

Transfer Learning Techniques Using Simulation Data For Machine Learning Automotive Radar Systems

Felix Rutz^{*}, Ralph Rasshofer[†] und Erwin Biebl[‡]

Abstract: For a reliable detection and classification of vulnerable road users in modern automotive radar systems, the latest research introduces machine-learning (ML) based algorithms. However, suitable training datasets for ML systems based on real-world radar measurements are rarely available or lack specific raw radar data. Different approaches based on transfer-learning methods from data generated by a simulation framework for the range-Doppler-representation of radar measurement data are researched. In particular, influences of dataset size and sample quality, as well as different transfer learning approaches concerning the performance of the ML system in the radar data domain, are examined.

Index Terms: Automotive Radar, Machine Learning, Radar Dataset Simulation, Transfer Learning

1 Introduction

The EU Vision Zero road traffic safety initiative seeks to decrease road injuries and fatalities by 2050 [1]. Therefore, advanced driver-assistance technologies are required in modern automobiles to achieve this goal. The foundation of these systems to perform successfully is a thorough sensing of the vehicle's surroundings. Hence, radar and lidar scanners, front cameras, and ultrasonic probes comprise a basic sensor ensemble enabling a multimodal data-driven depiction of the surrounding automobile environment. The sensor-specific information obtained is subsequently analyzed and optimized for different safety-related applications, such as autonomous emergency braking, blind spot detection, and higher automated driving capabilities [2].

With the implementation of advanced conditional and automated driving functions, the various raw sensor data are often combined and evaluated by elaborately trained machine learning systems. Their performance levels are unprecedented compared to standard signal-processing methods. For example, a deep learning system is used as an evaluation framework joining the vision and the motion planning domain for improved pedestrian detection [3]. However, the application of artificial intelligence-based functions is not limited to safety-related systems. In [4], a deep learning system regulates a complex automatic torque converter transmission, outperforming the control performance of classical control approaches.

^{*}Professur für Höchstfrequenztechnik, Technische Universität München (e-mail: felix.rutz@tum.de).

[†]BMW AG, München (e-mail: Ralph.Rasshofer@bmw.de).

[‡]Professur für Höchstfrequenztechnik, Technische Universität München (e-mail: biebl@tum.de).

2 Related Work

For improved protection of vulnerable road users, radar systems, camera devices, and lidar sensors have been investigated for reliable identification and categorization of target objects. Comparing these available measurement devices, the advantage of radar sensors is their least-imperishable characteristic in changing weather or lightning conditions, even on long detection ranges, as shown by [5]. However, the advantage is not limited to axial perception, as component supplier Bosch lately introduced synthetic aperture functionality to the vehicle for a detailed lateral perception [6].

The latest generation of commercially available vehicular long-range radars uses continuous frequency-modulated ramps for target detection in the 76 – 77 GHz band [7]. This frequency band allows the use of broadband waveforms, which have advantageous effects on range and velocity resolutions. Further, the angular resolution is enhanced by antenna arrays due to the small mechanical dimensions corresponding with the required wavelength.

In order to meet the requirements of self-driving vehicles in terms of radar sensor characteristics, scaling the radar parameters, for example, by higher bandwidths, is not sufficient. In order to meet the demands of environment sensing, fundamentally different approaches are required. For this purpose, research is being conducted in the field of signal processing on more complex algorithms. Approaches include the evaluation of large array structures using sophisticated models for improved angular resolution [8], as well as the use of alternative radar modulation techniques [9].

The rising capabilities of signal processing methods based on machine learning (ML) is another promising approach, further enabling the combined detection and classification of remote objects rather than only detecting the presence of an unclassified canonical target. In [10] the range-Doppler-representation of radar measurements is used for detection tasks, whereas [11] presents a detection algorithm using the range-azimuth spectrum instead.

