



# FeatureNMS: Non-Maximum Suppression by Learning Feature Embeddings

Niels Ole Salscheider<sup>[0000–0002–0886–0000]</sup>

FZI Forschungszentrum Informatik, Haid-und-Neu-Str. 10-14, 76131 Karlsruhe  
[salscheider@fzi.de](mailto:salscheider@fzi.de)

**Abstract.** Most state of the art object detectors output multiple detections per object. The duplicates are removed in a post-processing step called *Non-Maximum Suppression*. Classical Non-Maximum Suppression has shortcomings in scenes that contain objects with high overlap: The idea of this heuristic is that a high bounding box overlap corresponds to a high probability of having a duplicate. We propose FeatureNMS to solve this problem. FeatureNMS recognizes duplicates not only based on the intersection over union between bounding boxes, but also based on the difference of feature vectors. These feature vectors can encode more information like visual appearance. Our approach outperforms classical NMS and derived approaches and achieves state of the art performance.

**Keywords:** Non-Maximum Suppression, Object Detection, Feature Embeddings, Computer Vision

## 1 Introduction

Object detection is an important task in a huge variety of applications. In some of these applications, images can contain a lot of partially overlapping objects. One example are images of traffic scenes that contain crowds of humans. This scenario is common in autonomous driving or surveillance scenarios.

Most state of the art object detectors are based on Convolutional Neural Networks (CNN). There are single-stage detectors like YOLO [1,2,3], SSD [4] and RetinaNet [5] and two-stage detectors like R-CNN [6], Fast R-CNN [7] and Faster R-CNN [8]. Two-stage detectors first generate proposals and then a dedicated second stage decides if the proposal is in fact an object of interest. Single-stage detectors on the other hand directly perform object detection on the input image.

Both approaches have in common that they usually generate multiple detections per object. Duplicate detections are then removed in a postprocessing step called *Non-Maximum Suppression* (NMS). The widely used classical approach is a greedy heuristic. Detections are sorted by their scores in a decreasing order. Then each detection is checked against all following in the sorted list. If the *Intersection over Union* (IoU) with one of the following detections is larger than a certain threshold this detection is removed.

### 1.1 Related Work: Non-Maximum Suppression

There have been several proposals how to improve this heuristic. SoftNMS [9] does not remove overlapping detections but decreases the detection scores of duplicates. The factor by which it is decreased is a function of the IoU of the corresponding bounding boxes.

The idea of AdaptiveNMS [10] is to adjust the threshold for the greedy heuristic based on the detection density. The density for each detection is predicted by the neural network used for object detection.

Visibility Guided NMS [11] uses another approach. The detection network outputs two bounding boxes per object. One bounding box encloses the whole object while the other encloses only the visible part. The IoU of bounding boxes from different objects is usually smaller for the visible part. Because of that, classical NMS is performed on these bounding boxes. But the final output are the corresponding bounding boxes for the whole object.

Other works [12,13] try to work around the shortcomings of classical NMS in the detector training. The idea is to push bounding boxes of different objects far enough apart while achieving as much overlap as possible for detections that belong to the same object. This makes the task of NMS easier since the detections violate the assumptions of classical NMS less.

In [14], the authors propose to solve the NMS task with a CNN. The proposed network learns to rescore detections to suppress duplicates. In order to achieve this each block in the network has access to pairwise features of detections. These features include the IoU of both bounding boxes, normalized distances, as well as scale and aspect ratio differences.

Relation Networks [15] add a relation module to the detection network. This relation module learns to perform NMS inside the network. For this it can use geometric and appearance features of the detections.

### 1.2 Related Work: Embedding Learning

Learning of embeddings is used in a wide range of applications like zero-shot learning [16], visual search [17,18,19] or image comparison [20,21]. The underlying idea is conceptionally simple: The embedding vectors of positive image pairs (i. e. images that show the same object) should be similar. Embedding vectors of negative pairs on the other hand should be separated by a certain distance.

There are several loss functions that can be used to achieve this objective. Contrastive loss [22] is a widely used for this purpose. It consists of two terms: One term pulls the  $L^2$  distance of positive pairs as close to zero as possible. The other term pushes the  $L^2$  distance of negative pairs apart if it is below a certain margin.

