)$$

where  $W_1$  and  $W_2$  are the weights of the two fully-connected layers that make up the MLP, and  $\sigma$  is the ReLU nonlinear activation function.

For  $F \in R^{(s \times n)}$ , the feature mixer, owing to its isotropy architecture, produces an output " $Z$ "  $\in R^{(s \times n)}$  with the same shape and feeds it into the second feature mixer block, and so on, until  $L$  consecutive blocks have been traversed, as per Equation 4:



$$Z = FM_L \left( FM_{(L-1)} (\dots FM_1(F)) \right) \quad (4)$$

where  $Z$  and the feature map  $F$  have the same dimensions. To control the dimensions of the final global descriptor, two fully-connected layers are used to successively transform the channel and row dimensions. First, a depth projection is used to map  $Z$  from  $R^{(s \times n)}$  to  $R^{(d \times n)}$ , as given by Equation 5:

$$Z' = W_d(\text{Transpose}(Z)) \quad (5)$$

where  $W_d$  denotes the weight of the fully-connected layer. Subsequently, a row-wise projection is used to map the output  $Z'$  from  $R^{(d \times n)}$  to  $R^{(d \times r)}$ , as given by Equation 6:

$$O = W_r(\text{Transpose}(Z')) \quad (6)$$

where  $W_r$  denotes the weight of the fully-connected layer. The final output  $O$  has dimensions of  $d \times r$ , which are flattened, and  $L2$  is normalized to form a global feature vector.

## IV. EXPERIMENTS

### A. Implementation details

**Datasets:** Our model was trained using the GSV-Cities dataset [20]. The following six datasets were employed for evaluation purposes: Pittsburgh250k [16] (contains 8k queries and 83k reference images collected from Google Street View and Pittsburgh30k-test), Pittsburgh30k-test [16] (a subset of Pittsburgh250k, with 8k queries and 8k reference images), SF-XL-Val dataset [24], Tokyo 24/7 [25], Nordland [45], and SF-XL-Testv1 [24]. The datasets contain extreme variations in lighting, weather, and seasons. Specific information regarding these datasets is presented in Tab.2.

**Architecture:** We employed the PyTorch deep-learning framework to implement DINO-Mix. To enable a fair comparison with other methods in terms of accuracy, we conducted precision tests on various VPR frameworks such as NetVLAD, GeM, ConvAP, CosPlace, and MixVPR, and obtained the testing accuracy for other methods from their corresponding papers.

**Training:** Owing to the excellent pre-trained weights of the DINO-Mix backbone, most of the training weights were frozen during the training process. However, to make it more suitable for the VPR task, we fine-tuned the end of the backbone and trained the feature aggregation module. We trained DINO-Mix following the standard framework proposed in GSV-Cities [20], which introduces a high-precision dataset consisting of 67k locations described by 560 k images. The batch size  $B$  was flexibly adjusted based on the model parameter size, and each location was trained with four images, resulting in a mini-batch of  $B \times 4$  images. Stochastic gradient descent [43] (SGD) with a momentum of 0.9 and a weight decay of 0.001 was employed for optimization. The initial learning rate was set to 0.05 and was divided by three every five epochs. Finally, the model was trained using images resized to  $224 \times 224$  pixels over 50 epochs. Most existing VPR studies employ a triplet loss function based

on weak supervision [44] for network training; however, this approach requires significant GPU memory and has a high computational overhead. Thus, we utilized multi-similarity loss [46] as the training loss function. Multi-similarity loss mitigates the issues of large interclass distances and small intraclass distances in metric learning by considering multiple similarities. Instead of relying solely on absolute spatial distance as the sole metric, it uses the overall distance distribution of other sample pairs within a batch size to weigh the loss. This computational approach effectively promotes model convergence in the early stages, as expressed by Equation 7:

$$L_{MS} = \frac{1}{m} \sum_{i=1}^m \left\{ \frac{1}{\alpha} \log \left[ 1 + \sum_{k \in P_i} e^{-\alpha(S_{ik} - \lambda)} \right] + \frac{1}{\beta} \log \left[ 1 + \sum_{k \in N_i} e^{\beta(S_{ik} - \lambda)} \right] \right\} \quad (7)$$

where  $P_i$  represents the set of positive sample pairs for each instance in a batch size,  $N_i$  denotes the set of negative sample pairs for each instance in the same batch size,  $S_{ij}$  and  $S_{ik}$  represent the similarities between the two images, and  $\alpha$ ,  $\beta$ , and  $\lambda$  are hyperparameters.

**Evaluation:** In this work, we employed top-k accuracy [28] as a metric to evaluate the precision of the VPR methods. Top-k accuracy is a commonly used evaluation method in the VPR domain, where it is considered successful if at least one of the top-k localization results for a query image has a geographical distance of less than a threshold  $s$  from the true location. In our experiments, we set  $s$  to  $25 m$  to align with existing methods.

