# FeatureNMS: Non-Maximum Suppression by Learning Feature Embeddings

Niels Ole Salscheider

FZI Research Center for Information Technology

Haid-und-Neu-Str. 10-14

76131 Karlsruhe

Germany

Email: 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: This heuristic assumes that a high overlap between two bounding boxes corresponds to a high probability of one being a duplicate. We propose FeatureNMS to solve this problem. FeatureNMS recognizes duplicates not only based on the intersection over union between the 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.**

## I. 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 a set of proposals. A dedicated second stage then decides which proposals are 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 post-processing 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 the latter detection is removed.

This heuristic however has shortcomings in crowded scenes because the underlying assumption does not hold. In these scenes, distinct objects often have a high overlap. In this paper we propose *FeatureNMS* to solve this problem. Our approach recognizes duplicates based on their feature embeddings if a definite decision based on the IoU is not possible.

The remainder of this paper is structured as follows: Section II presents related work. Section III describes our proposed approach to Non-Maximum Suppression. The general idea is presented in Section III-A while Section III-B contains details about the necessary modifications to the object detector network. In Section IV we present our evaluation procedure and the results. Finally, Section V concludes the paper.

## II. Related Work

Both NMS and embedding learning have been studied in previous research. This section presents relevant and related work in these fields.

### A. Non-Maximum Suppression

There have been several proposals how to improve the classical NMS 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 local object density. This local object density is predicted by the object detection network for each 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. Given detections of two different objects, the IoU for the visible parts is usually smaller than the IoU for the whole objects. Because of that, classical NMS is performed on the bounding boxes of the visible parts. But the final output are the corresponding bounding boxes of the whole object.

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

In [14], the authors propose to solve the NMS task with a CNN. The proposed network learns to re-score detections to suppress duplicates. 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. It can use geometric and appearance features of the detections for this.

## B. 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 widely used for this purpose. It consists of two terms: One term pulls the $\ell^2$ distance of positive pairs as close to zero as possible. The other term pushes the $\ell^2$ distance of negative pairs apart if it is below a certain margin.

Choosing the margin parameter correctly can be challenging. It can be too difficult to push the embeddings of hard negative examples 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 to 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 not try to push the embeddings of all positive pairs to be as close to each other as possible. Instead, it just requires the distance to be below 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.

## III. APPROACH

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

### A. 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.

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

---

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

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

---

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 belong to 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.

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 larger than 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 $\ell^2$ distance of the feature embeddings of 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 III-B.

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 III-B).

### B. Detector Architecture and Training

We evaluate our approach with the RetinaNet [5] object detector. But it generalizes to many other detector architectures—the only required change is to learn an embedding vector

**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.

---

$\mathcal{P} \leftarrow$ GETPROPOSALS($image$)
$\mathcal{P} \leftarrow$ SORT($\mathcal{P}$)
$\mathcal{D} \leftarrow \emptyset$
**while** $\mathcal{P} \neq \emptyset$ **do**
    $p \leftarrow$ POP($\mathcal{P}$)
    $isDuplicate \leftarrow$ false
    **for** $d \in \mathcal{D}$ **do**
        $iou \leftarrow$ GETIOU($p, d$)
        **if** $iou > N_2$ **then**
            $isDuplicate \leftarrow$ true
        **else if** $iou > N_1$ **then**
            $embeddingDist \leftarrow$ GETEMBEDDINGDIST($p, d$)
            **if** $embeddingDist < T$ **then**
                $isDuplicate \leftarrow$ true
            **end if**
        **end if**
    **end for**
    **if** $\neg isDuplicate$ **then**
        PUSH($p, \mathcal{D}$)
    **end if**
**end while**

---

per detection. For this, we add one network head to each output pyramid level of the RetinaNet backbone. The head outputs an embedding vector for each anchor box. We chose an embedding of length of 32, but other lengths are possible. In our experiments, this proved to be a good trade-off between accuracy, computational overhead and memory consumption.

Like all other RetinaNet heads, 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 $\ell^2$-normalized along the embedding dimension (consisting of 32 values). This ensures 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:

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

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 can nearly always 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 5 000.

We weight the different losses during training according to [26]. This way, the weighting factors can adjust based on the training progress and do not have to be tuned manually.

## IV. 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. In this work, we use the visible body bounding box annotations. 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[1]. Our patches for this implementation that we used to perform the experiments are available online[2].

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.

Most of the training steps were performed at a reduced resolution of $768 \times 768$ pixels. The reason is that the limited VRAM of our GPU does not allow to train at 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

---

[1]https://github.com/tensorflow/models
[2]https://github.com/fzi-forschungszentrum-informatik/NNAD/tree/featurenms

an IoU of at least 0.5 between detection and ground truth bounding box. The second is the average precision at a minimum IoU of 0.75. The last metric that we use is the log-average miss rate [30]. 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 I. We also provide precision-recall curves for all approaches in Figure 1. All reported values are based on the output of the same detector network—only the NMS approach differs.