Choosing the margin parameter correctly can be challenging. For hard negative examples it might be too difficult to push the embeddings far enough apart while keeping small distances for positive pairs. Triplet loss [23] tries to solve this problem by using triplets of images: An anchor, a positive example and a negative example. It tries to ensure that the embedding of the anchor is closer



to all positive examples than any negative example. The authors also propose a sampling strategy to select suitable triplets for training.

Recently, Margin loss [24] has been proposed as an alternative to contrastive loss. It does try to push the embeddings of all positive pairs to be as close to each other as possible. Instead it just requires that the distance is lower than a certain margin, making the loss more robust. Together with a distance weighted sampling strategy it achieves state-of-the-art performance on multiple tasks.

The remainder of this paper is structured as follows: Section 2 describes our proposed approach to Non-Maximum Suppression. The general idea is presented in Section 2.1 while Section 2.2 contains details about the necessary modifications to the object detector network. In Section 3 we present our evaluation procedure and the results. Finally, Section 4 concludes the paper.

## 2 Approach

We first describe our proposed approach for Non-Maximum Suppression in Section 2.1. We then describe the necessary modifications to the object detector and the training procedure in Section 2.2.

### 2.1 Proposed Non-Maximum Suppression

With classical Non-Maximum Suppression, all detections are first sorted by their confidence scores and added to a proposal list  $\mathcal{P}$ . The list of final detections  $\mathcal{D}$  is empty in the beginning. Then the following step is executed iteratively until  $\mathcal{P}$  is empty: The proposal  $p$  with the highest confidence score in  $\mathcal{P}$  is removed from  $\mathcal{P}$  and compared to all detections in  $\mathcal{D}$ . If the intersection over union between  $p$  and all detections in  $\mathcal{D}$  is smaller than a threshold  $N$  then  $p$  is added to  $\mathcal{D}$ . Otherwise  $p$  is discarded. The pseudo code of this algorithm is given in Algorithm 1.

---

#### Algorithm 1 Classical Non-Maximum Suppression.

---

```

 $\mathcal{P} \leftarrow \text{GETPROPOSALS}(\text{image})$ 
 $\mathcal{P} \leftarrow \text{SORT}(\mathcal{P})$ 
 $\mathcal{D} \leftarrow \emptyset$ 
while  $\mathcal{P} \neq \emptyset$  do
   $p \leftarrow \text{POP}(\mathcal{P})$ 
  for  $d \in \mathcal{D}$  do
     $\text{iou} \leftarrow \text{GETIOU}(p, d)$ 
    if  $\text{iou} \leq N$  then
       $\text{PUSH}(p, \mathcal{D})$ 
    end if
  end for
end while

```

---

This approach has one parameter  $N$  which has to be tuned to achieve good performance. A common choice is  $N = 0.5$ . The key idea of this algorithm is

that bounding boxes with a high overlap are likely to belong to the same object. Bounding boxes with a low overlap on the other hand are likely to belong to different objects.

There are however situations where this assumption fails. Especially in images with a high number of objects and partial occlusions there are many overlapping bounding boxes that belong to different objects. One example of this are crowds of humans.

We propose a novel approach to decide if two bounding boxes contain the same object or not. We call our approach FeatureNMS since it is based on (appearance) features of the detections. The overall structure of the proposed algorithm is the same as classical Non-Maximum Suppression—but the rule whether to add  $p$  from  $\mathcal{P}$  to  $\mathcal{D}$  or not is adjusted. The pseudo code of our approach is given in Algorithm 2.

---

**Algorithm 2** Proposed Non-Maximum Suppression. If the calculated value of the intersection over union is in a range that does not allow to make a definite decision we use a feature embedding similarity.

