CLASS-RELATIONAL LABEL SMOOTHING FOR LIFELONG VISUAL PLACE RECOGNITION

Anonymous authors

004

010 011

012

013

014

015

016

017

018

019

021

025

026

027 028 029

030

Paper under double-blind review

ABSTRACT

Visual Place Recognition (VPR) is a task of estimating the location of a query image, predominantly executed through image retrieval using learned global descriptors from a reference database of geo-tagged images. While recent approaches have aimed to improve the scalability of VPR training by leveraging classification loss as a proxy task, this leads to a task gap between classification and retrieval – classification discretizes the feature space into distinct class regions, often overlooking visual differences between classes. This gap makes VPR systems particularly vulnerable to extreme visual changes such as lifelong variations. To remedy these problems, we propose a novel Class-Relational Label Smoothing (CRLS) that transforms one-hot labels into soft labels by considering visual information of inter-class relations. We further enhance this method by dynamically adjusting the influence of CRLS based on the stability of class weights, which is quantified by their magnitudes. Importantly, our findings suggest that the magnitude of class weights serves as an indicator of class stability, which is also supported by derivative analysis. We demonstrate that our method outperforms state-of-the-art methods on the most extensive 17 benchmarks, effectively bridging the task gap between classification and retrieval in visual place recognition. Code and trained weights will be made publicly available.

1 INTRODUCTION

031 Visual Place Recognition (VPR) is an important task in various applications, such as robotics (Stumm 032 et al., 2013), autonomous driving (Bresson et al., 2017) and navigation (Mirowski et al., 2018), where 033 its goal is to identify a location based on visual data. Typically, VPR is approached as an image 034 retrieval problem (Wang et al., 2022; Zhu et al., 2023; Shen et al., 2023; Leyva-Vallina et al., 2023; Hausler et al., 2021; Arandjelovic et al., 2016; Warburg et al., 2020; Torii et al., 2013a; Thoma et al., 2020; Keetha et al., 2023; Lu et al.; 2024; Izquierdo & Civera, 2024), employing nearest neighbor 037 search based on the similarity of descriptors between a query image and gallery images. This process 038 enables the identification of the gallery image most similar to the query image, and subsequently, localization is achieved using the geo-reference of the identified image. 039

040 Real-world applications of VPR face numerous challenges due to significant appearance changes in 041 various environments. These changes include well-known seasonal variations (Naseer et al., 2018; 042 Sünderhauf et al., 2013), weather conditions (Ros et al., 2016; Berton et al., 2021), illumination 043 changes (day/night) (Sattler et al., 2012; Maddern et al., 2017) and viewpoint changes (Carlevaris-044 Bianco et al., 2016; Berton et al., 2023). Subsequently, several benchmark works (Warburg et al., 2020; Ali-bey et al., 2022; Berton et al., 2022) have introduced lifelong datasets that include images collected over a long temporal span. We find that changes over extended periods, such as building 046 modifications and remodeling, present extreme challenges for VPR, which we refer to as *lifelong* 047 variations. For instance, Fig. 1a shows how two buildings evolve over time, with one highlighted in 048 green and the other in orange. These buildings, located in San Francisco, undergo significant changes due to frequent urban modifications, often requiring continuous updates to the VPR models. To achieve effective retrieval for VPR, the trained model should continuously capture visual differences, 051 leading to a continuous representation space that can deal with lifelong variations. 052

To learn such representation spaces, existing VPR methods (Arandjelovic et al., 2016; Berton et al., 2021; Hausler et al., 2021; Peng et al., 2021; Zhu et al., 2023; Leyva-Vallina et al., 2023; Wang



Figure 1: An example of lifelong visual place recognition from SF-XL test v1 (Berton et al., 2022).
(a) illustrates the evolution of buildings over time, using images of the same location from Google StreetView and Flickr, linked by provided GPS data. Lifelong scenarios often show significant changes due to frequent urban modifications. (b) and (c) compare retrieval results of our method and a state-of-the-art method, respectively. Positive images are highlighted in green, while negative images are highlighted in red.

054

056

059

060

061

068

076 et al., 2022) predominantly utilize metric learning losses such as contrastive or triplet loss, which 077 rely heavily on mining negative examples throughout the training database (Arandjelovic et al., 2018). This mining process is notably resource-intensive and becomes prohibitively expensive as the size of the dataset increases. To address this scalability issue, more recent methods (Berton et al., 079 2022; 2023) have adopted a classification loss (*i.e.*, CosFace (Wang et al., 2018)) as a proxy task, thereby streamlining training by eliminating the need for the exhaustive negative mining. However, 081 while this improves scalability, it introduces a task gap between classification and retrieval tasks. Specifically, classification tends to discretize the representation space into distinct class regions, 083 leading to overconfident predictions (Guo et al., 2017). This contrasts with retrieval tasks, which 084 require a continuous representation of subtle visual differences (Leyva-Vallina et al., 2021; Kim et al., 085 2021; 2022) to capture gradual and continuous changes in environments. A naïve approach to mitigate the issue can be label smoothing (Szegedy et al., 2016), which prevents overconfidence by bringing 087 all classes closer together in the representation space (Müller et al., 2019). While label smoothing can 088 partially relax the discretization of representation space, it assigns probabilities uniformly across all classes, disregarding the visual differences between them, making it less effective for retrieval tasks.

090 In this paper, we propose Class-Relational Label Smoothing (CRLS) for lifelong VPR to better address 091 the task gap between classification and retrieval by incorporating visual similarities between classes 092 into the label smoothing process, as illustrated in Fig. 2. Unlike conventional label smoothing, which 093 treats all classes uniformly, we adjust the label distribution based on inter-class relationships, resulting 094 in a continuous representation space where visually similar classes are closer, better reflecting subtle visual differences. However, the impact of CRLS may be influenced by fluctuations in class weights. 095 To address this issue, we further introduce Class Stability Weighting (CSW), which dynamically 096 adjusts the impact of CRLS based on the stability of class weights. Specifically, we empirically 097 observe that class weight magnitudes reflect the classification difficulty of each class, as further 098 supported by derivative analysis. Consequently, our method enables the model to learn a more robust and continuous representation that captures gradual visual differences over time, making it particularly 100 effective in handling lifelong variations. Extensive experiments on 17 benchmarks demonstrate that 101 our approach effectively bridges the task gap between classification and retrieval, achieving superior 102 performance compared to state-of-the-art methods, especially under lifelong variations. 103

Our contributions are summarized as follows:

104 105

To address the task gap and challenges posed by lifelong variations, we introduce Class-Relational Label Smoothing (CRLS) that transforms one-hot hard labels into soft labels while considering visual similarities between classes.



Figure 2: Conceptual comparison between previous and our works. (a) The target class y shows a street scene with buildings, which is visually similar to class j_1 , while class j_2 shows a markedly different scene inside a tunnel. (b) Hard label strictly assigns the probability to the target class y, treating the visually similar class j_1 as a negative class. (c) Label smoothing (Szegedy et al., 2016) distributes probability evenly across all classes, including the visually dissimilar class j_2 as a positive class. (d) Our CRLS smoothes the label distribution being aware of visual relationships, thus treating class j_1 as a as pseudo-positive class and class j_2 as a pseudo-negative class.

- Building on the insights from the behavior of class weights during training, we propose Class Stability Weighting (CSW), which dynamically adjust the label smoothing process according to the stability of class weights.
- We organize the lifelong category based on the given temporal span and conduct comprehensive experiments on a diverse set of 17 benchmarks. We achieve state-of-the-art performance across the majority of benchmarks, showing notable improvements on lifelong benchmarks.
- 2 RELATED WORK
- 136 137

129

130 131

132

133

134 135

138 Visual place recognition. The task of Visual Place Recognition (VPR) has been approached 139 through image retrieval techniques, where the location of a query image is determined by matching 140 to geo-tagged images within an extensive database. Traditionally, this involved the aggregation 141 of hand-crafted features, such as SIFT (Lowe, 2004) and SURF (Bay et al., 2008), for a task of 142 VPR (Knopp et al., 2010; Gronat et al., 2013; Torii et al., 2013a; 2015). With the advent of deep learning, many recent works (Wang et al., 2022; Zhu et al., 2023; Shen et al., 2023; Leyva-Vallina 143 et al., 2023; Hausler et al., 2021; Arandjelovic et al., 2016; Ge et al., 2020; Warburg et al., 2020; Wang 144 et al., 2023; Lu et al.) have exploited metric learning losses, such as contrastive and triplet losses, 145 to get more discriminative image representations. StructVPR (Shen et al., 2023) enhances RGB 146 global features with structural knowledge through segmentation masks and knowledge distillation. 147 R^2 Former (Zhu et al., 2023) integrates retrieval and reranking with a transformer-based framework, 148 offering a more computationally efficient alternative to other RANSAC-based methods (Hausler 149 et al., 2021; Wang et al., 2022) in the reranking stage. GCL (Leyva-Vallina et al., 2023) extends the 150 traditional contrastive loss in a generalized manner and constructs soft labels based on view overlap, 151 enabling the consideration of continuous relations of viewpoints between images while training VPR 152 models. The methods introduced above are primarily built upon ranking losses, with the use of a 153 mining strategy such as hard negative mining.