Common to machine-learning approaches is the limited availability of sufficiently large training datasets. Hence, the presented deep neural network algorithms are often tailored to work well on a specific input dataset limited to the situational events represented within the learning set. Therefore, a direct comparison of different approaches remains impossible, as datasets and networks are incompatible with each other [12].

For a reliable inference and generalization to unseen driving scenarios, a sufficient training database with ideally all possible situations as data samples is required for all ML-based approaches. As such datasets are generally unavailable, learning from computer-generated radar data is under research. We, therefore, create simulation datasets and use them to train a computer vision object detection system based on the YOLOv5 framework [13] on the radar data domain. Different training experiments focus on the effects of using various simulated dataset sizes. The results are evaluated for sample quantity and quality of the dataset. The corresponding ML models are then separately re-trained with input data based on radar sensor measurements using different transfer learning approaches to improve the system predictions.

Rather than solely adapting the ML model to a new domain by freezing specific layers as feature representations and re-training them with different datasets as depicted in [14], we also introduce modifications to the feature space so the model adapts well to a slightly different domain within the same training run.

3 Evaluation Metrics

For the evaluation of the machine learning system, the precision, recall, and mean average precision metrics are used following the remarks by [15]. The precision is defined as the ratio of true positives (TP) and the total number of predicted positives as a sum of the true positives and false positives (FP):

$$Precision = \frac{TP}{TP + FP}. \quad (1)$$

The precision metric expresses the accuracy of the neural network as the proportion of its correct predictions to all positive predictions. Recall or sensitivity is defined as the ratio of true positives and the total of ground truth positives as a sum of the true positives and false negatives (FN):

$$Recall = \frac{TP}{TP + FN}. \quad (2)$$

The recall metric is related to the ability of the machine learning model to find all the positive samples from the dataset. The classification system is characterized by merging both parameters into the precision over recall curve. The area under the curve is summarized as a single number used as a numeric performance indicator referred to as the average precision (AP) for each class. Deriving the mean of the AP over all distinct curves for each class leads to the mean average precision (mAP) as a general evaluation metric of an object detection system.

4 Dataset Generation from Simulation

The ML framework model performance is significantly determined by data quality and quantity used for supervised training. For the ML model examinations, we use two separate datasets. One dataset was derived from measurements with the radarbook experimental platform. The device provides a millimeter-wave radar measurement setup with raw data processing capabilities in the 76 – 77 GHz band, incorporating similar properties in frequency, modulation schemes, and bandwidth as standard automotive-grade radar sensors. The semi-manual labeling process is described in detail by [16].

The primary population of the processed measurement samples is split into a training and a test set for evaluating the ML model on unseen data. The training set is derived from 2445 pedestrian labels, 1527 bicycle annotations, 785 automobile boxes, and 562 frames with empty boxes. The test set contains 204 pedestrian samples, 553 bicycle annotations, and 95 automobile labels. The dataset is assumed to be relatively small. Regarding class distribution, the allocations of the target classes within the distinct subsets are pretty unbalanced.

Nevertheless, processing real-world sensor data and deriving large and balanced datasets is an elaborative and cost-intensive task. Identifying and labeling rare but dangerous driving scenes for their representation in the training dataset is a nearly unachievable undertaking, often also accompanied by legal implications.

Instead of manually labeling more data from an actual radar sensor, a simulated dataset is created based on the MATLAB simulation framework introduced by [17]. The cyclist, vehicle, and pedestrian objects are modeled as a sum of their characteristic reflection points. Each point target is assigned a unique position, velocity, and associated characteristic backscattering cross-section. Based on the different reflection models for pedestrians, bicyclists, and vehicles, the simulation framework calculates a radar time-domain baseband signal using the radar range formula. Multi-path propagation is excluded in order to reduce the complexity of the simulation.

The simulated sensor data is then translated into range-Doppler-map representations with appropriate ground-truth annotation labels. The synthetically generated dataset has three times as many samples as the sensor-based dataset derived from measurements, with 18 300 sample images.

