

Segment Augmentation and Differentiable Ranking for Logo Retrieval

Feyza Yavuz

Dept. of Computer Engineering
Middle East Technical University, Ankara, Turkey
Email: feyza.yavuz@metu.edu.tr

Sinan Kalkan

Dept. of Computer Engineering
Middle East Technical University, Ankara, Turkey
Email: skalkan@metu.edu.tr

Abstract—Logo retrieval is a challenging problem since the definition of similarity is more subjective than image retrieval, and the set of known similarities is very scarce. In this paper, to tackle this challenge, we propose a simple but effective segment-based augmentation strategy to introduce artificially similar logos for training deep networks for logo retrieval. In this novel augmentation strategy, we first find segments in a logo and apply transformations such as rotation, scaling, and color change, on the segments, unlike the conventional strategies that perform augmentation at the image level. Moreover, we evaluate suitability of using ranking-based losses (namely Smooth-AP) for learning similarity for logo retrieval. On the METU and the LLD datasets, we show that (i) our segment-based augmentation strategy improves retrieval performance compared to the baseline model or image-level augmentation strategies, and (ii) Smooth-AP indeed performs better than conventional losses for logo retrieval.

For a query logo, identifying the similar ones in a database of logos is a content-based image retrieval problem. With the rise in deep learning, there have been many studies that have used deep learning for logo retrieval problem [1]–[5]. Existing approaches generally rely on extracting features of logos and ranking them according to a suitable distance metric [1], [5], [6].

Logo retrieval is a challenging problem especially for two main reasons: (i) Similarity between logos is highly subjective, and similarity can occur at different levels, e.g., texture, color, segments and their combination etc. (ii) The amount of known similar logos is limited. We hypothesize that this has limited the use of more modern deep learning solutions, e.g. metric learning, contrastive learning, differentiable ranking, as they require tremendous amount of positive pairs (similar logos) as well as negative pairs (dissimilar logos) for training deep networks.

In this paper, we address these challenges by (i) proposing a segment-level augmentation to produce artificially similar logos and (ii) using metric learning (Triplet Loss [7]) and differentiable ranking (Smooth Average Precision (AP) Loss [8]) as a proof of concept that, with our novel segment-augmentation method, such data hungry techniques can be trained better.

Main Contributions. Our contributions are as follows:

- We propose a segment-level augmentation for producing artificial similarities between logos. To the best of our knowledge, ours is the first to introduce segment-level augmentation into deep learning. Unlike image-level augmentation methods that transform the overall image, we identify segments in a logo and make transformations at the segment level. Our results suggest that this is more suitable than image-level augmentation for logo retrieval.
- To showcase the use of such a tool to generate artificially similar logos, we use data-hungry deep learning methods, namely, Triplet Loss [7] and Smooth-AP Loss [8], to show that our novel segment-augmentation method can indeed yield better retrieval performance. To the best of our knowledge, ours is the first to use such methods for logo retrieval.

I. INTRODUCTION

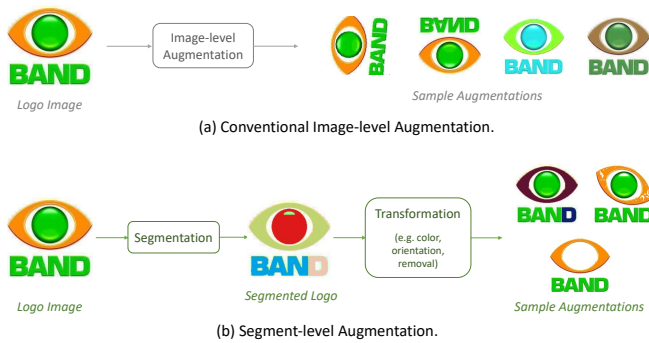


Fig. 1. (a) Conventional data augmentation approaches apply transformations at the image level. (b) We propose segment-level augmentation as a more suitable approach for problems like logo retrieval.

With the rapid increase in companies founded worldwide and the fierce competition among them globally, the identification of companies with their logos has become more pronounced, and it has become more paramount to check similarities between logos to prevent trademark infringements. Checking trademark infringements is generally performed manually by experts, which can be sub-optimal due to human-caused errors and time-consuming as it takes days to make a decision. Therefore, automatically identifying similar logos using content-based image processing techniques is crucial.

II. RELATED WORK

A. Logo Retrieval

Earlier studies in trademark retrieval [6] used hand-crafted features and deep features extracted using pre-trained networks and revealed that deep features obtained considerably better results. Perez *et al.* [5] improved the results by combining two CNNs trained on two different datasets. Later, Tursun *et al.* [1] achieved impressive results by introducing different attention methods to reduce the effect of text regions, and in their most recent work [4], they introduced different modifications and achieved state-of-the-art results.

B. Data Augmentation

Data augmentation [9], [10] is an essential and well-known technique in deep learning to make networks more robust to variations in data. Conventional augmentation methods perform geometric transformations such as zooming, flipping or cropping the entire image. Alternatively, adding noise, random erasing or synthesizing training data [11] are key approaches to improve overall model performance. Random Erasing [12] is a recently introduced method that obtains significant improvement on various recognition tasks. Although augmentation methods that focus on cutting and mixing windows [13]–[15] rather than the whole image are not widely used, they have shown significant gains in performance.

In logo retrieval, studies generally use conventional augmentation methods. For example, Tursun *et al.* [16] applied a reinforcement learning approach to learn an ensemble of test-time data augmentations for trademark retrieval. An exception to such an approach is the study by Tursun *et al.* [1], who proposed a method to remove text regions from logos while evaluating similarity.