154 More recently, CosPlace (Berton et al., 2022) newly proposes the extensive San Francisco eXtra 155 Large (SF-XL) dataset, showing the scalability issues of the mining strategy in existing VPR methods. 156 To overcome this problem, it leverages a classification loss as a proxy task, specifically Large Margin 157 Cosine Loss (LMCL) (Wang et al., 2018), achieving the superior performances on VPR benchmarks. 158 Furthermore, EigenPlaces (Berton et al., 2023) revises the labeling strategy of the SF-XL dataset to 159 incorporate classes with diverse views, effectively addressing a challenge of substantial viewpoint changes in VPR methods. However, among the VPR methods under the study that utilize classification 160 loss as a proxy task for training, there has been no exploration into the effectiveness of using label 161 smoothing or incorporating class similarity to address the task gap.

162 Label smoothing and its retrieval application. Label smoothing (Szegedy et al., 2016) is first 163 proposed to address the issue of overconfidence in model predictions, thereby enhancing the model's 164 generalization ability with smoother decision boundaries. Since its introduction, it has been widely 165 used for network calibration (Wang, 2023) from many works (Pereyra et al., 2017; Müller et al., 2019; 166 Liu et al., 2022; Park et al., 2023; Liu et al., 2023). MbLS (Liu et al., 2022) applies label smoothing selectively, using the distance over logits and a controllable margin for flexible generalization. 167 Building upon the work, ACLS (Park et al., 2023) further adjusts the label smoothing intensity 168 according to the logit distances. CALS (Liu et al., 2023) is another calibration technique that calculates class-wise penalty weights from the loss using the augmented Lagrangian multiplier 170 method. It adjusts label smoothing intensity based on these penalty weights, enabling class-wise 171 calibration. However, these methods were not proposed for image retrieval or VPR. 172

For a retrieval application of person re-identification (re-ID), Luo et al. (2019) proposed utilizing the 173 label smoothing as a recipe for establishing a strong baseline. It has since been adopted and modified 174 by several re-ID works to improve search accuracy (Zhu et al., 2020; Cho et al., 2022; Jia et al., 175 2022). Zhu et al. (2020) introduced an instance-wise adaptive label smoothing that adjusts smoothing 176 strengths based on network predictions and viewpoint variability within each identity class. Another 177 instance-wise label smoothing was proposed by Cho et al. (2022), which applies different smoothing 178 levels to local features considering the relationship between global and local features. However, the 179 existing works including calibration have not considered visual similarity-based class relations to 180 make representation space continuous over visual changes. 181

3 **METHODS**

3.1 PRELIMINARIES

Recent works (Berton et al., 2022; 2023) in VPR utilize a classification loss as a proxy task, employing CosFace (Wang et al., 2018) as the training loss. A margin-based logit function for any class jincluding a target class y can be formulated with a margin m and a scale factor s as:

$$l(\cos\theta_j) = \begin{cases} s(\cos\theta_j - m) & j = y\\ s\cos\theta_j & j \neq y \end{cases}.$$
 (1)

Predicted probability p_i are calculated using cosine similarity and softmax, which is specifically formulated as follows:

$$p_j = \frac{\exp\left(l\left(\cos\theta_j\right)\right)}{\sum\limits_k \exp\left(l\left(\cos\theta_k\right)\right)}, \ \cos\theta_k = \frac{W_k^T \cdot x}{||W_k|| \cdot ||x||},\tag{2}$$

where W_k is a weight vector of class k and x is a feature vector of an input image. These probabilities are then utilized to compute cross-entropy loss for CosFace loss as follows:

182

183

185 186

187

188

189 190 191

192

193 194

195 196

197

198

199 200 $\mathcal{L}_{\text{CosFace}}\left(q,p\right) = H\left(q,p\right) = -\sum_{k} q_k \log\left(p_k\right),$ (3)

where hard target distribution q_k is 1 for the target class and 0 for the rest. Subsequently, the model 201 trained on the classification loss is deployed for a task of image retrieval, specifically in the context 202 of VPR. 203

204 Label Smoothing (LS) (Szegedy et al., 2016) is a regularization technique that adjusts the hard target, 205 typically represented as one-hot encoded vectors, towards a smoother distribution. This adjustment 206 encourages the model to be less confident, thereby improving its generalization capabilities. The LS technique modifies the target distribution q according to the smoothing parameter α and the number 207 of classes K, as shown in the following equations: 208

209
210
211
212

$$f_{LS}(q_j) = \begin{cases} (1-\alpha), \quad j = y, \\ \frac{\alpha}{(K-1)}, \quad j \neq y, \end{cases}$$
(4)

$$\mathcal{L}_{LS}(q, p) = H(f_{LS}(q), p).$$

$$\mathcal{L}_{\mathsf{LS}}(q,p) = H(f_{\mathsf{LS}}(q),p)$$

213 In this work, we argue that assigning equal smoothing across all classes may be less effective for image retrieval task. For the task, establishing a continuous representation space that reflects 214 visual differences is crucial for extreme changes such as lifelong variations, while a discretized 215 representation space is usually sufficient for classification tasks.



Figure 3: Relationship between training accuracy and magnitude of class weight. (a)-(d) show histograms depicting this relationship at different training epochs. The x-axis represents the magnitude of the class weight, defined as the l_2 -norm of the weight vector, while the y-axis shows mean classification accuracy for classes within the specific magnitude range. It appears that relatively lower magnitudes consistently correspond to higher training accuracy throughout the overall training.

3.2 CLASS-RELATIONAL LABEL SMOOTHING

223

224

225

226

227

228

229

237

243 244 245

251

252 253

254

267

To better reflect the visual differences more within the representation space, we here introduce Class-Relational Label Smoothing (CRLS) extending LS by integrating visual information from inter-class relations. Thanks to the cosine-based classification loss demonstrated in Eq. 2 and 1, the l_2 -normed class weight vectors are aligned within the same distance space as l_2 -normed feature vectors, thereby enabling the calculation of visual similarity across classes. Hence, to capture the visual difference between classes, the similarity-based affinity between classes is calculated using the normalized dot product of their class weights $W \in \mathbb{R}^{K \times C}$ of dimensionality C, represented as:

$$A_{y,j} = \sigma(W_y)^{\top} \sigma(W_j), \tag{5}$$

where the l_2 -normalization function $\sigma(\mathbf{f}) = \mathbf{f}/||\mathbf{f}||$, and W_y and W_j are the weights of the classifier for the target class and any other classes, respectively. We then apply a softmax function to the class relations to obtain a label distribution. Additionally, we leveraged a temperature parameter τ to enhance a contrast in visual similarities among classes, as follows:

$$\hat{A}_{y,j} = \frac{\exp(A_{y,j}/\tau)}{\sum_{k \neq y} \exp(A_{y,k}/\tau)}, \quad \forall j \in \{1, \dots, K\}.$$
(6)

Finally, our loss using CRLS is reformulated with the constructed label distribution A_y as follows:

$$f_{\text{CRLS}}(q_j) = \begin{cases} (1-\alpha), & j=y, \\ \alpha \hat{A}_{y,j}, & j\neq y, \end{cases}$$
(7)

$$\mathcal{L}_{\text{CRLS}}(q, p) = H(f_{\text{CRLS}}(q), p).$$

We assign labels in a sequence that mirrors the visual similarity to the target class, thereby enhancing the model's ability to generalize across visually similar scenarios, *e.g.*, lifelong variation.

3.3 Loss Integration with Class Stability Weighting

255 Our CRLS approach builds a target distribution $f_{CRLS}(q)$ using class weights changing during training, 256 which may causes fluctuations in the target distribution during training. These rapid fluctuations can 257 destabilize the learning process. Classes that are more difficult to learn tend to require more updates for their weights during training, which can lead to larger fluctuations. These frequent updates 258 may result in higher weight magnitudes for such classes. Therefore, we investigate a relationship 259 between training accuracy and magnitude of class weight. Figure 3 empirically shows an inversely 260 proportional relationship throughout the entire training epochs; specifically, as the magnitude of 261 class weights increases, we observe a decline in training accuracy. A class achieving high training 262 accuracy tends to exhibit less fluctuation in its weight during training, thus giving a chance for the 263 stable application of CRLS. We therefore define the magnitude as a measure of class stability, and 264 use it for final loss function. We integrate two losses of LS and CRLS to make the training more 265 stable through Class Stability Weighting (CSW) as: 266

$$\mathcal{L} = \gamma_{y_i} \mathcal{L}_{\text{LS}} + (1 - \gamma_{y_i}) \mathcal{L}_{\text{CRLS}},\tag{8}$$

where γ_{y_i} is the loss weight for input image *i*, determined for each target class y_i by min-max normalization of the magnitudes over all classes. Integrating CSW with CRLS enables dynamic adjustment of each loss's contribution, effectively mitigating potential fluctuations during training.