A total of 12 200 samples from simulated data are used for training. The training example set contains multi-labeled samples with 8311 pedestrian labels, 10 309 bicycle annotations, and 9230 automobile boxes. The test set consists of the remaining 6100 samples. The sample distribution per class is more balanced within the simulated radar dataset than the measurement-based annotations. Table 1 compares the total number of range-Doppler-map samples and the total number of class labels for the three class entities, car, bicycle, and person, between the original sensor dataset and the simulated one.

Table 1: Comparison of manually labeled and simulated datasets

	Radar Sensor Dataset	Simulation Dataset
Total sample count	4889	18 300
No. of person labels	2649	13 219
No. of bicycle labels	2080	16 226
No. of car labels	880	12 946

5 Analysis of Simulation Data

The ability of a machine learning model to generalize and fit a dataset strongly correlates with the dataset size and data quality itself. Based on the Bagging method [18], N multiple training sets from the 12 200 simulation base population X have been drawn. However, the samples were drawn without replacement for data quantity and quality analysis concerning the ML model prediction, but with increasing the total subset size N in a range from 500 to 12 000. The step size between two dataset iterations is 500 additional samples. This prevents possibly faulty samples from being drawn into a training set multiple times.

For every subset, a separate ML model was trained from scratch. Fig. 1 depicts the results of the corresponding mAP of the resulting ML models in dependence on the different training datasets.

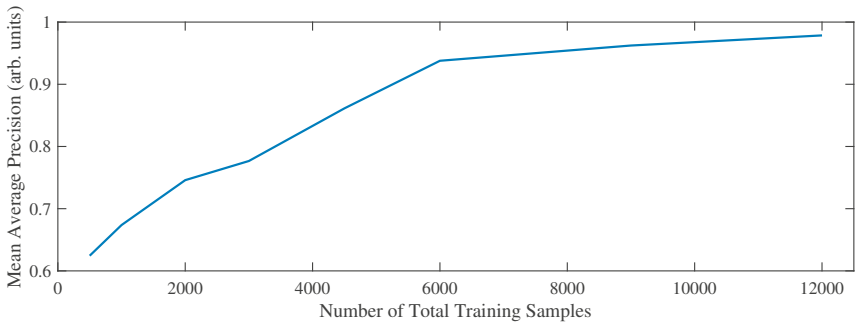
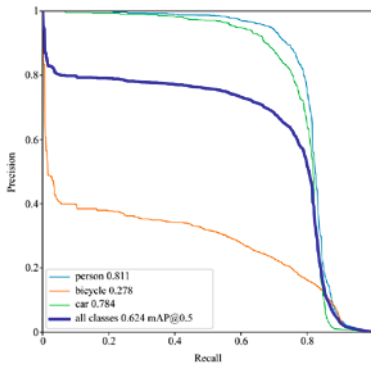


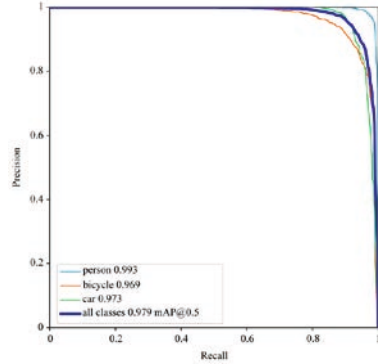
Figure 1: Evaluation of mean Average Precision with the total dataset sample size based on simulated data

On the smallest available dataset containing only 500 training samples, the ML system reaches a mean average precision of 0.62. The metric increases steeply with every increase of the training dataset until it reaches a pivot point in mAP of 0.95 using a training dataset of 6000 samples. Further experiments only correlate with minor improvements in the overall mAP score when increasing the training database size, reaching a mAP limit of 0.97 using 12 000 training examples.