---

```

 $\mathcal{P} \leftarrow \text{GETPROPOSALS}(\text{image})$ 
 $\mathcal{P} \leftarrow \text{SORT}(\mathcal{P})$ 
 $\mathcal{D} \leftarrow \emptyset$ 
while  $\mathcal{P} \neq \emptyset$  do
   $p \leftarrow \text{POP}(\mathcal{P})$ 
  for  $d \in \mathcal{D}$  do
     $\text{iou} \leftarrow \text{GETIOU}(p, d)$ 
    if  $\text{iou} \leq N_1$  then
       $\text{PUSH}(p, \mathcal{D})$ 
    else if  $\text{iou} \geq N_2$  then
       $\text{embeddingDistance} \leftarrow \text{GETEMBEDDINGDISTANCE}(p, d)$ 
      if  $\text{embeddingDistance} > T$  then
         $\text{PUSH}(p, \mathcal{D})$ 
      end if
    end if
  end for
end while

```

---

Again, each proposal  $p \in \mathcal{P}$  is compared to all detections  $d \in \mathcal{D}$ . The intersection over union between  $p$  and  $d$  is computed. If this value is less or equal than a threshold  $N_1$  we assume that the detections belong to different objects. If this value on the other hand is at least as large as another threshold  $N_2$  we assume that the detections must belong to the same object. In any other case the two bounding boxes might belong to the same or to different objects—the intersection over union alone cannot be used to make a final decision. In this case we calculate the  $L^2$  distance of feature embeddings for both bounding boxes. If this distance is larger than a threshold  $T$  we assume that the bounding boxes belong to different objects. Otherwise they are likely to belong to the same object. The



feature embeddings are an output of the CNN that we use for object detection. It is described in detail in Section 2.2.

We propose to choose  $N_1 = 0.1$  and  $N_2 = 0.9$  but other values are possible, depending on the application. The right value for  $T$  depends on the training objective of the detection network. In our work we use  $T = \beta = 1.0$  (cf. Section 2.2).

## 2.2 Detector Architecture and Training

We evaluate our approach with the RetinaNet [5] object detector but it generalizes to a lot of different detector architectures. We add one head to the RetinaNet backbone. This head outputs a feature embedding for each anchor box. We chose an embedding of length of 32, but other lengths are possible. In our experiments, this was a good tradeoff between accuracy, computational overhead and memory consumption.

The network head for the feature embeddings consists of four identical blocks. Each block is formed by a 2D convolution layer with 512 channels, a Batch Normalization [25] layer and a ReLU activation function. The output of the last block is L2-normalized along the feature embedding dimension (consisting of 32 values). This makes sure that all embeddings lie on a unit hypersphere which is a common choice for embedding learning [23].

The training objective for the feature embedding is based on Margin Loss [24]. The total loss can be calculated as follows:

$$L = \frac{\sum_{i \in \mathcal{A}} \sum_{j \in \mathcal{A} \setminus \{i\}} L'(i, j)}{|\mathcal{A}| \cdot (|\mathcal{A}| - 1)} \quad (1)$$

In this equation  $L'$  is the pairwise loss between two targets:

$$L'(i, j) = \begin{cases} \max(0, \|\mathbf{f}_i, \mathbf{f}_j\|_2 - (\beta - \alpha)), & \text{if } obj(i) = obj(j) \\ \max(0, (\beta + \alpha) - \|\mathbf{f}_i, \mathbf{f}_j\|_2), & \text{otherwise} \end{cases} \quad (2)$$

Here,  $\mathcal{A}$  is the set of anchor boxes that are assigned to ground truth bounding boxes. The vector  $\mathbf{f}_i$  is the embedding feature vector that belongs to the target (anchor box)  $i$ . The function  $obj(i)$  gives the object id of target  $i$ . The parameter  $\alpha$  determines the margin between positive and negative examples, and the parameter  $\beta$  determines the decision threshold. We chose  $\alpha = 0.2$  and  $\beta = 1.0$ .

Our sampling strategy is different from [24]. Since we only train on active target pairs within a single image the number of pairs is limited. This means that we usually use all possible pairs during a training step. Only if the number of pairs exceeds 5 000 we use uniform sampling to restrict the number of samples to this.

We weight the different losses of the different training objectives according to [26]. This way the weights can adjust depending on the progress of training and do not have to be tuned manually.