C. Differentiable Ranking

Image or logo retrieval are by definition ranking problems, though ranking is not differentiable. To address this limitation, many solutions have been proposed recently [8], [17], [18]. These approaches mainly optimize Average Precision (AP) with different approximations: For example, Cakir *et al.* [17] quantize distances between pairs of instances and use differentiable relaxations for these quantized distances. Rolinek *et al.* [18] consider non-differentiable ranking as a black box and use smoothing to estimate suitable gradients for training a network to rank. Finally, Brown *et al.* [8] propose smoothing AP itself to use differentiable operations to train a deep network to rank.

These approximations have been mainly applied to standard retrieval benchmarks. In this paper, we show that differentiable ranking-based loss functions can lead to a performance improvement for logo retrieval as well.

D. Summary

Looking at the studies in the literature, we observe that **(1)** No study has performed segment-level augmentation either for logo retrieval or for general recognition or retrieval problems. The closest study for this research direction is the study by Tursun *et al.* [1], which just removed text regions in

logos while evaluating similarity. **(2)** Promising deep learning approaches such as metric learning using e.g. Triplet Loss and differentiable ranking have not been employed for logo retrieval.

III. METHOD

In this section, after providing a definition for logo retrieval, we present our novel segment-based augmentation method and how we use it with deep metric learning and differentiable ranking approaches.

A. Problem Definition

Given an input query logo I_q , logo retrieval aims to rank all logos in a retrieval set $\Omega = I_i, i = \{0, 1, \dots, N\}$, based on their similarity to the query I_q . To be able to evaluate retrieval performance and to train a deep network that relies on known similarities, we require for each I_q to have a set of positive (similar) logos, $\Omega^+(I_q)$, and a set of negative (dissimilar) logos, $\Omega^-(I_q)$. Note that logo retrieval defined as such does not have the notion of classes of a classification setting.

B. Segment-level Augmentation for Logo Retrieval

We perform segment-level augmentation by following these steps: (i) Logo segmentation, (ii) segment selection, and (iii) segment transformation. See Figure 2 for some samples.

1) *Logo Segmentation:* There are many sophisticated segmentation approaches available in the literature. Since logo images have relatively simpler regions compared to images, we observed that a simple and computationally-cheap approach using standard connected-component labeling is sufficient for extracting logo segments. See Figure 3 for some samples, and Supp. Mat. Section S4 for more samples and a discussion on the effect of segmentation quality.

2) *Segment Selection:* The next step is to select n random segments to apply transformations on them. Segment selection is a process that should be evaluated carefully since the number of segments or the area for each segment is not the same

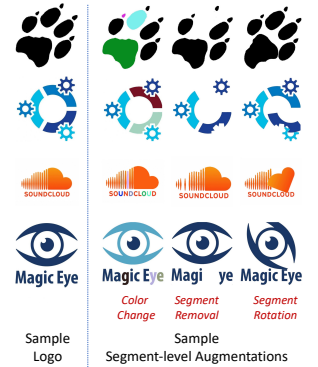


Fig. 2. Examples for our segment-level augmentation.



Fig. 3. Sample segmentation results. Each segment is represented with different color.

for each logo. Simplicity of logo instances also affects the number of available components and many logo instances have less than five components. Therefore, the choice of n can have drastic effects especially when the number of components in a logo is small, especially for the ‘segment removal’ transformation. For this reason, ‘segment removal’ is not applied to a segment with the largest area, and n is chosen to be small values. We present an ablation study to evaluate the effect of n on model performance for the introduced augmentation strategies. For the same reason, the background component is removed from available segments for augmentation.

3) *Segment Transformation*: For each selected segment S , the following are performed with probability p :

- ‘(Segment) Color change’: Every pixel in S is assigned to a randomly selected color.
- ‘Segment removal’: Pixel values in S are set to the same value of the background component.
- ‘Segment rotation’: We first select a segment and create a mask for the segment. The mask image and the corresponding segment pixels are rotated with a random angle in $[-90, 90]$. Then, the rotated segment is combined with the other segments. See also Figure 4 for an example.

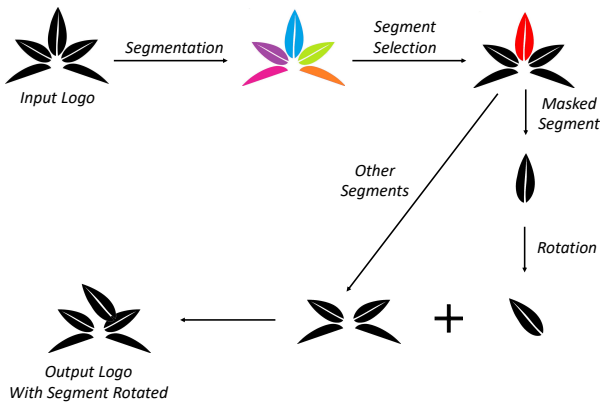


Fig. 4. The steps of rotating a segment.

See Figure 2 for some sample augmentations.

C. Adapting Ranking Losses for Logo Retrieval

1) *Mini-batch Sampling for Training*: For training the deep networks, we construct the batches as follows, similar to but different from [8] as we do not have classes: Each mini-batch B with size $|B|$ is constructed with two sub-sets: the similar set B^+ and the dissimilar set B^- . The similar logo set B^+ consists of logos that are known to be similar to each other (this information is available in the dataset [19]; logos with known similarities are provided as the ‘‘query set’’), and B^- contains logos that are dissimilar to the logos in B^+ (to be specific, logos other than the query set of the dataset are randomly sampled for B^-). The size of B^+ is set to 4, and that of B^- is $|B| - 4$. For training the network, every $I \in B^+$ has label as ‘‘1’’ and $I \in B^-$ has label ‘‘0’’.