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 make the final decision. 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. This is because in a few cases 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 a bad network design or training approach could distort the achieved accuracy: If the density estimation by the detector is not accurate it could reduce 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

| Method | AP @ 0.5IoU | AP @ 0.75IoU | log-average MR |
|---|---|---|---|
| FeatureNMS ($N_1 = 0.1$, $N_2 = 0.9$) | **0.6865** | **0.3030** | **0.7535** |
| 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 I
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) [30]. OUR APPROACH (FEATURENMS) OUTPERFORMS ALL OTHER APPROACHES USED FOR COMPARISON.

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 estimation 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 for the same object is low. Classical NMS fails to suppress the duplicates. FeatureNMS on the other hand is still able to correctly associate the detections based on the feature vector.

## V. 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 run-time overhead during inference is low: It performs the same operations as classical NMS. Additionally to these, it only requires to compute a feature vector per bounding box detection and to compare them 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.
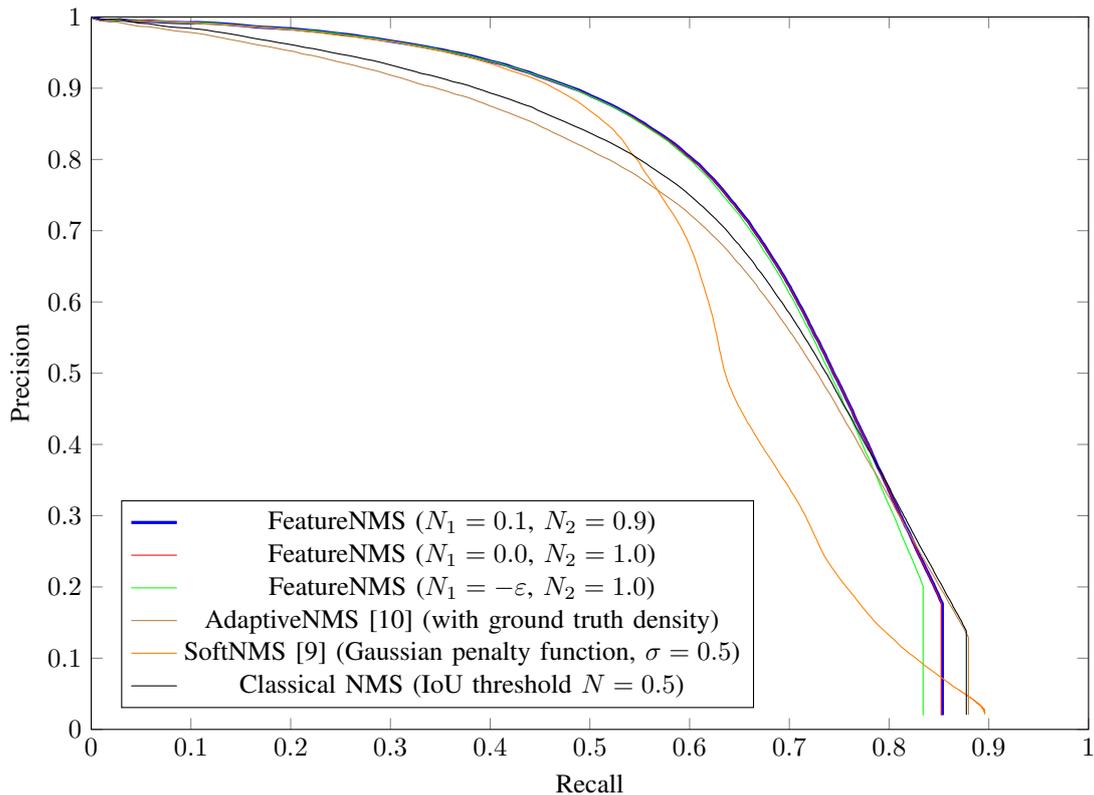
Fig. 1. Precision-Recall curves of different approaches for NMS on the CrowdHuman dataset [27].

REFERENCES

[1] J. Redmon, S. Divvala, R. Girshick, and A. Farhadi, "You Only Look Once: Unified, Real-Time Object Detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2016, pp. 779–788.

[2] J. Redmon and A. Farhadi, "YOLO9000: Better, Faster, Stronger," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 7263–7271.

[3] ——, "YOLOv3: An Incremental Improvement," *arXiv preprint arXiv:1804.02767*, 2018.

[4] W. Liu, D. Anguelov, D. Erhan, C. Szegedy, S. Reed, C.-Y. Fu, and A. C. Berg, "SSD: Single Shot MultiBox Detector," in *Proceedings of the European Conference on Computer Vision*. Springer, 2016, pp. 21–37.

[5] T.-Y. Lin, P. Goyal, R. B. Girshick, K. He, and P. Dollár, "Focal Loss for Dense Object Detection," in *ICCV*. IEEE Computer Society, 2017, pp. 2980–2988.

[6] R. Girshick, J. Donahue, T. Darrell, and J. Malik, "Rich Feature Hierarchies for Accurate Object Detection and Semantic Segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 580–587.

[7] R. Girshick, "Fast R-CNN," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1440–1448.