070						
212	Dataset	#Database	#Queries	Lifelong	Multi-View	Single-View
273	SF-XL test v1 (Berton et al., 2022)	2.8M	1000	 ✓	✓	
274	SF-XL test v2 (Berton et al., 2022)	2.8M	598	1	\checkmark	
275	MSLS Val (Warburg et al., 2020)	18.9k	740	√		\checkmark
210	MSLS Challenge (Warburg et al., 2020)	38.7k	27k	 ✓ 		\checkmark
276	AmsterTime (Yildiz et al., 2022)	1231	1231	 ✓ 	\checkmark	
277	Eynsham (Cummins & Newman, 2009)	23.9k	23.9k		\checkmark	
278	Pitts30k (Torii et al., 2013b)	6.8k	10k		\checkmark	
270	Pitts250k (Torii et al., 2013b)	8.3k	84k		\checkmark	
279	Tokyo 24/7 (Torii et al., 2015)	76k	315		\checkmark	
280	San Francisco Landmark (Chen et al., 2011)	1M	598		\checkmark	
001	SVOX Night (Berton et al., 2021)	17k	823			\checkmark
201	SVOX Overcast (Berton et al., 2021)	17k	872			\checkmark
282	SVOX Rain (Berton et al., 2021)	17k	937			\checkmark
283	SVOX Snow (Berton et al., 2021)	17k	870			\checkmark
004	SVOX Sun (Berton et al., 2021)	17k	854			\checkmark
204	St Lucia (Milford & Wyeth, 2008)	1549	1464			\checkmark
285	Nordland (Sünderhauf et al., 2013)	27.5k	27.5k			\checkmark

Table 1: **Overview of VPR benchmarks.** This table provides dataset statistics, including the number of images in the database and queries, and their categorization as lifelong, multi-view, or single-view.

Gradient analysis on class weight. We assume that weight magnitude closely reflects the cumulative gradient magnitude over training epochs. By analyzing gradient magnitude, we can understand the relationship between training accuracy and magnitude of class weight. The training accuracy is associated with a similarity to target class, and the similarity reflects the difficulty of a given instance (Meng et al., 2021). The magnitude of the derivative of the loss function with respect to the class weight W_{y_i} , where y_i is the target class of an input *i*, can be simplified with respect to $\cos \theta_{y_i}$ as follows¹:

$$\left\|\frac{\partial L_{CE}}{\partial W_{y_i}}\right\| = \frac{\partial L_{CE}}{\partial l(\cos \theta_{y_i})} \frac{\partial l(\cos \theta_{y_i})}{\partial \cos \theta_{y_i}} \left\|\frac{\partial \cos \theta_{y_i}}{\partial W_{y_i}}\right\| \\ \propto (1 - p_{y_i}) \sqrt{1 - \cos^2 \theta_{y_i}} \quad \text{(with respect to } \cos \theta_{y_i}\text{)}.$$
(9)

The derived equation consists of the product of two terms: $(1 - p_{y_i})$ and $\sqrt{1 - \cos^2 \theta_{y_i}}$. Both of these terms are inversely related to $\cos \theta_{y_i}$, which means $\cos \theta_{y_i}$ controls the emphasis on gradients based on the difficulty during training. The change in class weight while updating tends to depend on the magnitude of the gradient. Consequently, in CSW, the weight magnitude helps to identify saturated classes that require minimal change, which aids in stabilizing the training process.

4 EXPERIMENTS

295 296 297

298

299

300

301

302 303

304 305

306

4.1 DATASETS AND EVALUATION

307 We extensively evaluate our method on diverse benchmarks outlined in EigenPlaces (Berton et al., 2023) for comprehensive and fair comparisons. These benchmarks are initially categorized into 308 Multi-View, featuring image pairs that include both frontal and lateral perspectives of the road, and 309 Single-View, where image pairs consist solely of frontal views relative to the road. Additionally, we 310 introduce a new category termed *Lifelong*, where the benchmark contains datasets with significant 311 temporal variations, capturing changes over extended periods. Among datasets providing temporal 312 span information, we classify those with at least a two-year gap between query and database images 313 as lifelong, allowing for long-term changes in the captured environments. Detailed statistics and the 314 category of each benchmark are shown in Table 1. For fairness and reproducibility, we conduct all 315 experiments using publicly available repositories^{2,3} for VPR evaluation. 316

We summarize the details of five benchmarks in the lifelong category as follows:

SF-XL test v1 and v2 (Berton et al., 2022). The SF-XL database, shared for evaluations in both test v1 and test v2, encompasses the entire city of San Francisco with 2.8M testing images captured in 2013. The images are sourced from Google StreetView and provide a wide range of challenging

¹The detailed derivation can be found in the Appendix.

²https://github.com/gmberton/VPR-methods-evaluation

³https://github.com/gmberton/VPR-datasets-downloader

Method	Backbone	Dim.	SF- test	-XL t v1	SF- test	-XL t v2	MS V	SLS al	MS Chal	SLS lenge	Am	ister.
			R@1	R@5	R@1	R@5	R@1	R@5	R@1	R@5	R@1	R@5
TransVPR	-	256	12.0	21.0	31.8	47.0	70.9	85.0	49.0	68.7	10.6	21.1
StructVPR	MobileNetV2	448	-	-	-	-	83.0	91.0	64.5	80.4	-	-
R^2 Former	ViT-S	256	20.3	31.9	37.1	60.4	80.8	90.9	56.6	75.8	12.9	26.8
GCL	VGG-16	512	9.6	15.3	37.3	54.0	64.7	77.0	42.9	55.9	10.2	22.7
CosPlace	VGG-16	512	65.9	75.3	83.1	91.3	82.4	<u>90.4</u>	61.2	73.8	<u>38.7</u>	<u>61.3</u>
EigenPlaces	VGG-16	512	<u>69.4</u>	<u>78.4</u>	<u>86.3</u>	<u>93.6</u>	<u>84.6</u>	90.3	60.9	72.7	38.0	59.2
Ours	VGG-16	512	73.5	80.4	87.0	94.1	85.3	91.5	63.4	75.1	39.2	61.8
R^2 Former	ResNet-50	256	19.0	30.4	47.0	65.9	79.6	90.7	57.0	74.1	16.7	31.4
MixVPR	ResNet-50	512	61.2	72.0	85.6	91.8	83.5	92.0	60.0	73.6	35.7	53.0
MixVPR	ResNet-50	4096	72.5	79.3	88.6	94.5	88.4	93.5	64.3	76.5	40.8	58.9
GCL	ResNet-50	2048	11.4	20.6	42.3	57.0	66.2	77.8	43.2	59.7	14.6	30.1
CosPlace	ResNet-50	2048	76.4	83.3	88.8	95.0	87.3	<u>94.0</u>	67.5	<u>77.9</u>	47.7	<u>69.8</u>
EigenPlaces	ResNet-50	2048	84.1	<u>89.1</u>	<u>90.8</u>	<u>95.7</u>	89.1	93.8	<u>67.9</u>	77.7	<u>48.9</u>	69.5
Ours	ResNet-50	2048	86.0	90.4	92.3	96.0	90.1	94.1	68.8	78.9	51.1	72.5

Table 2: **Comparison on the lifelong category.** This table categorizes methods by backbone type into three sections: Others, VGG, and ResNet. The best results are highlighted in bold, and the second best results are underlined, excluding the Others. Recall@1 and Recall@5 (%) are reported.

scenarios such as significant viewpoint changes, with highly accurate GPS data. SF-XL test v1
consists of 1,000 query images sourced from Flickr, captured across various years from 2006 to
2020, leading to a temporal gap of up to 14 years. SF-XL test v2 uses a set of 592 queries from from
the San Francisco Landmark dataset (Chen et al., 2011), which was released in 2011, ensuring a
minimal temporal gap of two years since the SF-XL database images were captured in 2013. Given
the frequent modification of urban structures (*e.g.*, buildings), SF-XL provides a realistic test for
lifelong scenarios with significant temporal variability.

MSLS Val and Challenge (Warburg et al., 2020) is a crowdsourced dataset for lifelong visual place
 recognition, containing tens of thousands of images from 30 major cities across six continents. For
 the challenge set, GPS data is withheld to ensure the integrity of evaluations, which are conducted
 through an online competition platform (Pavao et al., 2023). Notably, it spans seven years of temporal
 coverage, making it particularly suitable for evaluating lifelong VPR scenarios.

AmsterTime (Yildiz et al., 2022) consists of a set of 1,231 grayscale historical images as queries and 1,231 contemporary photos as the gallery in Amsterdam. It contains matching labels between queries and gallery images, confirmed by human experts. This dataset presents one of the most challenging lifelong scenarios, with an extreme temporal gap of *over a century*.