### 3 Evaluation

We evaluate our approach on the CrowdHuman dataset [27]. This dataset contains 15 000 training images and 4 370 validation images. We use the validation images to compare the performance of the different NMS approaches, but we did not use it to tune any parameters. The dataset contains multiple annotations per person: A head bounding box, a visible region bounding box and a full body bounding box. We use the visible body bounding box annotation in this work. Before feeding the images into the network we resize them so that the longer side has a fixed amount of pixels. Then the image is padded with a fixed color value to obtain a square image.

Our implementation is based on the RetinaNet implementation from Tensorflow<sup>1</sup>. Our patches for this implementation that we used to perform the experiments are available online<sup>2</sup>. We use the default hyperparameters with the following exceptions:

- Batch size of 4
- 800 000 training steps
- LAMB optimizer [28]
- Learning rate
  - $1 \cdot 10^{-4}$  (step 0 - 100 000)
  - $5 \cdot 10^{-5}$  (step 100 000 - 200 000)
  - $1 \cdot 10^{-5}$  (step 200 000 - 400 000)
  - $5 \cdot 10^{-6}$  (step 400 000 - 800 000)
- Image size
  - $768 \times 768$  pixels (first 750 000 training steps)
  - $1024 \times 1024$  pixels (last 50 000 training steps and during testing)

We initialized the weights of our CNN backbone from a model that was pretrained on the COCO dataset [29]. During training we froze the weights of the first convolutional layer and the corresponding batch normalization layer.

We tried to train the proposed network with SGD (both with and without momentum) but the training did not converge. We also tried to use the Adam optimizer [30] but the trained network did not achieve good performance. The LAMB optimizer [28] on the other hand works well with our network and loss formulation. It combines the ideas of LARS [31] and Adam [30]. It stabilizes training by adjusting the learning rate per layer.

Most of the training steps were performed at a reduced resolution of  $768 \times 768$  pixels. The reason for that is that the GPU which we use for training does not have enough VRAM for higher resolutions with a batch size of 4. Afterwards we fine-tuned the network at full resolution on the CPU for 50 000 training steps.

We evaluate the different NMS approaches with three common metrics. The first is the average precision when requiring an IoU of at least 0.5 between detection and ground truth bounding box. The second is the average precision

<sup>1</sup> <https://github.com/tensorflow/models>

<sup>2</sup> <https://github.com/fzi-forschungszentrum-informatik/NNAD/tree/eccv2020>

at a minimum IoU of 0.75. The last metric we use is the log-average miss rate [32]. This metric is computed by averaging miss rates at 9 FPPI (false positives per image) values evenly spaced in log-space between  $10^{-2}$  and  $10^0$ . The IoU threshold used for this is 0.5.

The results can be found in Table 1. We also provide precision-recall curves for all approaches in Figure 1.

Our approach (FeatureNMS,  $N_1 = 0.1$ ,  $N_2 = 0.9$ ) outperforms all other approaches that we compared to. As an ablation study we evaluated our approach with different parameters and found that the performance does not change much. When using  $N_1 = 0.0$  and  $N_2 = 1.0$  the only assumption is that bounding boxes without any overlap can't belong to the same object. If there is any overlap the feature vector is always used to decide if they belong to the same object or not. When using  $N_1 = -\varepsilon$  and  $N_2 = 1.0$  even this assumption is given up. For each pair of detections in an image the feature vector is used to decide if a box should be suppressed. This experiment shows the discriminativeness of our feature vector. Even with these parameters precision and recall are high and our approach still performs better than the others.

The performance of classical NMS is below that of FeatureNMS except for very low detection score thresholds. Here the precision is low for both approaches but the recall of classical NMS is slightly higher. The reason for that is that there are a few cases where the feature vectors of detections that belong to different objects are too similar. These detections are erroneously suppressed by FeatureNMS but not by classical NMS.


SoftNMS [9] achieves similar precision as FeatureNMS at high detection score thresholds with low recall. But the precision at higher recall values is much lower.