2) *Smooth-AP Adaptation*: Smooth-AP [8] is a ranking-based differentiable loss function, approximating AP. The main aspect of this approximation is to replace discrete counting operation (the indicator function) in the non-differentiable AP with a Sigmoid function. Brown et al. [8] applied their study to standard retrieval benchmarks such as Stanford Online Products [20], VGGFace2 [21] and VehicleID [22]. However, the logo retrieval problem requires a dataset with a different structure as there is no notion of classes as in the Stanford Online Products [20], VGGFace2 [21] and VehicleID [22] datasets. Hence, Smooth-AP cannot be applied directly to our problem.

The first adaptation is about the structure of the mini-batch sampling. In Smooth-AP, Brown et al. explain their sampling as they have ‘‘formed each mini-batch by randomly sampling classes such that each represented class has P samples per class’’ [8]. Standard retrieval benchmarks have a notion of classes and are assumed to have sufficient instances per class to distribute among the mini-batches; however, there are not enough instances for known similarity ‘‘classes’’ in logo retrieval. This difference requires an adaption in both sampling and calculation of the loss. Smooth-AP Loss is calculated as follows [8]:

$$\mathcal{L}_{AP} = \frac{1}{C} \sum_{k=1}^C (1 - \tilde{A}P_k), \quad (1)$$

where C is the number of classes and $\tilde{A}P_k$ is the smoothed AP calculated for each class in the mini-batch with their Sigmoid-based smoothing method.

Our mini-batch sampling (Section III-C1) causes a natural contradiction because our batches only contain two classes: ‘‘similar’’ and ‘‘dissimilar’’; therewith the ‘‘dissimilar’’ class should not be included in the calculation of the loss. Dissimilar class instances have the same label (‘‘0’’), but that does not mean they have the same class; they are just not similar to the similar logo set B^+ in the mini-batch. Hence, the ranking among B^- does not matter in our case. This difference in the batch construction and the notion of classes lead to our second adaptation. In this adaptation, the only calculated AP approximation belongs to the known ‘‘similar’’ class (logos in B^+). Therefore, the loss calculation becomes:

$$\mathcal{L}_{AP}^+ = 1 - \tilde{A}P_+, \quad (2)$$

where $\tilde{A}P$ is calculated (approximated) in the same way as in the original paper [8].

3) *Triplet Loss Adaptation*: Triplet Loss [7] is a well-known loss function used in many computer vision problems. Triplet Loss is differentiable, but, unlike Smooth-AP Loss [8], rather than optimizing ranking, it optimizes the distances between positive pairs and negative pairs of instances. In this paper, for each mini-batch, triplets consist of one ‘‘anchor’’ instance, one positive instance, and one negative instance. For the same reasons discussed in Smooth-AP Loss [8], only the instances of known similarity classes can be used as the anchor instance. Optimizing the distances between dissimilar logo instances is not sensible because, as discussed in the previous

section, instances of the dissimilar logos do not have any known similarity between them. Thus, triplet loss calculation is limited to the triplets that contain known similar instances as the “anchor” instance.

IV. EXPERIMENTS AND RESULTS

We now evaluate the performance of the proposed segment-augmentation strategy and its use with Triplet Loss and Smooth-AP Loss.

A. Experimental and Implementation Details

1) *Dataset*: We use the METU Dataset [19], which is one of the largest publicly available logo retrieval datasets. The dataset is composed of more than 900K authentic logos belonging to actual companies worldwide. Moreover, it includes query sets, i.e. similar logos, of varying difficulties, allowing logo retrieval researchers to benchmark their methods against other methods. We have used 411K training images, 413K test images, and 418 query images.

2) *Training and Implementation Details*: For every experiment that will be discussed, we use ImageNet [11] pre-trained ResNet50 [23] as our backbone architecture which has a linear layer with 512 dimensions, rather than a final Softmax layer. We use the Adam optimizer with the hyper-parameters tuned as 10^{-7} for the learning rate and 256 for the batch size.

3) *Evaluation Measures*: Following the earlier studies [1], [19], we use Normalized Average Rank (NAR) and Recall@K for quantifying the performance of the methods. NAR is calculated as:

$$NAR = \frac{1}{N \times N_{rel}} \left(\sum_{i=1}^{N_{rel}} R_i - \frac{N_{rel}(N_{rel} + 1)}{2} \right), \quad (3)$$

where N_{rel} is the number of similar images for a particular query image; N is the size of the image set; and R_i is the rank of the i^{th} similar image. NAR lies in the range $[0, 1]$, where 0 denotes the perfect score, and 1 the worst. Recall@K (R@K) is recall for top-K similar logos.

TABLE I

THE EFFECT OF USING TRIPLET LOSS AND SMOOTH-AP LOSS FOR LOGO RETRIEVAL. NEITHER IMAGE-LEVEL NOR SEGMENT-LEVEL AUGMENTATION IS USED FOR ANY METHOD IN THIS TABLE.

Method	NAR ↓	Recall@1 ↑	Recall@8 ↑
Baseline	0.102	0.310	0.536
Triplet Loss	0.053	0.344	0.586
Smooth-AP Loss	0.046	0.339	0.581

B. Experiment 1: Effect of Ranking Losses

Before analyzing the effect of segment-level augmentation, in this section, we first provide a stand-alone analysis to illustrate the effect of the ranking losses. We compare Triplet Loss and Smooth-AP Loss with a baseline that compares features extracted with the pre-trained Resnet50 backbone using Cosine Similarity. For this analysis, no image-level or

segment-level augmentations are used, except for the Random Resized Crop to fit the images to the expected resolution of the network, i.e. 224×224 .

The results in Table I suggest that both loss adaptations provide a significant performance improvement in both NAR and Recall measures and Smooth-AP adaptation achieves the best performance. Applying Cosine Similarity on off-the-shelf ResNet50 features shows adequate results in no-text logo instances, however, it performs worse on logos with text (see Appendix E).