### B. Comparison to the State-Of-The-Art

Based on the conclusions drawn from the ablation studies in this work, we adopt the ViTb14 pre-trained model, which exhibits the best performance among the four models of DINOv2 as the backbone network for DINO-Mix in the VPR task, and modify the DINOv2 model by removing its Layer Norm and Head modules. We utilize Mix as the feature aggregation module to construct the model. During training, we update the parameters of the last three blocks of the backbone network and the entire Mix feature aggregation module. The number of Feature Mixer blocks in the Mix feature aggregation module is set to 2, and the dimensionality of the image features output by the model is 4096. By utilizing these optimal parameter settings, we conduct tests on six test sets for DINO-Mix and compare it with existing methods, as shown in Tab.3. In addition, this paper presents Fig.1 to more visually demonstrate the accuracy difference between DINO-Mix and other major VPR methods.

We adopted the ViTb14 pre-trained model, which exhibited the best performance among the four models of DINOv2, as the backbone network for DINO-Mix in the VPR task and modified the DINOv2 model by removing its layer norm and head modules. We used Mixer as a feature aggregation module to construct the model. During training, we updated

TABLE 2: *The parameter table of the training dataset and test dataset*

Dataset	Train/val database queries	Test database queries	Dataset size (GB)	Database type	Database image size	Urban	Appearance changes	
							Season	Day/Night
GSV-cities [20]	524701/0	0/0	21.7	panorama	480×640	✓	✓	✓
Pittsburgh250k [16]	170112/15432	83952/8280	9.5	panorama	300×400	✓	✗	✗
Pittsburgh30k [16]	20000/15024	10000/6816	2.0	panorama	480×640	✓	✗	✗
Tokyo 24/7 [25]	0/0	75984/315	4.2	panorama	480×640	✓	✗	✓
Nordland [45]	0/0	27592/27592	1.3	font-view	360×640	✗	✓	✗
SF-XL-Val [24]	0/0	8015/7993	0.7	panorama	512×512	✓	✗	✗
SF-XL-Testv1 [24]	0/0	27191/1000	1.3	panorama	512×512	✓	✗	✓

the parameters of the last three blocks of the backbone network and the entire mix feature aggregation module. The number of feature mixer blocks in the mix feature aggregation module was set to two, and the dimensionality of the image features output by the model was 4096. By utilizing these optimal parameter settings, we conducted tests on six test sets for DINO-Mix and compared them with the existing methods, as shown in Tab.3. In addition, Fig.1 illustrates the difference in accuracy between DINO-Mix and other major VPR methods.

As listed in Tab.3, the test accuracy of the DINO-Mix model proposed in this paper has comprehensively surpassed that of the SOTA method, with further improvement in the Pittsburgh250k, Pittsburgh30k, and SF-XL-Val test sets focusing on changes in viewpoints, and especially in the Tokyo24/7, Nordland, and SF-XL-Testv1 test sets with changes in complex appearance environments.

### C. ablation studies

1) *Hyperparameters*: In DINO-Mix, the number of layers  $L$  in the feature mixer is also a critical factor for image retrieval accuracy. To determine the optimal number of Mixer layers, we conducted tests on the Pitts30k-test, Pitts250k-test, Sf-xl-val, Tokyo24/7, Nordland, and Sf-xl-testv1 datasets with different numbers of mix layers  $L$  (1, 2, 3, 4, 5, 6, 7) for DINO-Mix using ViTb14 as the backbone network. The TOP-1 accuracy is depicted in Fig.5. A careful examination of the figure reveals that without any Mix layers, DINO-Mix exhibits a lackluster test accuracy across all six datasets. However, upon incorporating one Mix layer, there is a remarkable enhancement in test accuracy. This observation highlights the pivotal role played by the feature aggregation module in elevating the precision of DINO-Mix. As the number of Mix layers further increases up to two, there is a marginal improvement in test accuracy, culminating in a peak. Nevertheless, as the number of Mix layers continues to escalate, DINO-Mix’s test accuracy on the six datasets displays a slow decline and fluctuations, accompanied by a linear increase in parameters. Based on the above analysis, this study adopts a two-layer Mix scheme as the feature aggregator in DINO-Mix.

2) *Descriptor dimensionality*: We conducted an ablation study on the dimensionality of image feature vectors ex-