We also compared our approach to AdaptiveNMS [10]. AdaptiveNMS predicts the local object density for each detection and uses that to adjust the threshold of classical NMS. We did not want to adjust the detector network for this because then the contributions to the achieved performance would not be clear: If the density estimation by the detector is not accurate because of our training approach, then it would influence the performance of AdaptiveNMS. Because of that we decided to use the ground truth density as input to AdaptiveNMS. This also means that the density estimation performance is an overestimate—a real-world detector will not achieve a perfect estimation.

To our surprise we found that AdaptiveNMS performs slightly worse than classical NMS with this ground truth density. The precision is slightly below that of classical NMS on nearly all points of the precision-recall curve. This is because the threshold for NMS is increased in densely populated regions of the image, which also leads to more false positives in these regions. Our findings are in contrast to the results reported in [10]. There are several possible explanations for this: One is that the localization performance of our detector is lower than that of the detectors used in the original paper. A lower localization performance will result in more false positives when the NMS threshold is high. Another possible explanation is that the ground truth density is actually not the best threshold: The neural network might not output a good density estimation but

a smoothed one that is closer to the average. This could suppress some false positives in areas with high object densities.

We also visually compared the detection results of our approach to these of classical NMS. Figure 2 contains some example images. We found that there are two situations where FeatureNMS outperforms classical NMS. The first situation occurs in the first two example images. Here, there are detections with high overlap that belong to different objects. Classical NMS suppresses some of these detections while FeatureNMS can correctly separate these. The second situation occurs in the second two example images. Here, the bounding box detector outputs some detections with low localization accuracy. Because of that the IoU between multiple detections of the same object is low. Classical NMS fails to suppress the duplications. FeatureNMS on the other hand is still able to correctly associate the detections using the feature vector.



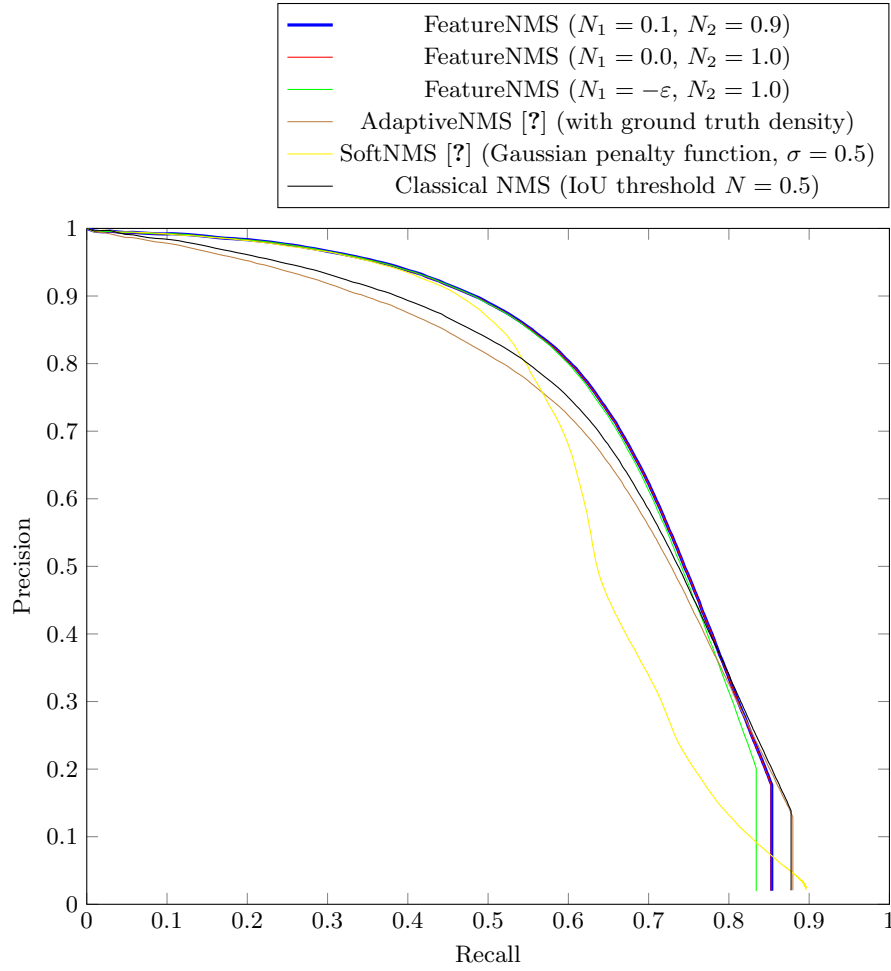
Method	AP @ 0.5IoU	AP @ 0.75IoU	log-average MR
FeatureNMS ( $N_1 = 0.1, N_2 = 0.9$ )	<b>0.6865</b>	<b>0.3030</b>	<b>0.7535</b>
FeatureNMS ( $N_1 = 0.0, N_2 = 1.0$ )	0.6860	0.3027	0.7545
FeatureNMS ( $N_1 = -\varepsilon, N_2 = 1.0$ )	0.6838	0.2996	0.7541
AdaptiveNMS [10] (with ground truth density)	0.6480	0.2843	0.8309
SoftNMS [9] (Gaussian, $\sigma = 0.5$ )	0.6280	0.2991	0.7582
Classical NMS (IoU threshold $N = 0.5$ )	0.6597	0.2855	0.8129