The corresponding parameter models for the before-mentioned supporting nodes using 500, 6000, and 12 000 training samples are selected for further experiments. The mean average precision shows a high deviation at both system boundaries. Therefore, the detailed precision-recall plots for each class for the smallest and largest datasets are examined. Deriving from the definition of precision and recall, the combined plot summarises the trade-off between the true positive rate and the positive predictive value for an ML model [19]. Fig. 2 depicts both plots for training datasets containing 500 and 12 000 total samples.



(a) Precision-Recall-Curve for 500 samples



(b) Precision-Recall-Curve for 12000 samples

Figure 2: Comparison of Precision-Recall-Curves for training sets including 500 and 12000 total samples based on simulation data

In a qualitative comparison of both graphics, the smaller dataset shows significant setbacks in the total area covered by the precision over recall graphs for all three classes. The most noticeable flaw is discovered for the bicycle class, resulting in a mAP value of 0.278 using 500 training samples. Contrarily, the area under the curve reaches its maximum when training with 12000 samples, leading to a nearly perfect precision over recall for all classes.

Therefore, the dataset for the small-sized training set has been further investigated. Every sample image from the small dataset was inspected by hand. In the total number of bicycle instances, corresponding bounding boxes for 14 samples have been identified as misplaced. Hence, the samples only contain background noise instead of the true target ground truth representation. With an increasing number of training samples, the influence of misplaced labels seems to be mitigated within the ML model. The reason for the misplaced bounding boxes remains unclear. A likely explanation is an error during the transformation of ground truth labels into a corresponding data format, as different ML algorithms require different representations of ground truth labeling parameters.

6 Transfer to Measurement Data Domain

To implement an ML-based system in a vehicle, the models generated from simulation data must also generalize in sensor measurement data. The underlying system parameter weights from simulation data are therefore used to improve the generalization capabilities in a slightly changed domain setting. The system needs to adapt to the environmental radar data domain rather than purely classify data based on calculations excluding real-world effects such as multi-path propagation. Assuming that many factors and model weights leading to the results in simulation data also apply to the radar sensor data domain, the adaptation process of the ML model is referred to as transfer learning [20].

In traditional machine learning, a specifically tailored algorithm is trained and implemented for one detailed task using a curated dataset representing the problem. However, this approach requires the unseen input data to share the exact distribution and feature space with the available training data. This requirement is not fulfilled in most real-world applications, as shown in the previous section comparing the sample distribution in the range-Doppler-map datasets. Using transfer learning (TL) instead, the ML system can re-use its knowledge and skills of a specific task to solve a problem in another target domain. Therefore, the need for high-quality data in the target domain is reduced.

6.1 Transfer Learning by Updating Feature Representation

Following the explanations in [21], scalar feature weights represent knowledge in a machine-learning model in a network of processing units. Each processing element implements a nonlinear function, altering its input data with a specific weight. Multiple units form a specific network structure, often consisting of subsequent layers. For transfer learning, as many weights as possible are re-used from training in the source domain and fine-tuned with a small amount of data in the target domain. During the target training process, the element weights are updated, enabling the model to acquire more information to represent the input data. Different weights need to be evolved depending on the chosen ML system architecture.

Refined in detail by [22], two key aspects affect the outcome of feature adaptation, one factor being the size of the target dataset. If the target dataset is small, overfitting in the network is avoided by freezing as many layers as possible. Hence, the ML model relies more on features extracted from the source set samples. The second factor is the similarity between the source and the target dataset. With limited variations in the dataset samples, faster fine-tuning results are achievable. Finding the optimal number of layers that need to be updated or fixed during training is an incidental challenge in transfer learning. From the presented model structure, different configurations of layers are being evolved during transfer training runs with the target sensor dataset researched.

6.2 Transfer Learning by Feature Space Modification

Another approach for domain adaptation is presented in [23]. Instead of changing the network weights using multiple training runs on various source and target datasets, the sample set is expanded. By simply merging samples from both database populations into a joint training dataset, the adjacent learning process by fitting an ML model from scratch is compelled to derive a more general network representation, as samples from both data domains are considered during the same training session. Due to the continuous weight updates within each training epoch, the ML system focuses on identifying underlying features common to all input samples. Hence, a suitable generalization of the final model is achieved, and simultaneous overall training time is reduced. This straightforward approach is implemented and compared to the procedure based on updating the feature representation using two distinct datasets.