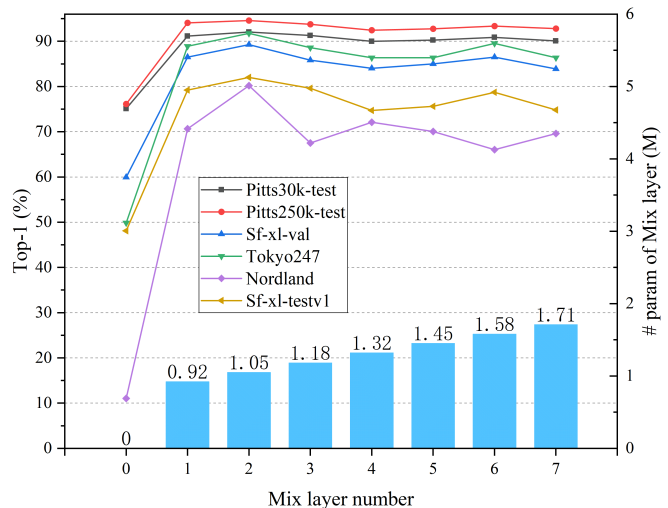


Fig. 5: *Ablation on the number of Feature-Mixer blocks*

tracted using the DINO-Mix. The experiment employed ViTb14, which exhibited the best performance as the backbone network, with two layers in the Mixer module. The test was used as the Pitts30k-test, Pitts250k-test, Sf-xl-val, Tokyo24/7, Nordland, and Sf-xl-testv1 datasets, and the image feature vector dimensionality was varied by changing the number of channels in the output vector of the Mixer module. The tested dimensions of the image feature vectors were 128, 256, 512, 1024, 2048, 4096, and 8192. As depicted in Fig.6, an increase in the dimensionality of image feature vectors is observed to have a positive impact on the overall Top-1 test accuracy of DINO-Mix across various datasets. This trend is particularly pronounced in Sf-xl-val, Tokyo247, Nordland, and Sf-xl-testv1 datasets, where there is a rapid rise in accuracy. Ultimately, the highest level of accuracy is achieved at a dimensionality of 4096. This phenomenon suggests that utilizing image feature vectors with too low dimensionality may result in reduced robustness to variations such as changes in illumination and seasonal shifts in VPR tasks. Consequently, this study adopts a final image feature dimensionality of 4096.

3) *Backbone architecture*: DINOv2 encompasses four ViT models, with ViTg14 (giant) being the largest. Through

TABLE 3: *Table of Test Results of Different Methods on Datasets with Changes in Viewpoint, Illumination, Season. DINO-Mix(ViTb14) is ours, Bolded numbers are optimal results, and underlined numbers are sub-optimal results*

Method	Training data	Vector dim	Test dataset					
			Pitts250k	Pitts30k	SF-XLval	Tokyo24/7	Nordland	SF-XLTestv1
MAX [14]	GSV-cities	1024	46.45	55.87	43.68	3.17	9.30	8.10
AVG [14]	GSV-cities	1024	51.85	64.20	45.43	8.25	14.99	13.90
SPOC [39]	GSV-cities	256	60.59	68.37	56.11	19.68	10.18	20.50
MAC [40]	GSV-cities	256	61.75	69.42	56.47	18.73	13.49	22.00
RMAC [41]	GSV-cities	256	71.3	75.94	63.74	32.70	13.30	29.40
RRM [42]	GSV-cities	256	88.14	87.49	81.60	57.46	46.00	53.40
GeM [11]	GSV-cities	256	76.01	79.61	68.79	34.60	23.8	33.70
NetVLAD [14]	Pitts-30k	16384	86.93	86.36	65.34	53.97	7.86	42.50
NetVLAD [14]	GSV-cities	16384	89.71	88.04	80.38	70.79	36.25	58.90
CRN [33]	GSV-cities	16384	90.60	89.03	81.83	70.16	38.58	64.40
MultiRes-NetVLAD [35]	Pitts-30k	32768	86.70	86.80	–	69.80	–	–
SARE [36]	Pitts-30k	4096	88.00	87.20	–	74.80	–	45.5
SERS [37]	Pitts-30k	4096	90.40	89.10	–	80.30	16.00	50.30
CosPlace [24]	GSV-cities	4096	89.89	88.54	84.01	63.17	47.62	54.30
ConvAP [20]	GSV-cities	4096	91.52	89.67	76.95	72.06	63.93	59.20
MixVPR [21](Resnet18)	GSV-cities	4096	91.75	89.57	80.68	75.24	64.75	64.90
MixVPR [21](Resnet50)	GSV-cities	4096	<u>94.13</u>	<u>91.52</u>	<u>85.40</u>	<u>85.40</u>	<u>76.12</u>	<u>75.70</u>
<b>DINO-Mix(ViTb14)(Ours)</b>	GSV-cities	4096	<b>94.58</b>	<b>92.03</b>	<b>89.25</b>	<b>91.75</b>	<b>80.18</b>	<b>82.00</b>

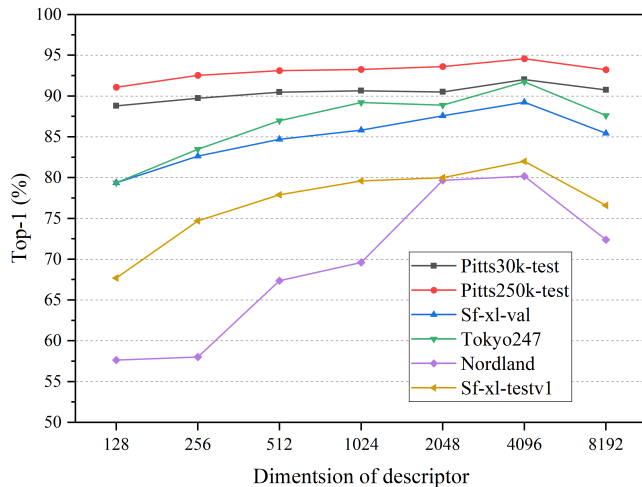


Fig. 6: *performance on pitts30k-test with different dimensionality configurations.*

model knowledge distillation, three smaller models were obtained from the distillation process, including ViTl14 (large), ViTb14 (base), and ViTs14 (small), as displayed in Tab.1. To evaluate the performance of these four models in the DINO-Mix framework, we conducted training on GSV-Cities as the training set and tested Pitts30k-test using ViTg14-Mix, ViTl14-Mix, ViTb14-Mix, and ViTs14-Mix. The feature mixer was fixed at two layers and the dimensionality of the image feature vectors was set to 4096. In addition, we trained and tested the DINO-Mix models with four different

backbone networks under six scenarios: updating the weights of the last one, two, three, six, and nine blocks, and not updating the weights of the backbone network (none). The results are shown in Fig.7.