It is evident that the improvement in Recall is not as visible as NAR. This difference states that the adapted loss functions highly affect the overall rankings of the similar known instances. However, these effects are not completely reflected by the Recall because of the selected $K=8$ value.

TABLE II

THE EFFECT OF IMAGE-LEVEL (H. FLIP, V. FLIP) AND SEGMENT-LEVEL AUGMENTATION. ONLY THE BEST AUGMENTATION STRATEGIES ARE REPORTED. SEE SECTION IV-E FOR AN ABLATION ANALYSIS.

Method	NAR ↓	Recall@1 ↑	Recall@8 ↑
Baseline (No augmentation)	0.102	0.310	0.536
Triplet Loss (No augmentation)	0.053	0.344	0.586
Triplet Loss (Image-level aug.)	0.051	0.354	0.596
Triplet Loss (S. Color, S. Removal)	0.046	0.374	0.640
Smooth-AP Loss (No augmentation)	0.046	0.339	0.581
Smooth-AP Loss (Image-level aug.)	0.044	0.339	0.596
Smooth-AP Loss (S. Color)	0.040	0.354	0.610

C. Experiment 2: Effect of Segment Augmentation

We now compare our segment-based augmentation methods with the conventional image-level augmentation techniques on the METU dataset [19]. In every experiment, we resize the images with Random Resized Crop to fit them to the expected resolution of the network, i.e. 224×224 . For both segment-based and image-level augmentation, the same number of images are augmented and the probability $p = 0.5$ is used for selecting a certain transformation.

We have provided a comparison between the best resulting methods for both image-level and segment-level augmentation methods in Table II. We see that image-level augmentation can improve ranking performance. However, the results suggest that segment-level augmentation provides a significantly better gain both in terms of NAR and R@K measures. Detailed comparison between image-level and segment-level methods is provided in ablation study.

TABLE III

NORMALIZED AVERAGE RANK (NAR) VALUES FOR PREVIOUS STATE-OF-THE-ART RESULTS ON THE METU DATASET. THE RESULTS ARE NOT COMPARABLE AS THE METHODS DIFFER IN THEIR BACKBONES, TRAINING DATASETS, OR TRAINING REGIME. FOR SOME METHODS, THESE DETAILS ARE NOT EVEN REPORTED. SEE THE TEXT FOR DETAILS.

Method	NAR ↓
Hand-crafted Features (Feng <i>et al.</i> [2])	0.083
Hand-crafted Features (Tursun <i>et al.</i> [19])	0.062
Off-the-shelf Deep Features (Tursun <i>et al.</i> [6])	0.086
Transfer Learning (Perez <i>et al.</i> [5])	0.047
Component-based attention (SPoC [24], [1])	0.120
Component-based attention (CRoW [25], [1])	0.140
Component-based attention (R-MAC [26], [1])	0.072
Component-based attention (MAC [26] [1])	0.120
Component-based attention (Jimenez [27], [1])	0.093
Component-based attention (CAM MAC [1])	0.064
Component-based attention (ATR MAC [1])	0.056
Component-based attention (ATR R-MAC [1])	0.063
Component-based attention (ATR CAM MAC [1])	0.040
MR-R-MAC w/UAR (Tursun <i>et al.</i> [4])	0.028
Segment-Augm. (Color Change) w Smooth-AP (Ours)	0.040

D. Experiment 3: Comparison with State of the Art

We compare our method and the state-of-the-art methods on the METU dataset [19]. It is important to note that a fair comparison between the methods is not possible because they differ in their backbones, training datasets or dataset splits, or training time. For some papers, even those details are missing; e.g. for [4], which reports the best NAR performance. Therefore, we list the results in Table III, and refrain from drawing conclusions.

TABLE IV

THE EFFECT OF PROBABILITY p FOR AUGMENTING A SELECTED SEGMENT, WITH SMOOTH-AP LOSS.

Color C.	Rotation	Removal	p	NAR ↓	R@1 ↑	R@8 ↑
	<i>Baseline</i>		0	0.046	0.339	0.581
✓			0.5	0.040	0.354	0.610
✓			0.75	0.043	0.354	0.601
✓			1.0	0.049	0.325	0.566
	✓		0.5	0.048	0.344	0.601
	✓		0.75	0.049	0.369	0.591
	✓		1.0	0.060	0.330	0.556
		✓	0.5	0.048	0.339	0.596
		✓	0.75	0.047	0.344	0.571
		✓	1.0	0.047	0.344	0.591

E. Experiment 4: Ablation Study

1) *Choice of Hyper-Parameters:* Our segment-level augmentation has two hyper-parameters: The number of segments, n , selected for augmentation, and the probability, p , of applying a selected augmentation. Table IV shows that the best performance is obtained with p as 0.5. A similar analysis for n (with values 1, 2, $L/3$ and $L/2$ where L is the number of segments in a logo) provided the best performance for n as $L/3$.

2) *Effects of Individual Augmentation Methods:* Tables V (Smooth-AP Loss) and VI (Triplet Loss) list the effects of both image-level and segment-level augmentation. The tables show that, among the segment-level augmentation methods, (Segment) ‘Color Change’ outperforms the others for both loss functions. With Triplet Loss adaptation, (Segment) ‘Removal’ and ‘Rotation’ provide slightly better NAR values than the baseline. Another point worth mentioning is that combining (Segment) ‘Rotation’ or ‘Removal’ degrades the NAR performance measure whereas the combination of (Segment) ‘Removal’ and ‘Color Change’ yields the best result at *Recall@8*.

3) *Experiments with a Different Backbone:* Section A in the Appendix provides an analysis using ConvNeXt [28], a recent, fast and strong backbone competing with transformer-based architectures. Our results without any hyper-parameter tuning are comparable to the baseline or better than the baseline with the R@8 measure.