For performance evaluation, we follow the evaluation protocol commonly adopted in previous VPR 357 studies (Arandjelovic et al., 2016; Ge et al., 2020; Warburg et al., 2020; Hausler et al., 2021; Wang 358 et al., 2022; Zhu et al., 2023; Leyva-Vallina et al., 2023). We use a distance threshold of 25 meters to 359 identify positive matches by calculating physical distances from the provided location data, except 360 for AmsterTime (Yildiz et al., 2022) and Nordland (Sünderhauf et al., 2013). AmsterTime provides 361 predefined query-positive pairs, and thus we use the labels directly without distance calculation. For 362 Nordland, which consists of aligned frames from four seasons, a query is considered accurately 363 localized if one of the top-N predictions falls within ten frames of its corresponding ground truth in 364 the database. As is standard in VPR, we utilize recall@K (R@K) as the evaluation metric, measuring the proportion of queries with at least one positive image among the top-K shortlisted results. 365

366 367

4.2 IMPLEMENTATION DETAILS

368 To validate the effectiveness of our method, we employ several deployable architectures, including 369 VGG-16 (Simonyan & Zisserman, 2014), ResNet-50 (He et al., 2016), DINOv2 (Oquab et al., 370 2023), as backbone networks. For CNN backbones, we extract global descriptors using GeM 371 pooling (Radenović et al., 2018) followed by a fully connected layer. For DINOv2, following 372 Izquierdo & Civera (2024), we apply two-layer MLPs for token embeddings, followed by a fully 373 connected layer for the dimensionality reduction, and fine-tune only the last four blocks. Following 374 the class labeling strategy by EigenPlaces (Berton et al., 2023), we utilize the SF-XL training 375 dataset, cropping a total of 6.72M images from 3.43M panoramas. These images are classified into approximately 263.9k classes and further organized into 18 groups, with each training epoch 376 involving two groups. For training stability, we temporarily freeze our CRLS during the initial 9 377 epochs as a warm-up, out of a total of 40 epochs. We set the smoothing parameter $\alpha = 0.2$ and

392

393

396 397

408 409

410

			Evn	sham	Pi	tts	Pi	tts	Tol	суо	Sa	ın.	
Method	Backbone	Dim.	Lyn.	Silaili	- 30)k	25	250k		24/7		Landmark	
			R@1	R@5									
TransVPR	-	256	79.8	87.9	60.7	80.6	54.5	74.8	33.3	51.7	26.4	42.0	
StructVPR	MobileNetV2	448	-	-	85.1	92.3	-	-	-	-	-	-	
R^2 Former	ViT-S	256	82.4	90.3	73.1	<u>88.7</u>	70.2	87.5	48.3	65.1	42.8	61.4	
GCL	VGG-16	512	69.0	79.4	61.7	80.1	53.4	72.5	37.1	57.5	35.8	51.3	
CosPlace	VGG-16	512	88.3	92.7	88.4	94.6	89.7	96.6	82.5	90.8	80.8	<u>87.5</u>	
EigenPlaces	VGG-16	512	<u>89.4</u>	<u>93.6</u>	<u>89.7</u>	<u>95.0</u>	<u>91.2</u>	<u>96.8</u>	82.5	<u>90.8</u>	<u>83.8</u>	90.6	
Ours	VGG-16	512	89.5	93.7	90.1	95.2	91.6	96.9	81.0	92.4	84.8	90.6	
R^2 Former	ResNet-50	256	84.9	91.3	76.5	90.3	72.5	88.1	51.7	70.2	50.5	63.5	
MixVPR	ResNet-50	512	87.8	92.1	90.6	95.6	93.2	98.0	79.4	88.6	79.8	86.3	
MixVPR	ResNet-50	4096	89.6	93.2	91.6	95.6	<u>94.1</u>	<u>98.1</u>	86.3	91.1	84.6	90.3	
GCL	ResNet-50	2048	71.3	82.1	72.0	87.5	68.0	84.4	43.2	59.7	41.3	57.0	
CosPlace	ResNet-50	2048	90.0	93.9	90.9	95.7	92.3	97.4	87.3	94.0	87.1	91.1	
EigenPlaces	ResNet-50	2048	<u>90.7</u>	<u>94.4</u>	92.5	96.8	<u>94.1</u>	97.9	<u>92.7</u>	<u>96.2</u>	<u>89.6</u>	<u>94.3</u>	
Ours	ResNet-50	2048	90.9	94.6	92.3	96.3	94.2	98.2	94.0	96.8	91.6	95.2	

Table 3: Comparison on the multi-view category. The layout is the same as in Table 2.

Table 4: **Comparison on the single-view category.** The layout is the same as in Table 2. (*) indicates the use of high computational resources. Recall@1 (%) is reported.

-	Method	Backbone	Dim.	SVOX Night	SVOX Overcast	SVOX Rain	SVOX Snow	SVOX Sun	St Lucia	Nordland
=	TransVPR	-	256	6.4	61.1	26.9	47.0	13.3	81.4	22.2
	StructVPR	MobileNetV2	448	-	-	-	-	-	-	56.1
	R^2 Former	ViT-S	256	13.5	75.7	47.6	60.7	28.1	93.4	24.6
-	GCL	VGG-16	512	4.4	57.2	32.4	48.0	9.0	59.1	13.3
	CosPlace	VGG-16	512	44.8	88.5	85.2	89.0	67.3	95.3	<u>58.5</u>
	EigenPlaces	VGG-16	512	42.3	89.4	83.5	89.2	69.7	95.4	54.5
	Ours	VGG-16	512	47.6	91.5	85.3	90.5	70.8	96.2	59.5
-	R^2 Former	ResNet-50	256	22.4	78.1	54.4	69.8	34.2	90.0	31.9
	MixVPR	ResNet-50	512	45.8	93.8	86.9	93.6	79.2	99.2	66.5
	MixVPR	ResNet-50	4096*	62.9	96.2	92.1	97.0	85.4	99.5	76.7
	GCL	ResNet-50	2048	8.4	54.5	34.4	47.0	11.0	74.8	13.9
	CosPlace	ResNet-50	2048	50.7	92.2	87.0	92.0	78.5	99.6	71.8
	EigenPlaces	ResNet-50	2048	58.9	93.1	90.0	93.1	86.4	99.6	71.2
	Ours	ResNet-50	2048	64.6	<u>94.0</u>	<u>90.3</u>	<u>94.1</u>	<u>85.5</u>	99.6	<u>73.1</u>

temperature $\tau = 0.1$. The model is trained with a batch size of 320 using Adam (Diederik, 2014) optimizer with a learning rate of 1×10^{-4} . All training is performed on two RTX 3090 GPUs, and a complete training takes approximately 17 hours.

4.3 COMPARISON WITH STATE-OF-THE-ARTS

411 We conduct comprehensive comparisons of our method with recent state-of-the-art methods, including 412 EigenPlaces (Berton et al., 2023), CosPlace (Berton et al., 2022), GCL (Leyva-Vallina et al., 2023), R^2 Former (Zhu et al., 2023), MixVPR (Ali-Bey et al., 2023), StructVPR (Shen et al., 2023) and 413 TransVPR (Wang et al., 2022). We utilize author-released pre-trained networks for benchmark 414 evaluations. For StructVPR, as its code or pre-trained network has not been publicly released, we 415 reference its performance as reported in the original paper. Both CosPlace and EigenPlaces are 416 trained on the SF-XL dataset. GCL, R^2 Former, and StructVPR utilize MSLS for training, while 417 MixVPR is trained on the Google StreetView (GSV) (Ali-bey et al., 2022) dataset. We also report the 418 performance of TransVPR, which is trained on MSLS and employs a custom-designed transformer-419 based backbone. The extensive results are reported in Table 2, 3, 4, representing lifelong, multi-view, 420 and single-view category, respectively. 421

For the lifelong category in Table 2, which is the primary focus of this paper, our approach consis-422 tently achieves state-of-the-art performance across all benchmarks. By effectively utilizing visual 423 relationships across classes, our method successfully handles challenging lifelong variations with 424 superior performance, whereas no single previous method consistently outperforms others in this 425 category. For the multi-view category in Table 3, our results are either comparable to or surpass 426 the current state-of-the-art, with our method achieving better performance in most cases. For the 427 single-view category in Table 4, we report Recall@1 for clarity, with Recall@5 provided in the 428 Appendix. While MixVPR, using ResNet-50 with 4,096 feature dimensions, shows strong results 429 compared to other methods, it requires significantly more computational resources, utilizing 4,096 dimensions versus our method's 2,048 dimensions. Nevertheless, our method demonstrates competitive 430 performance, and notably, we achieve superior results compared to other state-of-the-art methods 431 with the same feature dimensionality. Moreover, MixVPR does not achieve standout performance

432	Table 5: Comparison with the methods using foundation models on the lifelong category. DI	-
433	NOv2 (Oquab et al., 2023) is used as a backbone network. Recall@1 (%) is reported.	

-	Method	Backbone	Dim.	SF-XL test v1	SF-XL test v2	MSLS Val	MSLS Chall.	Amster.	Avg.
=	AnyLoc	DINOv2-G	49152	66.4	83.8	65.0	39.6	40.0	59.0
	SelaVPR	DINOv2-L	1024	55.1	72.1	87.7	69.6	36.9	64.3
	CricaVPR	DINOv2-B	10752	65.8	83.3	89.1	68.1	38.6	69.0
	SALAD	DINOv2-B	8448	88.7	94.5	92.0	<u>75.8</u>	<u>58.6</u>	81.9
	SALAD	DINOv2-B	2112	82.2	93.3	90.8	74.4	54.3	79.0
	Ours	DINOv2-B	2048	93.7	<u>94.0</u>	92.2	77.3	59.9	83.4
-		1 4 4 1	e con		W (D	11 (01)	C 1	1 1	· .1

Table 6: Ablation study of CRLS and CSW. We report Recall@1 (%) across five benchmarks in the lifelong category.