From the perspective of the four differently sized backbone networks, ViTb14-Mix exhibited higher accuracy than the other three models, with a maximum Top-1 accuracy of 92.03%. In contrast, ViTg14-Mix exhibited the worst overall performance. This suggests that ViTg14’s large parameter count extracts deeper features from images, which adversely affects subsequent feature aggregation in the feature mixer.

Models without parameter updates for the backbone network demonstrated poorer performance. As the number of updated blocks increased, both ViTb14-Mix and ViTl14-Mix showed a gradual improvement in test accuracy, reaching their highest values after updating the parameters of the last three blocks, and stabilizing thereafter. In contrast, ViTs14-Mix achieved the highest test accuracy and stability after updating the parameters of the last two blocks. However, for ViTg14-Mix, the block parameter updates did not significantly enhance accuracy. Starting from the last three blocks, the ViTg14-Mix test accuracy showed a downward trend. This indicates that excessively deep block parameter updates may extensively alter the original pre-trained parameters.

In summary, updating the parameters of the last three blocks of the backbone network yielded optimal results. Considering the parameter counts of the four DINO-Mix models shown in Fig.8, we selected ViTb14-Mix, which has a moderate parameter count and superior test accuracy, as the final model for DINO-Mix.

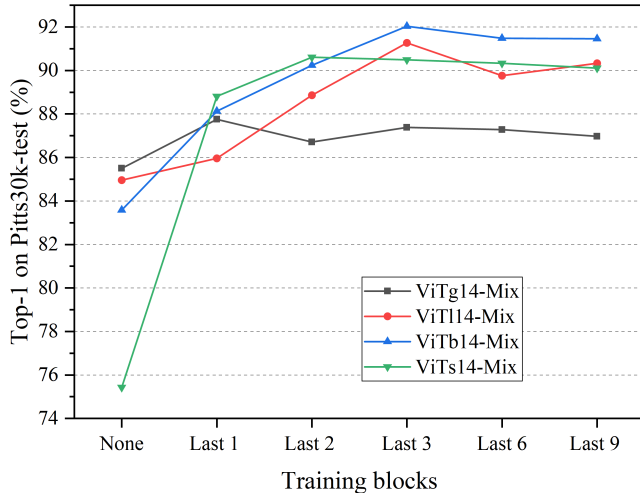


Fig. 7: *Test results of different DINO-Mix models with different weights for updating the number of layers.*

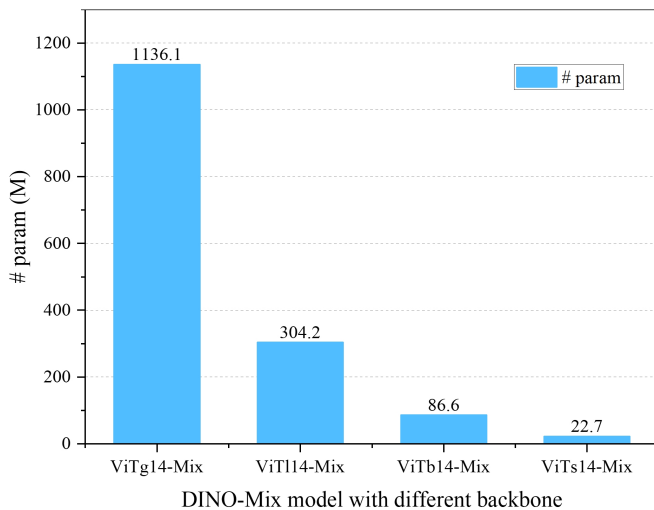


Fig. 8: *Parametric quantities of the four DINO-Mix models.*

## V. QUALITATIVE RESULTS

### A. Image retrieval comparison

In this study, we compared the performance of DINO-Mix with those of existing SOTA methods, including MixVPR, NetVLAD, ConvAP, and CosPlace, in image retrieval tasks. To demonstrate the robustness of DINO-Mix for VPR in complex environments, we selected several representative image retrieval cases from the Tokyo24/7, SF-XL-Testv1, and Nordland datasets. We presented four challenging scenarios: viewpoint changes, illumination changes, object occlusions, and seasonal variations. In cases where DINO-Mix succeeded, the other methods failed to accurately locate the query image, as displayed in Tab.4.