7 Transfer Learning Results

With the ML model based on 12 000 simulation samples, different transfer learning techniques have been applied using the radar sensor data. The best-considered results from the different approaches are compared with traditional ML training using only the sensor dataset as a default baseline model. Table 2 summarizes the results of the distinct transfer learning methods described before. For bicycle and car classes best results are achieved by training the ML model from scratch using a mixed dataset containing simulated and real-world data. Nonetheless, transfer learning using distinct datasets yields the best average precision for person class detection.

Table 2: Results of different transfer learning techniques

ML Training Approach	AP (Person)	AP (Bicycle)	AP (Car)	mAP@0.5
Trad. ML on Sensor Data	0.690	0.670	0.600	0.653
TL (Feature Rep.)	0.730	0.696	0.541	0.656
TL (Feature Space)	0.687	0.713	0.719	0.706

8 Conclusion

Different transfer learning techniques have been researched to improve ML model generalization in automotive radar domain data. The additional parameter adaptation from simulated radar domain data can mitigate the limitations of small sensor datasets for training. The simulation framework generates scalable datasets, including the corresponding ground-truth labels, for pre-training quickly and reliably.

However, the decrease in mAP when transferring the model weights from the simulation to the measurement data domain indicates that more underlying effective mechanisms exist in the sensor radar domain. Nevertheless, transfer learning methods significantly increase the algorithm’s performance compared to the default baseline model derived from the traditional system training approach.

Depending on the distinct object class, precision is increased up to 11.9 percent using transfer learning approaches for model generalization. The requirement for refining datasets from raw sensor data containing an even distribution of common and rare events is reduced by adding enormous samples from simulation frameworks.

Future research focuses on further improvements in generalization by data augmentation of sensor-based datasets by upsampling the total amount of available training pairs, which is directly related to further reducing data sample collection efforts.

References

- [1] European Climate, Infrastructure and Environment Executive Agency (CINEA), *EU Road Safety: Towards "Vision Zero"*, European Union, 2022.
- [2] H. Winner, S. Hakuli, F. Lotz, and C. Singer, Eds., *Handbuch Fahrerassistenzsysteme: Grundlagen, Komponenten und Systeme für aktive Sicherheit und Komfort*, 3rd ed., ser. ATZ-MTZ-Fachbuch. Wiesbaden: Springer Vieweg, 2015.
- [3] M. Lyssenko, C. Gladisch, C. Heinzemann, M. Woehle, and R. Triebel, "Towards safety-aware pedestrian detection in autonomous systems," in *2022 IEEE/RSJ International Conference on Intelligent Robots and Systems (IROS)*, 2022, pp. 293–300.
- [4] G. Gaiselmann, S. Altenburg, S. Studer, and S. Peters, "Deep Reinforcement Learning for Gearshift Controllers in Automatic Transmissions," *Array*, vol. 15, p. 100235, 2022.
- [5] A. Mukhtar, L. Xia, and T. B. Tang, "Vehicle Detection Techniques for Collision Avoidance Systems: A Review," *IEEE Transactions on Intelligent Transportation Systems*, vol. 16, no. 5, pp. 2318–2338, 2015.
- [6] M. Farhadi, R. Feger, J. Fink, T. Wagner, and A. Stelzer, "Automotive synthetic aperture radar imaging using tdm-mimo," in *2021 IEEE Radar Conference (Radar-Conf21)*, 2021, pp. 1–6.
- [7] K. Ramasubramanian and K. Ramaiah, "Moving from Legacy 24 GHz to State-of-the-Art 77-GHz Radar," *ATZelektronik worldwide*, vol. 13, no. 3, pp. 46–49, 2018. [Online]. Available: <https://doi.org/10.1007/s38314-018-0029-6>
- [8] G. Schnoering, C. Höller, T. Kawaguchi, K. Kawajiri, and S. Malterer, "Fast angular processing for sparse fmcw radar arrays with non-uniform fft," in *2023 24th International Radar Symposium (IRS)*, 2023, pp. 1–11.
- [9] G. Hakobyan and B. Yang, "High-Performance Automotive Radar: A Review of Signal Processing Algorithms and Modulation Schemes," *IEEE Signal Processing Magazine*, vol. 36, no. 5, pp. 32–44, 2019.
- [10] W. Ng, G. Wang, Siddhartha, Z. Lin, and B. J. Dutta, "Range-doppler detection in automotive radar with deep learning," in *2020 International Joint Conference on Neural Networks (IJCNN)*, 2020, pp. 1–8.
- [11] K. Patel, K. Rambach, T. Visentin, D. Rusev, M. Pfeiffer, and B. Yang, "Deep learning-based object classification on automotive radar spectra," in *2019 IEEE Radar Conference (RadarConf)*, 2019, pp. 1–6.
- [12] D. Feng, C. Haase-Schütz, L. Rosenbaum, H. Hertlein, C. Gläser, F. Timm, W. Wiesbeck, and K. Dietmayer, "Deep Multi-Modal Object Detection and Semantic Segmentation for Autonomous Driving: Datasets, Methods, and Challenges," *IEEE Transactions on Intelligent Transportation Systems*, vol. 22, no. 3, pp. 1341–1360, 2021.