CRLS	CSW	LS	SF-XL test v1	SF-XL test v2	MSLS Val	MSLS Chall.	Amster.
			83.8	90.6	88.4	66.3	47.8
\checkmark			85.0	91.1	89.5	67.8	48.8
\checkmark	\checkmark		85.6	91.8	90.1	68.9	50.4
\checkmark		\checkmark	85.2	90.9	89.6	67.8	48.4
\checkmark	\checkmark	\checkmark	86.0	92.3	90.1	68.8	51.1

in the more challenging lifelong and multi-view benchmarks. In summary, our method consistently
 achieves state-of-the-art or highly competitive performance across all three categories, demonstrating
 its robustness and superiority compared to recent state-of-the-art approaches.

451 Given recent advancements in VPR through the discriminative power of foundation models, we 452 further evaluate our method built on DINOv2 (Oquab et al., 2023) as the backbone, with the 453 results presented in Table 5. We compare our method with other foundation model-based methods, 454 including AnyLoc (Keetha et al., 2023), SelaVPR (Lu et al.), CricaVPR (Lu et al., 2024), and 455 SALAD (Izquierdo & Civera, 2024). Anyloc is a zero-shot method without fine-tuning. For SALAD, 456 with a dimensionality of 2,112, we reproduce the method using the author-released code and provided 457 parameters. The experimental results show that our method, based on DINOv2-B with a compact dimensionality of 2,048, achieves state-of-the-art performance across the lifelong benchmarks. 458 Specifically, our method outperforms all others in terms of average R@1, achieving 83.4% across 459 the five lifelong datasets. It surpasses methods using larger models (e.g., DINOv2-L,G) and higher 460 dimensionalities (e.g., 8,448 for SALAD). Notably, on the challenging SF-XL test v1, our method 461 achieves a substantial 5% improvement over the second-best, SALAD. 462

463 464

441

4.4 ABLATION STUDY

465

466 Effectiveness of CRLS and CSW. We explore the impact of CRLS, CSW, and naïve LS on retrieval performance in Table 6. As a baseline, we use the CosFace loss with one-hot encoded hard labels 467 as EigenPlaces. To solely assess the effectiveness of each component, we employ the vanilla 468 $\mathcal{L}_{CosFace}$ for cases without LS. Specifically, the loss of CRLS with only CSW is calculated as 469 $\mathcal{L} = \gamma_{y_i} \mathcal{L}_{\text{CosFace}} + (1 - \gamma_{y_i}) \mathcal{L}_{\text{CRLS}}$. These experiments are conducted using a ResNet-50 architecture 470 with a feature dimensionality of 2048. We adopt five evaluation datasets that reflect lifelong scenarios: 471 SF-XL test v1, SF-XL test v2, MSLS Val, MSLS Challenge, and AmsterTime. Table 6 demonstrates 472 that while using CRLS alone does provide a performance boost over the baseline, it achieves better 473 results when combined with CSW. This can be attributed to the fact that CRLS is more effective for 474 classes with low weight magnitudes, as these classes tend to exhibit less fluctuation during training. 475 The best performance is obtained under our final loss formulation with LS, where LS appears to 476 compensate for the classes having high magnitudes.

477 To further evaluate the impact of CSW on CRLS, we apply 478 the computed CSW weights in a reversed manner (e.g., 479 $\gamma_{y_i} \rightarrow 1 - \gamma_{y_i}$). This reversal assigns higher weights to 480 classes with high magnitudes, which are more prone to 481 fluctuations during training. The performance degradation 482 observed in Table 8 demonstrates that the contrapositive 483 approach is harmful to the training process. Conversely, this result validates the design choice implemented in 484 CSW, where weights inversely proportional to the magni-485 tude are assigned to stabilize CRLS.

Table 8: Experiments on contrapositive
approach of CSW. R-CSW refers to the
reverse CSW. Recall@1 (%) is reported.

Method	M Val	SLS Chall.	Amster.
Ours	90.1	68.8	51.1 48.5
+ R-CSW	88.9	66.6	

Method	SF-XL test v1	SF-XL test v2	MSLS Val	MSLS Chall.	Amster	
LR (Lienen & Hüllermeier, 2021)	84.0	91.0	88.9	67.9	48.6	
ACLS (Park et al., 2023)	82.2	91.1	87.7	67.2	47.4	
CRLS (Ours)	86.0	92.3	90.1	90.1 68.8		

Table 7: Ablation study of label smoothing strategies used as substitutes for CRLS in our training.

(a) Images within two highest-norm classes

(b) Images within two lowest-norm classes

Figure 4: **Examples of classes according to magnitude.** (a) tends to be hard to be trained and less informative, while (b) shows pleasant classes for appropriate training.

Comparison with other label smoothing strategies. To further validate the effectiveness of CRLS, we substitute it with alternative techniques in our final loss formulation and evaluate their performance.
 Table 7 presents the results of this experiment, where we compare against LR (Lienen & Hüllermeier, 2021) and ACLS (Park et al., 2023), both variants of label smoothing designed for network calibration.
 The results demonstrate that CRLS consistently outperforms these alternatives. While the other techniques are optimized primarily for network calibration, CRLS is specifically tailored for VPR, leveraging visual relationships between classes to learn a continuous representation space for handling diverse appearance variations.

Qualitative analysis on the class weight magnitude. Fig. 4 illustrates that classes with the highest magnitude tend to be less informative, while those with the lowest magnitude are more suitable for training⁴. This distinction in weight magnitudes reflects diverse levels of informativeness and trainability among the training classes. These properties of weight magnitude could be diversely applied in robust training, such as in curriculum learning. Additional qualitative results such as CRLS and retrieval results are presented in Appendix.

516 517

518

486

499

500

501

5 CONCLUSIONS

In this paper, we introduced a novel technique, CRLS, that effectively bridges the task gap between classification and retrieval in visual place recognition, which is particularly vulnerable to extreme visual changes such as lifelong variations. We further enhanced CRLS with CSW, a method that dynamically adjusts the influence of CRLS based on the stability of class weights, quantified by their magnitudes. Our findings, supported by derivative analysis, suggest that the magnitude of class weights serves as an indicator of class stability.

525 We evaluated our approach through extensive experiments on 17 diverse benchmarks, covering a 526 wide range of scenarios including lifelong, multi-view, and single-view settings. In the majority of 527 cases, our method outperformed existing state-of-the-art methods, while in the remaining cases, it 528 achieved highly competitive performance. Particularly in the most challenging lifelong benchmarks, 529 our approach demonstrated state-of-the-art performance with substantial improvements over existing 530 methods. Furthermore, our method showed its effectiveness in leveraging foundation models, such as 531 DINOv2, for VPR. By integrating ours with DINOv2 backbone, we achieved superior performance 532 compared to other methods utilizing the same backbone, while maintaining a compact feature representation. These results show the robustness and adaptability of our method in tackling lifelong 533 challenges in real-world environments, and in various backbones. 534

Interestingly, our analysis revealed that class weight magnitudes indicate the informativeness and
 trainability of classes, which could potentially be utilized in robust training techniques such as
 curriculum learning. Unfortunately, our analysis is based on the CosFace loss function, while our
 analytical approach may facilitate such extensions to other losses.

⁵³⁹

⁴More examples can be found in the Appendix

540 REFERENCES

547

- Amar Ali-bey, Brahim Chaib-draa, and Philippe Giguère. Gsv-cities: Toward appropriate supervised visual place recognition. *Neurocomputing*, 513:194–203, 2022.
- Amar Ali-Bey, Brahim Chaib-Draa, and Philippe Giguere. Mixvpr: Feature mixing for visual place
 recognition. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 2998–3007, 2023.
- 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*, pp. 5297–5307, 2016.
- Relja Arandjelovic, Petr Gronat, Akihiko Torii, Tomas Pajdla, and Josef Sivic. Netvlad: Cnn architecture for weakly supervised place recognition. *IEEE Transactions on Pattern Analysis & Machine Intelligence*, 40(06):1437–1451, 2018.
- Herbert Bay, Andreas Ess, Tinne Tuytelaars, and Luc Van Gool. Speeded-up robust features (surf).
 Computer vision and image understanding, 110(3):346–359, 2008.
- Gabriele Berton, Carlo Masone, and Barbara Caputo. Rethinking visual geo-localization for large scale applications. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pp. 4878–4888, 2022.
- Gabriele Berton, Gabriele Trivigno, Barbara Caputo, and Carlo Masone. Eigenplaces: Training view-point robust models for visual place recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 11080–11090, 2023.
- Gabriele Moreno Berton, Valerio Paolicelli, Carlo Masone, and Barbara Caputo. Adaptive-attentive
 geolocalization from few queries: A hybrid approach. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*, pp. 2918–2927, 2021.
- Guillaume Bresson, Zayed Alsayed, Li Yu, and Sébastien Glaser. Simultaneous localization and
 mapping: A survey of current trends in autonomous driving. *IEEE Transactions on Intelligent Vehicles*, 2(3):194–220, 2017.
- Nicholas Carlevaris-Bianco, Arash K Ushani, and Ryan M Eustice. University of michigan north
 campus long-term vision and lidar dataset. *The International Journal of Robotics Research*, 35(9): 1023–1035, 2016.
- 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*, pp. 737–744. IEEE, 2011.
- 577
 578
 578
 579
 579
 580
 570
 570
 570
 571
 572
 573
 574
 575
 575
 575
 576
 577
 577
 578
 579
 579
 570
 570
 570
 571
 571
 572
 573
 573
 574
 575
 576
 576
 577
 576
 577
 578
 578
 578
 578
 579
 578
 579
 578
 579
 578
 579
 578
 579
 578
 579
 578
 579
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
 578
- Mark Cummins and Paul Newman. Highly scalable appearance-only slam-fab-map 2.0. In *Robotics: Science and systems*, volume 5, pp. 17. Seattle, USA, 2009.
- ⁵⁸³ P Kingma Diederik. Adam: A method for stochastic optimization. (*No Title*), 2014.
- Yixiao Ge, Haibo Wang, Feng Zhu, Rui Zhao, and Hongsheng Li. Self-supervising fine-grained region similarities for large-scale image localization. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part IV 16*, pp. 369–386.
 Springer, 2020.
- Petr Gronat, Guillaume Obozinski, Josef Sivic, and Tomas Pajdla. Learning and calibrating perlocation classifiers for visual place recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pp. 907–914, 2013.
- 593 Chuan Guo, Geoff Pleiss, Yu Sun, and Kilian Q Weinberger. On calibration of modern neural networks. In *International conference on machine learning*, pp. 1321–1330. PMLR, 2017.