**Viewpoint Change:** Viewpoint changes encompass variations in the field angle and field range, posing challenges for image retrieval. Rows 1 and 2 in Tab.4 show examples of viewpoint changes in terms of field angle and field range.

Notably, only DINO-Mix resists interference caused by view-point changes and retrieves the correct image, whereas the other methods retrieve similar buildings or scenes.

**Illumination Change:** Illumination changes significantly affect image retrieval accuracy. Dim lighting conditions can blur textures in images, adversely affecting feature extraction and, consequently, image retrieval accuracy. Rows 3 and 4 in Tab.4 depict the image retrieval cases under dark conditions. Rows 5 and 6 present nighttime scenarios with artificial and natural light variations, respectively. DINO-Mix exhibits strong robustness against illumination changes, whereas the other methods suffer from the effects of lighting variations and fail to retrieve accurate results.

**Occlusion:** Image retrieval focuses primarily on objects, such as buildings, facilities, and natural landscapes. However, pedestrians, vehicles, and other objects can interfere with the semantic information in an image, posing challenges for image retrieval. As shown in rows 7 and 8 in Tab.4, where a large number of pedestrians are present in the query images, and in row 9, where the influence of buildings is significant, these occlusions pose significant difficulties for image retrieval. MixVPR retrieved the correct content but exceeded the threshold  $s$  (25  $m$ ) in the localization results. In contrast, DINO-Mix successfully extracted the correct features from the images and retrieved accurate results despite these challenges.

**Season Change:** The appearance characteristics of a location undergo significant changes in different seasons, such as heavy snowfall in winter (as illustrated in row 10 of Tab.4), and leaves falling from trees (row 11). These seasonal variations also have a profound impact on the image retrieval accuracy. Under such challenging circumstances, DINO-Mix overcame the drastic contrast caused by seasonal changes and achieved satisfactory results.

### B. Attention map visualization

To provide a more intuitive demonstration of the superiority of DINO-Mix over other VPR methods, we visualized their attention maps as presented in 5. The attention scores are represented by varying colors from blue to green to red, indicating low to high attention levels. Our analysis reveals that DINO-Mix can focus more on buildings, object contours, and textures, which are crucial factors for image retrieval. In contrast, it effectively excludes negative elements such as pedestrians, cars, and occlusions. This suggests that DINO-Mix has a greater ability to capture essential features and extract more robust image representations.

## VI. CONCLUSIONS














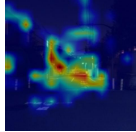
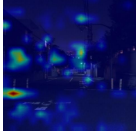
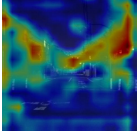
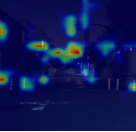
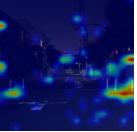

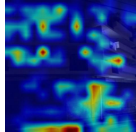
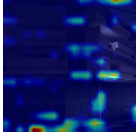
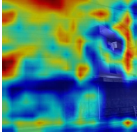
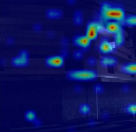
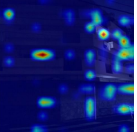


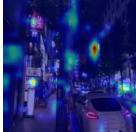
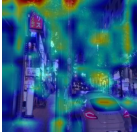



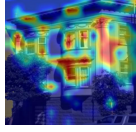

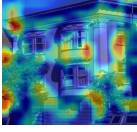







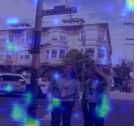





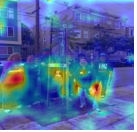

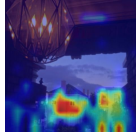
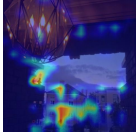
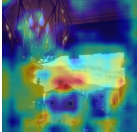
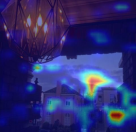
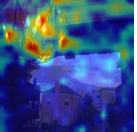

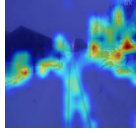
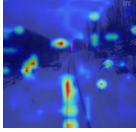
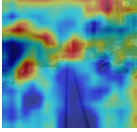

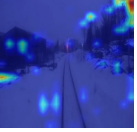


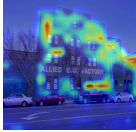


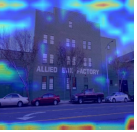
In this paper, we proposed a novel VPR framework, called DINO-Mix. First, we modified and fine-tuned the structure of the DINOv2 model. We then converted the extracted features from the backbone network into feature maps and employed the mixed feature aggregation module to aggregate these feature maps to obtain global feature vectors. The experimental results on various test sets demonstrate that the proposed DINO-Mix model outperforms SOTA methods



TABLE 4: Comparison of image retrieval results of DINO-Mix with other methods in difficult cases (Top-1). The green and red boxes in the table represent image retrieval success and failure, respectively, and the yellow box represents that the image content should be correct but the localization distance exceeds the threshold  $s$ .