**Table 1.** Comparison of different approaches for NMS on the CrowdHuman dataset [27]. We evaluated the average precision (AP) at a minimum IoU of 0.5 and 0.75, as well as the log-average miss rate (MR) [32]. Our approach (FeatureNMS) outperforms all other approaches used for comparison.

## 4 Conclusion

FeatureNMS is a simple yet effective approach to Non-Maximum Suppression. It outperforms all approaches that we used for comparison on the CrowdHuman dataset [27]. At the same time the inference run-time overhead is low: Additionally to the operations of classical NMS it only requires to compute a feature vector per bounding box detection and to compare these for overlapping bounding boxes. The necessary changes in the object detector network are minor and the approach can be used with most CNN detector architectures.





**Fig. 1.** Precision-Recall curves of different approaches for NMS on the CrowdHuman dataset [27].



**Fig. 2.** Comparison of example images when applying FeatureNMS and classical NMS.

## References

1. Redmon, J., Divvala, S., Girshick, R., Farhadi, A.: You Only Look Once: Unified, Real-Time Object Detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2016) 779–788
2. Redmon, J., Farhadi, A.: YOLO9000: Better, Faster, Stronger. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2017) 7263–7271
3. Redmon, J., Farhadi, A.: YOLOv3: An Incremental Improvement. arXiv preprint arXiv:1804.02767 (2018)
4. Liu, W., Anguelov, D., Erhan, D., Szegedy, C., Reed, S., Fu, C.Y., Berg, A.C.: SSD: Single Shot MultiBox Detector. In: Proceedings of the European Conference on Computer Vision, Springer (2016) 21–37
5. Lin, T.Y., Goyal, P., Girshick, R.B., He, K., Dollár, P.: Focal Loss for Dense Object Detection. In: ICCV, IEEE Computer Society (2017) 2980–2988
6. Girshick, R., Donahue, J., Darrell, T., Malik, J.: Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2014) 580–587
7. Girshick, R.: Fast R-CNN. In: Proceedings of the IEEE International Conference on Computer Vision. (2015) 1440–1448
8. Ren, S., He, K., Girshick, R.B., Sun, J.: Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks. In: Proceedings of the IEEE International Conference on Computer Vision. (2015) 1440–1448
9. Bodla, N., Singh, B., Chellappa, R., Davis, L.S.: Soft-NMS — Improving Object Detection With One Line of Code. In: Proceedings of the IEEE International Conference on Computer Vision. (2017) 5561–5569
10. Liu, S., Huang, D., Wang, Y. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2019) 6459–6468
11. Gählert, N., Hanselmann, N., Franke, U., Denzler, J.: Visibility Guided NMS: Efficient Boosting of Amodal Object Detection in Crowded Traffic Scenes. (2019)
12. Wang, X., Xiao, T., Jiang, Y., Shao, S., Sun, J., Shen, C.: Repulsion Loss: Detecting Pedestrians in a Crowd. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2018) 7774–7783
13. Zhang, S., Wen, L., Bian, X., Lei, Z., Li, S.Z.: Occlusion-aware R-CNN: Detecting Pedestrians in a Crowd. In: Proceedings of the European Conference on Computer Vision. (2018) 637–653
14. Hosang, J., Benenson, R., Schiele, B.: Learning Non-Maximum Suppression. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2017) 4507–4515
15. Hu, H., Gu, J., Zhang, Z., Dai, J., Wei, Y.: Relation Networks for Object Detection. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2018) 3588–3597
16. Bucher, M., Herbin, S., Jurie, F.: Improving Semantic Embedding Consistency by Metric Learning for Zero-Shot Classification. In: Proceedings of the European Conference on Computer Vision, Springer (2016) 730–746
17. Wang, J., Song, Y., Leung, T., Rosenberg, C., Wang, J., Philbin, J., Chen, B., Wu, Y.: Learning Fine-grained Image Similarity with Deep Ranking. In: Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition. (2014) 1386–1393