[8] S. Ren, K. He, R. B. Girshick, and J. Sun, "Faster R-CNN: Towards Real-Time Object Detection with Region Proposal Networks," in *Proceedings of the IEEE International Conference on Computer Vision*, 2015, pp. 1440–1448.

[9] N. Bodla, B. Singh, R. Chellappa, and L. S. Davis, "Soft-NMS — Improving Object Detection With One Line of Code," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 5561–5569.

[10] S. Liu, D. Huang, and Y. Wang, "Adaptive nms: Refining pedestrian detection in a crowd," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2019, pp. 6459–6468.

[11] N. Gählert, N. Hanselmann, U. Franke, and J. Denzler, "Visibility Guided NMS: Efficient Boosting of Amodal Object Detection in Crowded Traffic Scenes," in *Proceedings of Conference on Neural Information Processing Systems*, 2019.

[12] X. Wang, T. Xiao, Y. Jiang, S. Shao, J. Sun, and C. Shen, "Repulsion Loss: Detecting Pedestrians in a Crowd," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 7774–7783.

[13] S. Zhang, L. Wen, X. Bian, Z. Lei, and S. Z. Li, "Occlusion-aware R-CNN: Detecting Pedestrians in a Crowd," in *Proceedings of the European Conference on Computer Vision*, 2018, pp. 637–653.

[14] J. Hosang, R. Benenson, and B. Schiele, "Learning Non-Maximum Suppression," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2017, pp. 4507–4515.

[15] H. Hu, J. Gu, Z. Zhang, J. Dai, and Y. Wei, "Relation Networks for Object Detection," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2018, pp. 3588–3597.

[16] M. Bucher, S. Herbin, and F. Jurie, "Improving Semantic Embedding Consistency by Metric Learning for Zero-Shot Classification," in *Proceedings of the European Conference on Computer Vision*. Springer, 2016, pp. 730–746.

[17] J. Wang, Y. Song, T. Leung, C. Rosenberg, J. Wang, J. Philbin, B. Chen, and Y. Wu, "Learning Fine-grained Image Similarity with Deep Ranking," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2014, pp. 1386–1393.

[18] S. Bell and K. Bala, "Learning visual similarity for product design with convolutional neural networks," *ACM Transactions on Graphics*, vol. 34, no. 4, pp. 1–10, 2015.

[19] D. Shankar, S. Narumanchi, H. Ananya, P. Kompalli, and K. Chaudhury, "Deep Learning based Large Scale Visual Recommendation and Search for E-Commerce," *arXiv preprint arXiv:1703.02344*, 2017.

[20] E. Hoffer and N. Ailon, "Deep metric learning using Triplet network," in *International Workshop on Similarity-Based Pattern Recognition*. Springer, 2015, pp. 84–92.

[21] S. Dey, A. Dutta, J. I. Toledo, S. K. Ghosh, J. Lladós, and U. Pal, "SigNet: Convolutional Siamese Network for Writer Independent Offline Signature Verification," *arXiv preprint arXiv:1707.02131*, 2017.

FeatureNMS

Classical NMS

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

[22] R. Hadsell, S. Chopra, and Y. LeCun, "Dimensionality Reduction by Learning an Invariant Mapping," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, vol. 2. IEEE, 2006, pp. 1735–1742.

[23] F. Schroff, D. Kalenichenko, and J. Philbin, "FaceNet: A Unified Embedding for Face Recognition and Clustering," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2015, pp. 815–823.

[24] C.-Y. Wu, R. Manmatha, A. J. Smola, and P. Krahenbuhl, "Sampling Matters in Deep Embedding Learning," in *Proceedings of the IEEE International Conference on Computer Vision*, 2017, pp. 2840–2848.

[25] S. Ioffe and C. Szegedy, "Batch Normalization: Accelerating Deep Network Training by Reducing Internal Covariate Shift," *arXiv preprint arXiv:1502.03167*, 2015.

[26] A. Kendall, Y. Gal, and R. Cipolla, "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, pp. 7482–7491.

[27] S. Shao, Z. Zhao, B. Li, T. Xiao, G. Yu, X. Zhang, and J. Sun, "CrowdHuman: A Benchmark for Detecting Human in a Crowd," *arXiv preprint arXiv:1805.00123*, 2018.

[28] Y. You, J. Li, S. Reddi, J. Hseu, S. Kumar, S. Bhojanapalli, X. Song, J. Demmel, and C.-J. Hsieh, "Large Batch Optimization for Deep Learning: Training BERT in 76 minutes," *arXiv preprint arXiv:1904.00962*, 2019.

[29] T.-Y. Lin, M. Maire, S. Belongie, L. Bourdev, R. Girshick, J. Hays, P. Perona, D. Ramanan, C. L. Zitnick, and P. Dollár, "Microsoft COCO: Common Objects in Context," in *Proceedings of the European Conference on Computer Vision*. Springer, 2014, pp. 740–755.

[30] P. Dollar, C. Wojek, B. Schiele, and P. Perona, "Pedestrian Detection: An Evaluation of the State of the Art," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 34, no. 4, pp. 743–761, 2011.