Category	Query	DINO-Mix	MixVPR	NetVLAD	ConvAP	CosPlace
Viewpoint Change						
Illumination Change						
Occlusions						
Season Change						

TABLE 5: *The attention map visualization of the query images.*

Query	DINO-Mix	MixVPR	NetVLAD	ConvAP	CosPlace
					
					
					
					
					
					
					
					
					
					
					

in terms of VPR accuracy, with an average improvement of 5.14% across test sets containing challenging conditions. Furthermore, through a series of image retrieval examples

under difficult circumstances, we demonstrated that the performance of the DINO-Mix architecture significantly surpasses that of current SOTA architectures. We identified



that changes in illumination pose a significant challenge. Therefore, our future work will focus on enhancing images with poor lighting conditions through image augmentation techniques to further improve the VPR accuracy.

## REFERENCES

- [1] A. Dosovitskiy, L. Beyer, A. Kolesnikov, D. Weissenborn, X. Zhai, T. Unterthiner, M. Dehghani, M. Minderer, G. Heigold, S. Gelly, J. Uszkoreit, and N. Houlsby, “An image is worth 16x16 words: Transformers for image recognition at scale.”
- [2] A. R. Zamir, S. Ardeshtir, and M. Shah, “GPS-tag refinement using random walks with an adaptive damping factor,” in *2014 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, pp. 4280–4287. [Online]. Available: <https://ieeexplore.ieee.org/document/6909941>
- [3] D. Lowe, “Distinctive image features from scale-invariant keypoints,” vol. 60, no. 2, pp. 91–110. [Online]. Available: <https://api.semanticscholar.org/CorpusID:174065>
- [4] A. R. Zamir and M. Shah, “Accurate image localization based on google maps street view,” in *Computer Vision – ECCV 2010*, K. Daniilidis, P. Maragos, and N. Paragios, Eds., vol. 6314. Springer, pp. 255–268.
- [5] W. Zhang and J. Kosecka, “Image based localization in urban environments,” in *Third International Symposium on 3D Data Processing, Visualization, and Transmission, Proceedings*, M. Pollefeys and K. Daniilidis, Eds., pp. 33–40.
- [6] Martin A. Fischler and Robert C. Bolles, “Random sample consensus: a paradigm for model fitting with applications to image analysis and automated cartography,” vol. 24, no. 6, pp. 381–395. [Online]. Available: <http://dx.doi.org/10.1145/358669.358692>
- [7] L. Q. Yan, Y. M. Cui, Y. J. Chen, and D. F. Liu, “Hierarchical attention fusion for geo-localization,” in *2021 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP 2021)*. IEEE, pp. 2220–2224, WOS:000704288402094 abstractTranslation:. [Online]. Available: <https://www.webofscience.com/wos/alldb/summary/0c2e17a1-e132-41d5-892a-91e8febe8286-3c218449/relevance/1>
- [8] T. Y. Chu, Y. M. Chen, L. h. Huang, Z. G. Xu, and H. Y. Tan, “A grid feature-point selection method for large-scale street view image retrieval based on deep local features,” vol. 12, no. 23, p. 3978. [Online]. Available: <https://www.webofscience.com/wos/alldb/summary/0c2e17a1-e132-41d5-892a-91e8febe8286-3c218449/date-descending/1>
- [9] T. Y. Chu, Y. M. Chen, L. H. Huang, H. Y. Tan, J. P. Cao, and Z. Q. Xu, “Street view image retrieval with average pooling features,” in *IGARSS 2020 - 2020 IEEE International Geoscience and Remote Sensing Symposium*. IEEE, pp. 1205–1208. [Online]. Available: <https://ieeexplore.ieee.org/document/9323667/>
- [10] H. Noh, A. Araujo, J. Sim, T. Weyand, and B. Han, “Large-scale image retrieval with attentive deep local features,” in *2017 IEEE International Conference on Computer Vision (ICCV)*. IEEE, pp. 3476–3485. [Online]. Available: <http://ieeexplore.ieee.org/document/8237636/>
- [11] F. Radenovic, G. Tolias, and O. Chum, “Fine-tuning CNN image retrieval with no human annotation,” vol. 41, no. 7, pp. 1655–1668. [Online]. Available: <https://ieeexplore.ieee.org/document/8382272/>
- [12] D. Mishkin, M. Perdoch, and J. Matas, “Place recognition with WxBS retrieval,” in *CVPR 2015 Workshop on Visual Place Recognition in Changing Environments*, vol. 30, p. 9.
- [13] N. Dalal and B. Triggs, “Histograms of oriented gradients for human detection,” in *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR’05)*, vol. 1. IEEE, pp. 886–893. [Online]. Available: <http://ieeexplore.ieee.org/document/1467360/>
- [14] A. Relja, Petr Gronat, Akihiko Torii, Tomas Pajdla, and Josef Sivic, “NetVLAD: CNN architecture for weakly supervised place recognition,” vol. 40, no. 6, pp. 1437–1451. [Online]. Available: <https://ieeexplore.ieee.org/document/7937898/>
- [15] T. Y. Chu, Y. M. Chen, H. Su, Z. Z. Xu, G. D. Chen, and A. N. Zhou, “A news picture geo-localization pipeline based on deep learning and street view images,” vol. 15, no. 1, pp. 1485–1505. [Online]. Available: <https://www.tandfonline.com/doi/full/10.1080/17538947.2022.2121437>
- [16] A. Torii, J. Sivic, M. Okutomi, and T. Pajdla, “Visual place recognition with repetitive structures,” vol. 37, no. 11, pp. 2346–2359. [Online]. Available: <http://ieeexplore.ieee.org/document/7054472/>