4) *Experiments with a Different Dataset:* Section B in the Appendix reports results on the LLD dataset [29] that confirm our analysis on the METU dataset: We observe that segment-level augmentation provides significant gains for all measures.

F. Experiment 5: Visual Results

Section F in the Appendix provides sample retrieval results for several query logos for the baseline as well the adaptations of Triplet Loss and Smooth-AP Loss with our segment-level augmentation methods. The visual results also confirm that segment augmentation with our Smooth-AP adaptation performs best.

V. CONCLUSION

We introduced a novel data augmentation method based on image segments for training neural networks for logo retrieval. We performed segment-level augmentation by identifying segments in a logo and do transformations on selected segments. Experiments were conducted on the METU [19] and LLD [29] datasets with ResNet [23] and ConvNeXt [28] backbones and suggest significant improvements on two evaluation measures of ranking performance. Moreover, we use metric learning and differentiable ranking with the proposed segment-augmentation method to demonstrate that our method can lead to a further boost in ranking performance.

We note that our segment-level augmentation strategy generates similarities between logos that are rather simplistic: It is based on the assumption that two similar logos differ from each other in terms of certain segments having differences in color, orientation and presence. An important research

TABLE V
NORMALIZED AVERAGE RANK (NAR) AND RECALL@K VALUES FOR DATA AUGMENTATION EXPERIMENTS WITH SMOOTH-AP LOSS.

Image Level					Segment Level			NAR ↓	R@1 ↑	R@8 ↑
Resized Crop	Hor. Flip	Vert. Flip	Rotation	Color Jitter	S. Color Change	S. Rotation	S. Removal			
<i>Baseline</i>								0.102	0.310	0.536
✓								0.046	0.339	0.581
✓	✓							0.049	0.325	0.591
✓		✓						0.044	0.325	0.596
✓	✓	✓						0.044	0.339	0.596
✓			✓					0.045	0.344	0.551
✓				✓				0.049	0.349	0.576
✓					✓			0.040	0.354	0.610
✓						✓		0.048	0.344	0.601
✓							✓	0.048	0.339	0.596
✓					✓	✓		0.050	0.354	0.586
✓					✓		✓	0.046	0.374	0.625
✓						✓	✓	0.044	0.354	0.591
✓					✓	✓	✓	0.047	0.374	0.605

TABLE VI
NORMALIZED AVERAGE RANK (NAR) AND RECALL@K VALUES FOR DATA AUGMENTATION EXPERIMENTS WITH TRIPLET LOSS.

Image-Level					Segment-Level			NAR ↓	R@1 ↑	R@8 ↑
Resized Crop	Hor. Flip	Vert. Flip	Rotation	Color Jitter	S. Color Change	S. Rotation	S. Removal			
<i>Baseline</i>								0.102	0.310	0.536
✓								0.053	0.339	0.586
✓	✓							0.052	0.349	0.586
✓		✓						0.052	0.364	0.586
✓	✓	✓						0.051	0.354	0.596
✓			✓					0.054	0.320	0.546
✓				✓				0.053	0.344	0.591
✓					✓			0.048	0.369	0.601
✓						✓		0.052	0.354	0.596
✓							✓	0.052	0.354	0.596
✓					✓	✓		0.050	0.354	0.586
✓					✓		✓	0.046	0.374	0.640
✓						✓	✓	0.048	0.354	0.571
✓					✓	✓	✓	0.046	0.369	0.581

direction is exploring more sophisticated augmentation strategies for introducing artificial similarities. However, our results suggest that even such a simplistic strategy can improve the retrieval performance significantly and therefore, our study can be considered as a first step towards developing better segment/part-level augmentation strategies.