- 594 Stephen Hausler, Sourav Garg, Ming Xu, Michael Milford, and Tobias Fischer. Patch-netvlad: Multi-595 scale fusion of locally-global descriptors for place recognition. In Proceedings of the IEEE/CVF 596 Conference on Computer Vision and Pattern Recognition, pp. 14141–14152, 2021. 597 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image 598 recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 770-778, 2016. 600 601 Mike Izbicki, Evangelos E Papalexakis, and Vassilis J Tsotras. Exploiting the earth's spherical geom-602 etry to geolocate images. In Machine Learning and Knowledge Discovery in Databases: European 603 Conference, ECML PKDD 2019, Würzburg, Germany, September 16–20, 2019, Proceedings, Part 604 II, pp. 3–19. Springer, 2020. 605 Sergio Izquierdo and Javier Civera. Optimal transport aggregation for visual place recognition. 606 In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 607 17658-17668, 2024. 608 609 Mengxi Jia, Xinhua Cheng, Shijian Lu, and Jian Zhang. Learning disentangled representation 610 implicitly via transformer for occluded person re-identification. IEEE Transactions on Multimedia, 611 25:1294–1305, 2022. 612 613 Nikhil Keetha, Avneesh Mishra, Jay Karhade, Krishna Murthy Jatavallabhula, Sebastian Scherer, Madhava Krishna, and Sourav Garg. Anyloc: Towards universal visual place recognition. IEEE 614 Robotics and Automation Letters, 2023. 615 616 Sungyeon Kim, Dongwon Kim, Minsu Cho, and Suha Kwak. Embedding transfer with label 617 relaxation for improved metric learning. In Proceedings of the IEEE/CVF conference on computer 618 vision and pattern recognition, pp. 3967–3976, 2021. 619 620 Sungyeon Kim, Dongwon Kim, Minsu Cho, and Suha Kwak. Self-taught metric learning without 621 labels. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, 622 pp. 7431-7441, 2022. 623 Jan Knopp, Josef Sivic, and Tomas Pajdla. Avoiding confusing features in place recognition. In 624 Computer Vision–ECCV 2010: 11th European Conference on Computer Vision, Heraklion, Crete, 625 Greece, September 5-11, 2010, Proceedings, Part I 11, pp. 748-761. Springer, 2010. 626 627 Giorgos Kordopatis-Zilos, Panagiotis Galopoulos, Symeon Papadopoulos, and Ioannis Kompatsiaris. 628 Leveraging efficientnet and contrastive learning for accurate global-scale location estimation. In 629 Proceedings of the 2021 International Conference on Multimedia Retrieval, pp. 155–163, 2021. 630 María Leyva-Vallina, Nicola Strisciuglio, and Nicolai Petkov. Generalized contrastive optimization 631 of siamese networks for place recognition. arXiv preprint arXiv:2103.06638, 2021. 632 633 María Leyva-Vallina, Nicola Strisciuglio, and Nicolai Petkov. Data-efficient large scale place 634 recognition with graded similarity supervision. In Proceedings of the IEEE/CVF Conference on 635 Computer Vision and Pattern Recognition, pp. 23487–23496, 2023. 636 637 Julian Lienen and Eyke Hüllermeier. From label smoothing to label relaxation. In Proceedings of the 638 AAAI conference on artificial intelligence, volume 35, pp. 8583–8591, 2021. 639 Bingyuan Liu, Ismail Ben Ayed, Adrian Galdran, and Jose Dolz. The devil is in the margin: Margin-640 based label smoothing for network calibration. In Proceedings of the IEEE/CVF Conference on 641 Computer Vision and Pattern Recognition, pp. 80–88, 2022. 642 643 Bingyuan Liu, Jérôme Rony, Adrian Galdran, Jose Dolz, and Ismail Ben Ayed. Class adaptive 644 network calibration. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern 645 Recognition, pp. 16070–16079, 2023. 646 David G Lowe. Distinctive image features from scale-invariant keypoints. International journal of
- 647 David G Lowe. Distinctive image features from scale-invariant keypoints. *International journal of computer vision*, 60:91–110, 2004.

659

662

668

677

684

692

- 648 Feng Lu, Lijun Zhang, Xiangyuan Lan, Shuting Dong, Yaowei Wang, and Chun Yuan. Towards 649 seamless adaptation of pre-trained models for visual place recognition. In The Twelfth International 650 Conference on Learning Representations. 651
- Feng Lu, Xiangyuan Lan, Lijun Zhang, Dongmei Jiang, Yaowei Wang, and Chun Yuan. Cricavpr: 652 Cross-image correlation-aware representation learning for visual place recognition. In Proceedings 653 of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 16772–16782, 654 2024. 655
- 656 Hao Luo, Wei Jiang, Youzhi Gu, Fuxu Liu, Xingyu Liao, Shenqi Lai, and Jianyang Gu. A strong baseline and batch normalization neck for deep person re-identification. *IEEE Transactions on* 658 Multimedia, 22(10):2597-2609, 2019.
- Will Maddern, Geoffrey Pascoe, Chris Linegar, and Paul Newman. 1 year, 1000 km: The oxford 660 robotcar dataset. The International Journal of Robotics Research, 36(1):3-15, 2017. 661
- Qiang Meng, Shichao Zhao, Zhida Huang, and Feng Zhou. Magface: A universal representation for 663 face recognition and quality assessment. In Proceedings of the IEEE/CVF conference on computer 664 vision and pattern recognition, pp. 14225-14234, 2021. 665
- 666 Michael J Milford and Gordon F Wyeth. Mapping a suburb with a single camera using a biologically 667 inspired slam system. IEEE Transactions on Robotics, 24(5):1038-1053, 2008.
- Piotr Mirowski, Matt Grimes, Mateusz Malinowski, Karl Moritz Hermann, Keith Anderson, Denis 669 Teplyashin, Karen Simonyan, Andrew Zisserman, Raia Hadsell, et al. Learning to navigate in 670 cities without a map. Advances in neural information processing systems, 31, 2018. 671
- 672 Rafael Müller, Simon Kornblith, and Geoffrey E Hinton. When does label smoothing help? Advances 673 in neural information processing systems, 32, 2019. 674
- 675 Tayyab Naseer, Wolfram Burgard, and Cyrill Stachniss. Robust visual localization across seasons. IEEE Transactions on Robotics, 34(2):289–302, 2018. 676
- Maxime Oquab, Timothée Darcet, Théo Moutakanni, Huy Vo, Marc Szafraniec, Vasil Khalidov, 678 Pierre Fernandez, Daniel Haziza, Francisco Massa, Alaaeldin El-Nouby, et al. Dinov2: Learning 679 robust visual features without supervision. arXiv preprint arXiv:2304.07193, 2023. 680
- 681 Hyekang Park, Jongyoun Noh, Youngmin Oh, Donghyeon Baek, and Bumsub Ham. Acls: Adap-682 tive and conditional label smoothing for network calibration. In Proceedings of the IEEE/CVF 683 International Conference on Computer Vision, pp. 3936–3945, 2023.
- Adrien Pavao, Isabelle Guyon, Anne-Catherine Letournel, Dinh-Tuan Tran, Xavier Baro, Hugo Jair 685 Escalante, Sergio Escalera, Tyler Thomas, and Zhen Xu. Codalab competitions: An open source 686 platform to organize scientific challenges. Journal of Machine Learning Research, 24(198):1-6, 687 2023. URL http://jmlr.org/papers/v24/21-1436.html. 688
- 689 Guohao Peng, Jun Zhang, Heshan Li, and Danwei Wang. Attentional pyramid pooling of salient 690 visual residuals for place recognition. In Proceedings of the IEEE/CVF International Conference 691 on Computer Vision, pp. 885-894, 2021.
- Gabriel Pereyra, George Tucker, Jan Chorowski, Łukasz Kaiser, and Geoffrey Hinton. Regularizing 693 neural networks by penalizing confident output distributions. arXiv preprint arXiv:1701.06548, 694 2017.
- 696 Shraman Pramanick, Ewa M Nowara, Joshua Gleason, Carlos D Castillo, and Rama Chellappa. 697 Where in the world is this image? transformer-based geo-localization in the wild. In European 698 Conference on Computer Vision, pp. 196-215. Springer, 2022.
- Filip Radenović, Giorgos Tolias, and Ondřej Chum. Fine-tuning cnn image retrieval with no human 700 annotation. *IEEE transactions on pattern analysis and machine intelligence*, 41(7):1655–1668, 701 2018.