- [17] H. Bay, T. Tuytelaars, and L. Van Gool, “SURF: Speeded up robust features,” in *Computer Vision – ECCV 2006*, A. Leonardis, H. Bischof, and A. Pinz, Eds. Springer Berlin Heidelberg, vol. 3951, pp. 404–417, series Title: Lecture Notes in Computer Science. [Online]. Available: [http://link.springer.com/10.1007/11744023\\_32](http://link.springer.com/10.1007/11744023_32)
- [18] E. Rublee, V. Rabaud, K. Konolige, and G. Bradski, “ORB: An efficient alternative to SIFT or SURF,” in *2011 International Conference on Computer Vision*. IEEE, pp. 2564–2571. [Online]. Available: <http://ieeexplore.ieee.org/document/6126544/>
- [19] R. T. Wang, Y. Q. Shen, W. L. Zuo, S. P. Zhou, and N. N. Zheng, “TransVPR: Transformer-based place recognition with multi-level attention aggregation,” in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, pp. 13 638–13 647. [Online]. Available: <https://ieeexplore.ieee.org/document/9879296/>
- [20] A. Ali-bey, B. Chaib-draa, and P. Giguère, “GSV-cities: Toward appropriate supervised visual place recognition,” vol. 513, pp. 194–203. [Online]. Available: <https://linkinghub.elsevier.com/retrieve/pii/S0925231222012188>
- [21] A. Ali-Bey, B. Chaib-Draa, and P. Giguere, “MixVPR: Feature mixing for visual place recognition,” in *2023 IEEE/CVF Winter Conference on Applications of Computer Vision (WACV)*. IEEE, pp. 2997–3006. [Online]. Available: <https://ieeexplore.ieee.org/document/10030191/>
- [22] T. Ng, V. Balntas, Y. Tian, and K. Mikolajczyk, “SOLAR: Second-order loss and attention for image retrieval,” in *Computer Vision – ECCV 2020: 16th European Conference Part XXV 16*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds. Springer International Publishing, pp. 253–270, abstractTranslation:.
- [23] G. Tolias, T. Jenicek, and O. Chum, “Learning and aggregating deep local descriptors for instance-level recognition,” in *Computer Vision – ECCV 2020*, A. Vedaldi, H. Bischof, T. Brox, and J.-M. Frahm, Eds. Springer International Publishing, pp. 460–477.
- [24] G. Berton, C. Masone, and B. Caputo, “Rethinking visual geo-localization for large-scale applications,” in *2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR 2022)*, ser. IEEE Conference on Computer Vision and Pattern Recognition. IEEE Computer Soc, pp. 4868–4878, abstractTranslation: Backup Publisher: IEEE; CVF; IEEE Comp Soc ISSN: 1063-6919 Type: Proceedings Paper abstractTranslation:.
- [25] A. Torii, R. Arandjelovic, J. Sivic, M. Okutomi, and T. Pajdla, “24/7 place recognition by view synthesis,” vol. 40, no. 2, pp. 257–271, place: 10662 LOS VAQUEROS CIRCLE, PO BOX 3014, LOS ALAMITOS, CA 90720-1314 USA Publisher: IEEE COMPUTER SOC Type: Article.
- [26] M. Oquab, T. Darcet, T. Moutakanni, H. Vo, M. Szafraniec, V. Khalidov, P. Fernandez, D. Haziza, F. Massa, A. El-Nouby, M. Assran, N. Ballas, W. Galuba, R. Howes, P.-Y. Huang, S.-W. Li, I. Misra, M. Rabbat, V. Sharma, G. Synnaeve, H. Xu, H. Jegou, J. Mairal, P. Labatut, A. Joulin, and P. Bojanowski, “DINOv2: Learning robust visual features without supervision.” [Online]. Available: <http://arxiv.org/abs/2304.07193>
- [27] I. O. Tolstikhin, N. Houlsby, A. Kolesnikov, L. Beyer, X. Zhai, T. Unterthiner, J. Yung, A. Steiner, D. Keysers, J. Uszkoreit, M. Lucic, and A. Dosovitskiy, “MLP-mixer: An all-MLP architecture for vision,” in *Advances in Neural Information Processing Systems*, M. Ranzato, A. Beygelzimer, Y. Dauphin, P. S. Liang, and J. W. Vaughan, Eds., vol. 34. Curran Associates, Inc., pp. 24 261–24 272. [Online]. Available: [https://proceedings.neurips.cc/paper\\_files/paper/2021/file/cba0a4ee5ccd02fda0fe3f9a3e7b89fe-Paper.pdf](https://proceedings.neurips.cc/paper_files/paper/2021/file/cba0a4ee5ccd02fda0fe3f9a3e7b89fe-Paper.pdf)
- [28] G. S. Huang, Y. Zhou, X. F. Hu, L. Y. Zhao, and C. L. Zhang, “A survey of the research progress in image geo-localization,” vol. 25, no. 7, pp. 1336–1362, abstractTranslation:.
- [29] C. Masone and B. Caputo, “A survey on deep visual place recognition,” vol. 9, pp. 19 516–19 547. [Online]. Available: <https://ieeexplore.ieee.org/document/9336674/>
- [30] K. Tang, F. F. Li, and D. Koller, “Learning latent temporal structure for complex event detection,” in *2012 IEEE Conference on Computer Vision and Pattern Recognition*. IEEE, pp. 1250–1257. [Online]. Available: <http://ieeexplore.ieee.org/document/6247808/>
- [31] H. Jegou, M. Douze, C. Schmid, and P. Perez, “Aggregating local descriptors into a compact image representation,” in *2010 IEEE Computer Society Conference on Computer Vision and Pattern Recognition*. IEEE, pp. 3304–3311. [Online]. Available: <http://ieeexplore.ieee.org/document/5540039/>
- [32] H. Jegou, F. Perronnin, M. Douze, J. Sanchez, P. Perez, and C. Schmid, “Aggregating local image descriptors into compact