REFERENCES

- [1] O. Tursun, S. Denman, S. Sivapalan, S. Sridharan, C. Fookes, and S. Mau, "Component-based attention for large-scale trademark retrieval," *IEEE Transactions on Information Forensics and Security*, 2019.
- [2] Y. Feng, C. Shi, C. Qi, J. Xu, B. Xiao, and C. Wang, "Aggregation of reversal invariant features from edge images for large-scale trademark retrieval," in *4th International Conference on Control, Automation and Robotics (ICCAR)*, 2018, pp. 384–388.
- [3] J. Cao, Y. Huang, Q. Dai, and W.-K. Ling, "Unsupervised trademark retrieval method based on attention mechanism," *Sensors*, vol. 21, no. 5, 2021.
- [4] O. Tursun, S. Denman, S. Sridharan, and C. Fookes, "Learning regional attention over multi-resolution deep convolutional features for trademark retrieval," in *IEEE International Conference on Image Processing (ICIP)*, 2021, pp. 2393–2397.
- [5] C. A. Perez, P. A. Estévez, F. J. Galdames, D. A. Schulz, J. P. Perez, D. Bastias, and D. R. Vilar, "Trademark image retrieval using a combination of deep convolutional neural networks," in *International Joint Conference on Neural Networks (IJCNN)*, 2018, pp. 1–7.
- [6] C. Aker, O. Tursun, and S. Kalkan, "Analyzing deep features for trademark retrieval," in *25th Signal Processing and Communications Applications Conference (SIU)*, 2017, pp. 1–4.
- [7] K. Q. Weinberger and L. K. Saul, "Distance metric learning for large margin nearest neighbor classification," *Journal of Machine Learning Research*, vol. 10, no. 2, 2009.
- [8] A. Brown, W. Xie, V. Kalogeiton, and A. Zisserman, "Smooth-ap: Smoothing the path towards large-scale image retrieval," in *European Conference on Computer Vision (ECCV)*, 2020.
- [9] A. Mikołajczyk and M. Grochowski, "Data augmentation for improving deep learning in image classification problem," in *International Interdisciplinary PhD Workshop (IIPhDW)*, 2018, pp. 117–122.
- [10] C. Shorten and T. Khoshgoftaar, "A survey on image data augmentation for deep learning," *Journal of Big Data*, vol. 6, 07 2019.
- [11] O. Russakovsky, J. Deng, H. Su, J. Krause, S. Satheesh, S. Ma, Z. Huang, A. Karpathy, A. Khosla, M. Bernstein *et al.*, "Imagenet large scale visual recognition challenge," *International journal of computer vision*, vol. 115, no. 3, pp. 211–252, 2015.
- [12] Z. Zhong, L. Zheng, G. Kang, S. Li, and Y. Yang, "Random erasing data augmentation," *AAAI Conference on Artificial Intelligence (AAAI)*, vol. 34, 08 2017.
- [13] S. Yun, D. Han, S. Chun, S. Oh, Y. Yoo, and J. Choe, "Cutmix: Regularization strategy to train strong classifiers with localizable features," in *IEEE/CVF International Conference on Computer Vision (ICCV)*, 2019, pp. 6022–6031.
- [14] H. Zhang, M. Cissé, Y. N. Dauphin, and D. Lopez-Paz, "mixup: Beyond empirical risk minimization," *International Conference on Learning Representations (ICLR)*, 2018.
- [15] Y. Tokozume, Y. Ushiku, and T. Harada, "Between-class learning for image classification," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2018, pp. 5486–5494.
- [16] O. Tursun, S. Denman, S. Sridharan, and C. Fookes, "Learning test-time augmentation for content-based image retrieval," *arXiv preprint arXiv:2002.01642*, 2020.
- [17] F. Cakir, K. He, X. Xia, B. Kulis, and S. Sclaroff, "Deep metric learning to rank," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2019.
- [18] M. Rolinek, V. Musil, A. Paulus, M. Vlastelica, C. Michaelis, and G. Martius, "Optimizing rank-based metrics with blackbox differentiation," in *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2020.
- [19] O. Tursun, C. Aker, and S. Kalkan, "A large-scale dataset and benchmark for similar trademark retrieval," *CoRR*, vol. abs/1701.05766, 2017.
- [20] H. O. Song, Y. Xiang, S. Jegelka, and S. Savarese, "Deep metric learning via lifted structured feature embedding," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [21] Q. Cao, L. Shen, W. Xie, O. M. Parkhi, and A. Zisserman, "Vggface2: A dataset for recognising faces across pose and age," in *IEEE International Conference on Automatic Face & Gesture Recognition (FG)*, 2018.
- [22] H. Liu, Y. Tian, Y. Wang, L. Pang, and T. Huang, "Deep relative distance learning: Tell the difference between similar vehicles," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 2167–2175.
- [23] K. He, X. Zhang, S. Ren, and J. Sun, "Deep residual learning for image recognition," in *IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2016.
- [24] A. B. Yandex and V. Lempitsky, "Aggregating local deep features for image retrieval," in *IEEE International Conference on Computer Vision (ICCV)*, 2015.
- [25] Y. Kalantidis, C. Mellina, and S. Osindero, "Cross-dimensional weighting for aggregated deep convolutional features," in *European Conference on Computer Vision (ECCV)*, 2016.
- [26] G. Toliás, R. Sivic, and H. Jégou, "Particular object retrieval with integral max-pooling of CNN activations," in *International Conference on Learning Representations (ICLR)*, Y. Bengio and Y. LeCun, Eds., 2016.
- [27] A. Jimenez, J. M. Alvarez, and X. Giró-i-Nieto, "Class-weighted convolutional features for visual instance search," *CoRR*, vol. abs/1707.02581, 2017.
- [28] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, "A convnet for the 2020s," *IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*, 2022.
- [29] A. Sage, E. Agustsson, R. Timofte, and L. Van Gool, "Lld - large logo dataset - version 0.1," <https://data.vision.ee.ethz.ch/cvl/lld>, 2017.
- [30] A. Buslaev, V. I. Iglovikov, E. Khvedchenya, A. Parinov, M. Druzhinin, and A. A. Kalinin, "Albumentations: Fast and flexible image augmentations," *Information*, vol. 11, no. 2, 2020. [Online]. Available: <https://www.mdpi.com/2078-2489/11/2/125>
- [31] O. Tursun, S. Denman, S. Sridharan, and C. Fookes, "Learning test-time augmentation for content-based image retrieval," *arXiv preprint arXiv:2002.01642*, 2020.
- [32] Y. Zhu, A. Fathi, and L. Fei-Fei, "Reasoning about object affordances in a knowledge base representation," in *European Conference on Computer Vision (ECCV)*. Springer, 2014, pp. 408–424.
- [33] A. Myers, C. L. Teo, C. Fermüller, and Y. Aloimonos, "Affordance detection of tool parts from geometric features," in *IEEE International Conference on Robotics and Automation (ICRA)*, 2015.
- [34] A. Myers, A. Kanazawa, C. Fermüller, and Y. Aloimonos, "Affordance of object parts from geometric features," in *Workshop on Vision meets Cognition, CVPR*, vol. 9, 2014.
- [35] K. Leonard, G. Morin, S. Hahmann, and A. Carlier, "A 2d shape structure for decomposition and part similarity," in *23rd International Conference on Pattern Recognition (ICPR)*, 2016.
- [36] L. J. Latecki and R. Lakamper, "Shape similarity measure based on correspondence of visual parts," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 22, no. 10, pp. 1185–1190, 2000.

APPENDIX

A. Experiments with ConvNeXt, a Stronger Backbone

In Table VII, we evaluate our novel segment-level augmentation strategy using ConvNeXt [28], a stronger, recent, fast backbone that can compete with transformer-based architectures. Although ConvNeXt provides significantly better performance than the ResNet baseline, Table VII confirms our experiments with ResNet that segment-level augmentation can improve logo retrieval performance.