702 703 704 705	German Ros, Laura Sellart, Joanna Materzynska, David Vazquez, and Antonio M Lopez. The synthia dataset: A large collection of synthetic images for semantic segmentation of urban scenes. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 3234–3243, 2016.
706 707 708	Torsten Sattler, Tobias Weyand, Bastian Leibe, and Leif Kobbelt. Image retrieval for image-based localization revisited. In <i>BMVC</i> , pp. 4, 2012.
709 710 711 712 713	Yanqing Shen, Sanping Zhou, Jingwen Fu, Ruotong Wang, Shitao Chen, and Nanning Zheng. Structvpr: Distill structural knowledge with weighting samples for visual place recognition. In Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition, pp. 11217–11226, 2023.
714 715	Karen Simonyan and Andrew Zisserman. Very deep convolutional networks for large-scale image recognition. <i>arXiv preprint arXiv:1409.1556</i> , 2014.
716 717 718 719	Elena Stumm, Christopher Mei, and Simon Lacroix. Probabilistic place recognition with covisibility maps. In 2013 IEEE/RSJ International Conference on Intelligent Robots and Systems, pp. 4158–4163. IEEE, 2013.
720 721 722	Niko Sünderhauf, Peer Neubert, and Peter Protzel. Are we there yet? challenging seqslam on a 3000 km journey across all four seasons. In <i>Proc. of workshop on long-term autonomy, IEEE international conference on robotics and automation (ICRA)</i> , pp. 2013. Citeseer, 2013.
723 724 725 726	Christian Szegedy, Vincent Vanhoucke, Sergey Ioffe, Jon Shlens, and Zbigniew Wojna. Rethinking the inception architecture for computer vision. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 2818–2826, 2016.
727 728 729	Janine Thoma, Danda Pani Paudel, and Luc V Gool. Soft contrastive learning for visual localization. Advances in Neural Information Processing Systems, 33:11119–11130, 2020.
730 731 732	Akihiko Torii, Josef Sivic, Tomas Pajdla, and Masatoshi Okutomi. Visual place recognition with repetitive structures. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 883–890, 2013a.
733 734 735 736	Akihiko Torii, Josef Sivic, Tomas Pajdla, and Masatoshi Okutomi. Visual place recognition with repetitive structures. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 883–890, 2013b.
737 738 739	Akihiko Torii, Relja Arandjelovic, Josef Sivic, Masatoshi Okutomi, and Tomas Pajdla. 24/7 place recognition by view synthesis. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 1808–1817, 2015.
740 741 742 743	Gabriele Trivigno, Gabriele Berton, Juan Aragon, Barbara Caputo, and Carlo Masone. Di- vide&classify: Fine-grained classification for city-wide visual place recognition. In 2023 IEEE/CVF International Conference on Computer Vision (ICCV), pp. 11108–11118. IEEE, 2023.
744 745 746	Vicente Vivanco Cepeda, Gaurav Kumar Nayak, and Mubarak Shah. Geoclip: Clip-inspired alignment between locations and images for effective worldwide geo-localization. <i>Advances in Neural Information Processing Systems</i> , 36, 2023.
747 748 749	Cheng Wang. Calibration in deep learning: A survey of the state-of-the-art. <i>arXiv preprint arXiv:2308.01222</i> , 2023.
750 751 752 753	Hao Wang, Yitong Wang, Zheng Zhou, Xing Ji, Dihong Gong, Jingchao Zhou, Zhifeng Li, and Wei Liu. Cosface: Large margin cosine loss for deep face recognition. In <i>Proceedings of the IEEE conference on computer vision and pattern recognition</i> , pp. 5265–5274, 2018.
754 755	Ruotong Wang, Yanqing Shen, Weiliang Zuo, Sanping Zhou, and Nanning Zheng. Transvpr: Transformer-based place recognition with multi-level attention aggregation. In <i>Proceedings of the</i> <i>IEEE/CVF Conference on Computer Vision and Pattern Recognition</i> , pp. 13648–13657, 2022.

756 757 758	Xiaolong Wang, Runsen Xu, Zhuofan Cui, Zeyu Wan, and Yu Zhang. Fine-grained cross-view geo-localization using a correlation-aware homography estimator. <i>Advances in Neural Information Processing Systems</i> , 36, 2023.
759 760 761 762	Frederik Warburg, Soren Hauberg, Manuel Lopez-Antequera, Pau Gargallo, Yubin Kuang, and Javier Civera. Mapillary street-level sequences: A dataset for lifelong place recognition. In <i>Proceedings</i> of the IEEE/CVF conference on computer vision and pattern recognition, pp. 2626–2635, 2020.
763 764 765	Tobias Weyand, Ilya Kostrikov, and James Philbin. Planet-photo geolocation with convolutional neural networks. In <i>Computer Vision–ECCV 2016: 14th European Conference, Amsterdam, The Netherlands, October 11-14, 2016, Proceedings, Part VIII 14</i> , pp. 37–55. Springer, 2016.
766 767 768 769	Burak Yildiz, Seyran Khademi, Ronald Maria Siebes, and Jan Van Gemert. Amstertime: A visual place recognition benchmark dataset for severe domain shift. In 2022 26th International Conference on Pattern Recognition (ICPR), pp. 2749–2755. IEEE, 2022.
770 771 772	Sijie Zhu, Linjie Yang, Chen Chen, Mubarak Shah, Xiaohui Shen, and Heng Wang. R2former: Unified retrieval and reranking transformer for place recognition. In <i>Proceedings of the IEEE/CVF</i> <i>Conference on Computer Vision and Pattern Recognition</i> , pp. 19370–19380, 2023.
773 774 775 776	Zhihui Zhu, Xinyang Jiang, Feng Zheng, Xiaowei Guo, Feiyue Huang, Xing Sun, and Weishi Zheng. Aware loss with angular regularization for person re-identification. In <i>Proceedings of the AAAI conference on artificial intelligence</i> , pp. 13114–13121, 2020.
777	
778	
779	
780	
781	
782	
783	
784	
785	
786	
787	
788	
789	
790	
791	
792	
793	
794	
796	
797	
798	
799	
800	
801	
802	
803	
804	
805	
806	
807	
808	
809	

810 APPENDIX

A DERIVATION OF CLASS WEIGHT

The derivative of L_{CE} w.r.t W_{y_i} is computed by chain rule as follows:

$$\frac{\partial L_{CE}}{\partial W_{y_i}} = \frac{\partial L_{CE}}{\partial l(\cos \theta_{y_i})} \frac{\partial l(\cos \theta_{y_i})}{\partial \cos \theta_{y_i}} \frac{\partial \cos \theta_{y_i}}{\partial W_{y_i}}.$$
(10)

The derivatives w.r.t logits and cosine similarity can be computed as:

$$\frac{\partial L_{CE}}{\partial l(\cos \theta_{y_i})} = p_{y_i} - 1, \ \frac{\partial l(\cos \theta_{y_i})}{\partial \cos \theta_{y_i}} = s.$$
(11)

By the quotient rule, the derivative of cosine similarity w.r.t the class weight is:

$$\frac{\partial \cos \theta_{y_i}}{\partial W_{y_i}} = \frac{\partial}{\partial W_{y_i}} \left(\frac{W_{y_i}^T x_i}{\|W_{y_i}\| \|x_i\|} \right) = \frac{\|W_{y_i}\| \|x_i\| x_i - \frac{W_{y_i}}{\|W_{y_i}\|} \|x_i\| W_{y_i}^T x_i}{(\|W_{y_i}\| \|x_i\|)^2} \\
= \frac{1}{\|W_{y_i}\|} \left(\frac{x_i}{\|x_i\|} - \cos \theta_{y_i} \frac{W_{y_i}}{\|W_{y_i}\|} \right).$$
(12)