- codes,” vol. 34, no. 9, pp. 1704–1716. [Online]. Available: <http://ieeexplore.ieee.org/document/6104058/>
- [33] H. J. Kim, E. Dunn, and J.-M. Frahm, “Learned contextual feature reweighting for image geo-localization,” in *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, pp. 3251–3260. [Online]. Available: <http://ieeexplore.ieee.org/document/8099829/>
- [34] J. Yu, C. Y. Zhu, J. Zhang, Q. M. Huang, and D. C. Tao, “Spatial pyramid-enhanced NetVLAD with weighted triplet loss for place recognition,” vol. 31, no. 2, pp. 661–674. [Online]. Available: <https://ieeexplore.ieee.org/document/8700608/>
- [35] A. Khaliq, M. Milford, and S. Garg, “MultiRes-NetVLAD: Augmenting place recognition training with low-resolution imagery,” vol. 7, no. 2, pp. 3882–3889. [Online]. Available: <http://arxiv.org/abs/2202.09146>
- [36] L. Liu, H. D. Li, and Y. C. Dai, “Stochastic attraction-repulsion embedding for large scale image localization,” in *2019 IEEE/CVF International Conference on Computer Vision (ICCV)*. IEEE, pp. 2570–2579. [Online]. Available: <https://ieeexplore.ieee.org/document/9010658/>
- [37] Y. x. Ge, H. b. Wang, F. Zhu, R. Zhao, and H. s. Li, “Self-supervising fine-grained region similarities for large-scale image localization,” vol. 12349. Springer International Publishing, pp. 369–386.
- [38] A. Kirillov, E. Mintun, N. Ravi, H. Mao, C. Rolland, L. Gustafson, T. Xiao, S. Whitehead, A. C. Berg, W.-Y. Lo, P. Dollár, and R. Girshick, “Segment anything.” [Online]. Available: <http://arxiv.org/abs/2304.02643>
- [39] A. B. Yandex and V. Lempitsky, “Aggregating deep convolutional features for image retrieval,” in *2015 IEEE International Conference on Computer Vision (ICCV)*. IEEE, pp. 1269–1277. [Online]. Available: <http://ieeexplore.ieee.org/document/7410507/>
- [40] A. S. Razavian, J. Sullivan, S. Carlsson, and A. Maki, “Visual instance retrieval with deep convolutional networks,” vol. 4, no. 3, pp. 251–258.
- [41] G. Toliás, R. Sivic, and H. Jégou, “Particular object retrieval with integral max-pooling of CNN activations.” [Online]. Available: <http://arxiv.org/abs/1511.05879>
- [42] G. Kordopatis-Zilos, P. Galopoulos, S. Papadopoulos, and I. Kompatziaris, “Leveraging EfficientNet and contrastive learning for accurate global-scale location estimation,” abstractTranslation:. [Online]. Available: <https://api.semanticscholar.org/CorpusID:234742169>
- [43] S. Ruder, “An overview of gradient descent optimization algorithms,” abstractTranslation: abstractTranslation:. [Online]. Available: <http://arxiv.org/abs/1609.04747>
- [44] A. Hermans, L. Beyer, and B. Leibe, “In defense of the triplet loss for person re-identification,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. IEEE, abstractTranslation:. [Online]. Available: <https://api.semanticscholar.org/CorpusID:52292078>
- [45] N. Sünderhauf, P. Neubert, and P. Protzel, “Are we there yet? challenging SeqSLAM on a 3000 km journey across all four seasons,” in *Proc. of Workshop on Long-Term Autonomy, IEEE International Conference on Robotics and Automation (ICRA)*, p. 2013.
- [46] X. Wang, X. Han, W. Huang, D. Dong, and M. R. Scott, “Multi-similarity loss with general pair weighting for deep metric learning,” in *2019 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*. IEEE, pp. 5017–5025. [Online]. Available: <https://ieeexplore.ieee.org/document/8954016/>