TABLE VII
EXPERIMENTS ON THE CONVNEXT ARCHITECTURE [28] USING
SMOOTH-AP LOSS.

Color C.	Rotation	Removal	p	NAR ↓	R@1 ↑	R@8 ↑
<i>Baseline</i>			0	0.039	0.438	0.704
✓			0.5	0.039	0.413	0.709
	✓		0.5	0.041	0.413	0.694
		✓	0.5	0.040	0.403	0.699
✓	✓		0.5	0.046	0.379	0.661
✓		✓	0.5	0.045	0.374	0.684
	✓	✓	0.5	0.051	0.369	0.655
✓	✓	✓	0.5	0.055	0.354	0.615

B. Experiments with a Different Dataset

Without tuning, we repeat our experiments on using a different logo retrieval dataset, namely the Large Logo Dataset (LLD) [29]. LLD has 61380 logos for training and 61540 logos for testing. The LLD dataset does not have a query set for providing known similarities and therefore, for the experiments on LLD, we use the query set of the METU dataset to find similarities in the LLD dataset.

The results in Table VIII show that segment-level augmentation performs better than image-level augmentation in R@1 measure and on par in terms of NAR and R@8 measures. Considering that LLD is a smaller dataset than METU, we believe that segment-level augmentation can provide a larger margin in performance if tuned.

C. Comparisons with Different Image-level Augmentations

We now compare segment-level augmentation with more image-level augmentation methods. We have selected most commonly used image-level augmentation methods (Cutout, Elastic Transform and Channel Shuffle) from [30]. For this experiment, we use the same selection probability and setup as in Section IV.C. of the main paper. The results in Table IX show that, **without any tuning**, segment-level augmentation performs better than this new set of image-level augmentations in terms of NAR and R@8 measures whereas is inferior in terms of R@1 measure. With tuning, segment-level augmentation has potential to perform better.

D. An Analysis of Segmentation Maps

Logos are simplistic images composed of regions that are generally homogeneous in color. Therefore, an off-the-shelf simple segmentation algorithm works generally well for most logos (see Figure 5 for some examples). Logo segmentation can produce spurious segments especially on regions which have strong color gradients. However, this is not a problem for us because segments corresponding to over-segmentation or under-segmentation do function as a form of segment-level augmentation and this is still useful for training.

Figure 5 displays edge examples where our segment-level augmentation produces drastically different logos for which computing similarity with the original logos is highly challenging. These cases happen if there are only a few segments in a logo with comparable size and one of them is removed, or if the segmentation method under-segments the image and a segment that groups multiple regions is removed, or if the background segment is selected for augmentation and rotated. We did not try to address these edge cases as they are not frequent. We believe that they are useful stochastic perturbations and helpful to the training dynamics. We leave an analysis of this and improving the quality of segmentation & augmentation as future work.



Fig. 5. Edge cases in our segment-level augmentation. The first is a result of under-segmentation where the connected petals are with similar color are grouped in a single segment, and therefore, changing its color produces a distinct logo. The second case happens because the logo has two main segments and removing one results in a drastically different logo. The third is a case where a background segment is rotated and the outcome is an overlay of two segments.

E. The Effect of Text in Logos

In this experiment, we analyze the effect of text on logo similarity. The visual results in Table X and the quantitative analysis in Table XI show that, when used as queries, logos with text are more difficult to match to the logos in the dataset and the best and average ranks of similar logos are very high. This is not surprising since representations in deep networks do capture text as well and unless additional mechanisms are used to ignore them, they are taken into account while measuring similarity. The effect of text on logo retrieval is

TABLE VIII
NORMALIZED AVERAGE RANK (NAR) AND RECALL@K VALUES FOR DATA AUGMENTATION EXPERIMENTS FROM ALBUMENTATIONS LIBRARY WITH RESNET ARCHITECTURE ON LLD DATASET.

Image Level					Segment Level			NAR ↓	R@1 ↑	R@8 ↑
Resized Crop	Random Rotation	R. Horizontal Flip	Color Jitter	S. Color Change	S. Rotation	S. Removal				
✓							0.044	0.556	0.724	
✓	✓						0.068	0.527	0.679	
✓		✓					0.038	0.551	0.738	
✓			✓				0.039	0.551	0.719	
✓				✓			0.039	0.581	0.738	
✓					✓		0.058	0.527	0.687	
✓						✓	0.051	0.532	0.704	
✓				✓	✓		0.054	0.561	0.719	
✓				✓		✓	0.056	0.517	0.704	
✓					✓	✓	0.068	0.507	0.684	
✓				✓	✓	✓	0.064	0.536	0.684	

TABLE IX
NORMALIZED AVERAGE RANK (NAR) AND RECALL@K VALUES FOR DATA AUGMENTATION EXPERIMENTS FROM ALBUMENTATIONS LIBRARY WITH SMOOTH-AP LOSS ON RESNET.

Image Level					Segment Level			NAR ↓	R@1 ↑	R@8 ↑
Resized Crop	Cutout	Elastic Transform	Channel Shuffle	S. Color Change	S. Rotation	S. Removal				
							<i>Baseline</i>	0.102	0.310	0.536
✓	✓							0.052	0.325	0.581
✓		✓						0.060	0.379	0.596
✓			✓					0.0451	0.344	0.596
✓				✓				0.040	0.354	0.610
✓					✓			0.048	0.344	0.601
✓						✓		0.048	0.339	0.596
✓				✓	✓			0.050	0.354	0.586
✓				✓		✓		0.046	0.374	0.625
✓					✓	✓		0.044	0.354	0.591
✓				✓	✓	✓		0.047	0.374	0.605

already known in the literature and soft or hard mechanisms can be used to remove them from logos – see e.g. [19], [31].