Incidentally, the gradient vector results in a tangent vector at a point W_{y_i} on a unit sphere, as $\left(\frac{x_i}{\|x_i\|} - \cos \theta_{y_i} \frac{W_{y_i}}{\|W_{y_i}\|}\right) \cdot \frac{W_{y_i}}{\|W_{y_i}\|} = 0.$ Since the terms in Eq. 11 are constants, to compute the magnitude of the derivative, we calculate the norm of $\frac{\partial \cos \theta_{y_i}}{\partial W_{y_i}}$ as follows:

$$\left\|\frac{\partial\cos\theta_{y_i}}{\partial W_{y_i}}\right\| = \frac{1}{\|W_{y_i}\|} \sqrt{\left(\frac{x_i}{\|x_i\|} - \cos\theta_{y_i}\frac{W_{y_i}}{\|W_{y_i}\|}\right)^T \left(\frac{x_i}{\|x_i\|} - \cos\theta_{y_i}\frac{W_{y_i}}{\|W_{y_i}\|}\right)} = \frac{1}{\|W_{y_i}\|} \sqrt{1 - \cos^2\theta_{y_i}}.$$
(13)

Finally, we can express the magnitude of the derivative of the loss function w.r.t the class weight as:

$$\left\|\frac{\partial L_{CE}}{\partial W_{y_i}}\right\| = \frac{\partial L_{CE}}{\partial l(\cos \theta_{y_i})} \frac{\partial l(\cos \theta_{y_i})}{\partial \cos \theta_{y_i}} \left\|\frac{\partial \cos \theta_{y_i}}{\partial W_{y_i}}\right\|$$

$$= (1 - p_{y_i}) \sqrt{1 - \cos^2 \theta_{y_i}} \frac{s}{\|W_{y_i}\|} \quad \because p_{y_i} \in [0, 1],$$
(14)

where the overwhelming majority of $\cos \theta_{y_i}$ are greater than zero, as empirically shown in Fig. 5.



Figure 5: Histogram of $\cos \theta_{y_i}$ values calculated for all input *i*, with the vertical red dashed line indicating $\cos \theta_{y_i} = 0$.

888

889

894

909 910



(b) An anchor class of forest scene

Figure 6: Qualitative analysis of the impact by CRLS. The class weights are considered as embedding features, with the class weight of the anchor class treated as the query feature and the weights of other classes treated as database features. The results of top-1, 2 are obtained through the execution of the retrieval procedure at the class level.

Table 9: Comparison on the single-view category. The layout is the same as in Table 2. (*) indicates the use of high computational resources. Recall@5 (%) is reported.

205										
035	Mathad	Dealthona	Dim	SVOX	SVOX	SVOX	SVOX	SVOX	St Lucio	Nordland
896	Wiethou	DackDone	Diili.	Night	Overcast	Rain	Snow	Sun	St Lucia	Norulaliu
897	TransVPR	-	256	15.2	80.5	49.3	72.0	29.2	90.4	37.6
898	StructVPR	MobileNetV2	448	-	-	-	-	-	-	75.5
899	R^2 Former	ViT-S	256	30.4	91.2	68.8	83.8	46.6	98.1	38.8
000	GCL	VGG-16	512	12.8	74.5	46.9	69.4	15.9	75.3	22.6
900	CosPlace	VGG-16	512	63.5	93.9	91.7	94.0	79.2	98.1	73.7
901	EigenPlaces	VGG-16	512	61.0	<u>94.4</u>	91.6	<u>94.4</u>	82.2	<u>98.3</u>	70.1
902	Ours	VGG-16	512	66.8	96.1	93.5	95.1	<u>81.4</u>	98.7	74.9
903	R^2 Former	ResNet-50	256	43.4	92.9	73.7	88.5	55.6	95.5	48.6
004	MixVPR	ResNet-50	512	62.7	97.8	93.8	97.6	90.7	<u>99.9</u>	80.8
904	MixVPR	ResNet-50	4096*	79.8	98.2	96.8	98.3	93.0	100.0	87.1
905	GCL	ResNet-50	2048	16.4	74.2	56.1	67.7	25.6	86.6	24.7
906	CosPlace	ResNet-50	2048	67.4	97.7	95.1	98.4	89.7	<u>99.9</u>	83.8
907	EigenPlaces	ResNet-50	2048	76.9	97.9	96.4	97.6	95.0	<u>99.9</u>	83.8
908	Ours	ResNet-50	2048	83.0	<u>98.0</u>	<u>96.5</u>	98.5	<u>94.6</u>	<u>99.9</u>	<u>84.3</u>

Table 10: Comparison with the methods using foundation models on the lifelong category. DINOv2 (Oquab et al., 2023) is employed as a backbone network. Recall@5 (%) is reported.

912	Method	Backbone	Dim.	SF-XL test v1	SF-XL test v2	MSLS Val	MSLS Chall.	Amster.
913	AnvLoc	DINOv2-G	49152	78.1	92.3	75.4	53.1	63.8
914	SelaVPR	DINOv2-L	1024	68.1	87.1	95.8	86.9	59.8
915	CricaVPR	DINOv2-B	10752	76.5	91.0	95.0	80.0	60.0
916	SALAD	DINOv2-B	8448	<u>93.5</u>	<u>97.4</u>	96.2	89.2	<u>78.9</u>
017	SALAD	DINOv2-B	2112	89.4	97.3	<u>96.1</u>	87.1	75.1
317	Ours	DINOv2-B	2048	96.6	97.5	95.7	<u>87.9</u>	80.7

Table 11: Ablation studies on the hyper-parameters. We report Recall@1 (%) across five benchmarks in the lifelong category.

(a) Experiments on the smoothing intensity α used in LS and CRLS.

α	SF-XL test v1	SF-XL test v2	MSLS Val	MSLS Chall.	Amster.
0.0	83.8	90.6	88.4	66.3	47.8
0.1	85.1	91.6	89.2	68.1	49.4
0.2	86.0	92.3	90.1	68.8	51.1
0.3	84.5	91.3	89.2	68.7	51.6

(b) Experiments on the temperature parameter $1/\tau$.

$1/\tau$	SF-XL test v1	SF-XL test v2	MSLS Val	MSLS Chall.	Amster.
0	85.1	91.1	88.7	67.8	50.1
10	86.0	92.3	90.1	68.8	51.1
20	84.6	91.8	90.3	68.7	49.8
30	84.9	91.5	88.2	68.7	49.5

B FURTHER QUALITATIVE AND QUANTITATIVE RESULTS

We prepared examples of the lifelong scenarios and retrieval results to show the effectiveness of
our approach. As shown in Fig. 7, AmsterTime uses historical images as queries, which leads to a
scenario where buildings in the query and positive images have undergone extreme changes. For
EigenPlace, which is trained using classification loss, a goal of the loss is to find an exactly same
object. This makes EigenPlace hard to find the matching pairs with extreme changes. In contrast, our
method learns in a visually similar-aware manner, and thus the lifelong examples can be handled.

942To further validate the effectiveness of CRLS, we conduct retrieval at the class level using class943weights in Fig. 6. Upon comparison with the baseline, our method demonstrates that our class weights944on representation space better capture semantic or visually-similar information while considering945visual differences.

We provide additional quantitative results to further demonstrate the effectiveness of our method.
Table 9 presents the Recall@5 performance on the single-view category benchmarks, complementing
the Recall@1 results shown in the main paper. Similarly, Table 10 reports the Recall@5 performance
of our method and other approaches utilizing the DINOv2 foundation model on the lifelong category
benchmarks, providing a more comprehensive evaluation of their performance.

951 952

953 954

955

956

957

958

920

934 935

C ABLATIONS ON HYPER-PARAMETERS

In addition, we study varying the hyper-parameter α , representing the intensity of smoothing as shown in Table 11a. The ablation study finds that $\alpha = 0.2$ delivers the best performances in most cases, validating its selection as the optimal setting for our experiments. We also examine the effect of the temperature parameter τ in our CRLS approach, with results shown in Table 11b. The experiments reveal that a value of $1/\tau = 10$ ($\tau = 0.1$) generally yields the best performance across the five lifelong datasets, providing an optimal balance.

959 960 961

962

D CLASSIFICATION-BASED APPROACHES IN VISUAL GEO-LOCALIZATION

963 Contrasting with retrieval-based approaches, another significant area of research in Visual Place 964 Recognition (VPR) is focused on classification-based approaches Weyand et al. (2016); Vi-965 vanco Cepeda et al. (2023); Izbicki et al. (2020); Pramanick et al. (2022); Trivigno et al. (2023); 966 Kordopatis-Zilos et al. (2021). Unlike retrieval methods that strive to match a query image with a 967 large database of reference images, classification-based approaches divide the geographic area into 968 discrete cells or regions, with each cell treated as a separate class. This framework transforms the task into a classification problem, where the objective is to identify the correct geographic cell for a given 969 query image. These methods are advantageous for their faster inference times, as they bypass the need 970 for extensive similarity searches required by retrieval-based methods. However, despite the efficiency 971 of classification-based methods, retrieval methods can leverage the fine-grained similarities between

query and database images, enabling more precise localization. Moreover, they are not limited by predefined classes and can potentially localize images at any location covered in the database.



Figure 7: Examples of lifelong scenarios and their top-1 retrieval results from the AmsterTime and SF-XL test v1 datasets. The query and positive images demonstrate significant changes in structure and appearance due to building remodeling over time. We compare the retrieval results of our method with those of EigenPlace, highlighting our method's ability to handle these challenging scenarios effectively.

1021

972

- 1022
- 1023
- 1024
- 1025



Under review as a conference paper at ICLR 2025