18. Bell, S., Bala, K.: Learning visual similarity for product design with convolutional neural networks. *ACM Transactions on Graphics* **34**(4) (2015) 1–10
19. Shankar, D., Narumanchi, S., Ananya, H., Kompalli, P., Chaudhury, K.: Deep Learning based Large Scale Visual Recommendation and Search for E-Commerce. *arXiv preprint arXiv:1703.02344* (2017)
20. Hoffer, E., Ailon, N.: Deep metric learning using Triplet network. In: *International Workshop on Similarity-Based Pattern Recognition*, Springer (2015) 84–92
21. Dey, S., Dutta, A., Toledo, J.I., Ghosh, S.K., Lladós, J., Pal, U.: SigNet: Convolutional Siamese Network for Writer Independent Offline Signature Verification. *arXiv preprint arXiv:1707.02131* (2017)
22. Hadsell, R., Chopra, S., LeCun, Y.: Dimensionality Reduction by Learning an Invariant Mapping. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. Volume 2., IEEE (2006) 1735–1742
23. Schroff, F., Kalenichenko, D., Philbin, J.: FaceNet: A Unified Embedding for Face Recognition and Clustering. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2015) 815–823
24. Wu, C.Y., Manmatha, R., Smola, A.J., Krahenbuhl, P.: Sampling Matters in Deep Embedding Learning. In: *Proceedings of the IEEE International Conference on Computer Vision*. (2017) 2840–2848
25. Ioffe, S., Szegedy, C.: Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift. *arXiv preprint arXiv:1502.03167* (2015)
26. Kendall, A., Gal, Y., Cipolla, R.: Multi-Task Learning Using Uncertainty to Weigh Losses for Scene Geometry and Semantics. In: *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. (2018) 7482–7491
27. Shao, S., Zhao, Z., Li, B., Xiao, T., Yu, G., Zhang, X., Sun, J.: CrowdHuman: A Benchmark for Detecting Human in a Crowd. *arXiv preprint arXiv:1805.00123* (2018)
28. You, Y., Li, J., Reddi, S., Hseu, J., Kumar, S., Bhojanapalli, S., Song, X., Demmel, J., Hsieh, C.J.: Large Batch Optimization for Deep Learning: Training BERT in 76 minutes. *arXiv preprint arXiv:1904.00962* (2019)
29. Lin, T.Y., Maire, M., Belongie, S., Bourdev, L., Girshick, R., Hays, J., Perona, P., Ramanan, D., Zitnick, C.L., Dollár, P.: Microsoft COCO: Common Objects in Context. In: *Proceedings of the European Conference on Computer Vision*, Springer (2014) 740–755
30. Kingma, D.P., Ba, J.: Adam: A Method for Stochastic Optimization. *arXiv preprint arXiv:1412.6980* (2014)
31. You, Y., Gitman, I., Ginsburg, B.: Large Batch Training of Convolutional Networks. *arXiv preprint arXiv:1708.03888* (2017)
32. Dollar, P., Wojek, C., Schiele, B., Perona, P.: Pedestrian Detection: An Evaluation of the State of the Art. *IEEE Transactions on Pattern Analysis and Machine Intelligence* **34**(4) (2011) 743–761