F. Visual Retrieval Results

Figure 6 displays sample retrieval results for different queries. We see in Figure 6(a) and (b) that segment-level augmentation is able to provide better retrieval than its competitors. Figure 6(c) displays a failure case where a query with an unusual color distribution is used. We see that all methods are adversely affected by this; however, segment-level augmentation is able to retrieve logos with similar color distributions.

G. Running-Time Analysis

Table XII provides a running-time analysis of segment-level and image-level augmentation strategies. We observe that the time spent on segmentation (0.4ms) and Segment Color Change (2.4ms) is comparable to image-level transformation Horizontal Flip (2.8ms). However, Segment Removal and Segment Rotation takes significantly more time. It is important to that our implementation is not optimized for efficiency.

H. Discussion

1) *Comparing Triplet Loss and Smooth-AP Loss:* Our results on the METU and LLD datasets lead to two interesting

TABLE X

EFFECT OF TEXT ON LOGO RETRIEVAL PERFORMANCE. WHEN USED AS QUERIES, LOGOS WITH TEXT ARE MORE DIFFICULT TO MATCH TO THE LOGOS IN THE DATASET AND THE BEST AND AVERAGE RANKS OF SIMILAR LOGOS ARE VERY LOW.









Query								
Best Rank	39k	2	25	5	97k	1	159k	1
Average Rank	253k	2605	53k	10k	239k	25k	239k	1263

TABLE XI

EFFECT OF TEXT ON LOGO RETRIEVAL PERFORMANCE USING THE 213 QUERY LOGOS IN THE METU TRADEMARK DATASET.

Logo Type	Number of Logos	NAR
Logos w/o text	117	0.015
Logos w text	86	0.071

TABLE XII

TIME SPENT ON DIFFERENT STEPS OF SEGMENT-LEVEL AND IMAGE-LEVEL AUGMENTATION.

Method	Segmentation	Transformation
Color Jitter	-	1.2ms
R. Horizontal Flip	-	2.8ms
Segment Color	0.4ms	2.4ms
Segment Removal	0.4ms	4.3ms
Segment Rotation	0.4ms	14.0ms

findings, which we discuss below:

- *Finding 1: Triplet Loss is better in terms of R@1 and R@8 measures whereas Smooth-AP is better in terms of NAR measure.*
- *Finding 2: Smooth-AP provides its best NAR performance with Color Change whereas Triplet Loss provides its best with Color Change & Segment Removal.*

Triplet Loss is by definition optimizing or learning a distance metric and considered a surrogate ranking objective. On the other hand, Smooth-AP directly aims to optimize Average Precision, a ranking measure, and therefore, it pertains to a more global ranking objective than Triplet Loss, which works at the level of triplets only. See also Brown et al. [8] who discussed and contrasted Triplet Loss with Smooth-AP Loss.

The two objectives provides different inductive biases to the learning mechanism and therefore, we see differences in terms of their performances with respect to different performance measures and augmentation strategies. For example, Finding 1 is likely to be because NAR considers all logos whereas R@1 and R8 only consider logos on the top of the ranking which can be easily learned to be ranked at the top using

local similarity arrangements as in Triplet Loss. Moreover, Finding 2 occurs because different augmentation strategies incur different ranking among logos and the inductive biases of Triplet Loss and Smooth-AP handle them differently.

We believe that this is an interesting research question that we leave as future work.

2) *Different Effects and Selection Probabilities for Segment-Augmentations:* The results in Table IV of the main paper show that different augmentations contribute to performance differently and with different selection probabilities. For example, we see that Color Change (with $p = 0.5$) gives the best performance in terms of NAR and R@8 whereas Rotation (with $p = 0.75$) provides the best performance for R@1 and Removal improves over the baseline only in terms of R@8 measure.

We attribute these differences to the fact that the different segment-level augmentations incur different biases: Color Change enforces invariance to perturbations in color differences at the segment-level whereas Segment Rotation and Removal encourage invariance to changes to the spatial layout of the shape.

3) *Applicability to Other Problems:* We agree that our analysis is limited to logo retrieval. However, the idea of segment-level augmentation is a viable approach for reasoning about similarities at object part levels and require transfer of knowledge at the level of object parts. One good example is reasoning about affordances of objects [32]–[34], where supported functions of object parts can be transferred across objects having similar parts. Another example is reasoning about similarity between shapes that have partial overlap [35], [36], where correspondences between parts of shapes need to be calculated. In either example, the specific segment-level augmentation methods may have to be adjusted to the specific problem. For example, performing affine transformations on the segments may be helpful for problems with real-world

Method	Query	Top Retrievals (Yellow bar: Relevant. Red bar: Irrelevant)
Baseline		
Triplet Loss (Segment Level - color + removal)		
Smooth AP Loss (Image Level - Channel Shuffle)		
Smooth AP Loss (Segment Level - Color Change)		

(a)

Method	Query	Top Retrievals (Yellow bar: Relevant. Red bar: Irrelevant)
Baseline		
Triplet Loss (Segment Level - color + removal)		
Smooth AP Loss (Image Level - Channel Shuffle)		
Smooth AP Loss (Segment Level - Color Change)		

(b)

Method	Query	Top Retrievals (Yellow bar: Relevant. Red bar: Irrelevant)
Baseline		
Triplet Loss (Segment Level - color + removal)		
Smooth AP Loss (Image Level - Channel Shuffle)		
Smooth AP Loss (Segment Level - Color Change)		

(c)

Fig. 6. Visual results on the METU dataset for the methods at their best settings. (a-b) Example cases where segment-level augmentation produces better retrieval than image-level augmentation. (c) A negative result for a query with an unusual intensity distribution, which affects all methods adversely. Though we observe that segment-level augmentation is able to retrieve logos with similar color distributions.