- [13] G. Jocher *et al.*, “ultralytics/yolov5: v6.2 - YOLOv5 Classification Models, Apple M1, Reproducibility, ClearML and Deci.ai integrations,” Aug. 2022. [Online]. Available: <https://doi.org/10.5281/zenodo.7002879>
- [14] F. Rutz, E. Biebl, and J.-P. Konrad, “Transfer Learning in ML-based Radar Systems for Automotive Applications,” in *2022 Kleinheubach Conference*, 2022, pp. 1–3.
- [15] K. P. Murphy, *Machine learning: A probabilistic perspective*, 4th ed., ser. Adaptive computation and machine learning series. Cambridge, Mass.: MIT Press, 2013.
- [16] R. Pérez, F. Schubert, R. Rasshofer, and E. Biebl, “Deep learning radar object detection and classification for urban automotive scenarios,” in *2019 Kleinheubach Conference*, 2019, pp. 1–4.
- [17] T. Wengerter, R. Pérez, E. Biebl, J. Worms, and D. O’Hagan, “Simulation of urban automotive radar measurements for deep learning target detection,” in *2022 IEEE Intelligent Vehicles Symposium (IV)*, 2022, pp. 309–314.
- [18] E. Alpaydin, *Introduction to Machine Learning*, 4th ed., ser. Adaptive Computation and Machine Learning. Cambridge, Massachusetts: The MIT Press, 2020.
- [19] T. Saito and M. Rehmsmeier, “The precision-recall plot is more informative than the roc plot when evaluating binary classifiers on imbalanced datasets,” *PloS one*, vol. 10, no. 3, p. e0118432, 2015.
- [20] I. Goodfellow, Y. Bengio, and A. Courville, *Deep learning*, ser. Adaptive Computation and Machine Learning. Cambridge, Massachusetts: The MIT Press, 2016.
- [21] R. Reed and R. J. Marks, *Neural Smithing: Supervised Learning in Feedforward Artificial Neural Networks*. The MIT Press, 02 1999. [Online]. Available: <https://doi.org/10.7551/mitpress/4937.001.0001>
- [22] M. Elgendy, *Deep learning for vision systems*. Shelter Island: Manning, 2020.
- [23] H. Daumé, “Frustratingly Easy Domain Adaptation,” 2009. [Online]. Available: <https://arxiv.org/abs/0907.1815>