

# A Study on the Faithfulness of Feature Attribution Explanations in Pruned Vision-Based Multi-Task Learning

Xenia Demetriou<sup>1†</sup>, Sophie Sananikone<sup>1†</sup>, Vojislav Tobias Westmoreland<sup>1†</sup>,  
Matthia Sabatelli<sup>1</sup> and Marco Zullich<sup>1,\*</sup>

<sup>1</sup>Department of AI, Faculty of Science and Engineering, University of Groningen, Nijenborg 9, 9747 AG, Groningen, the Netherlands

## Abstract

As Artificial Neural Networks (ANNs) continue growing in size to enhance predictive accuracy, model compression techniques, such as pruning, counteract the increase in energy and computational costs. Multi-Task Learning (MTL) consists in training a single model on multiple tasks, providing regularization and increased generalization, especially in applications like robotics and autonomous driving. While compressed models often match the accuracy of their larger counterparts, they are usually evaluated on metrics related to task completion or efficiency, overlooking critical aspects such as fairness, and transparency—key requirements for high-stakes applications. Specifically, Explainable AI offers tools for making model predictions more transparent, for instance, by means of feature importance. However, especially when AI models are complex, these tools can lead to unfaithful or unreliable explanations, potentially undermining trust in these models. In the present work, we propose to investigate whether pruning applied to vision-based MTL significantly affects the faithfulness of the explanations generated for the tasks the models are trained on. We train and prune different models on the benchmark datasets NYUv2 and CityScapes. Despite the hurdle of generating feature importance for tasks such as surface normal prediction and depth estimation, our results show that unstructured pruning maintains faithfulness across different sparsity percentages. Structured pruning with milder sparsity percentages preserves faithfulness, but can decrease more rapidly at higher sparsity percentages. Overall, sparsity up to a certain threshold (90%) does not compromise explanation faithfulness. Beyond this point, both faithfulness and performance decline significantly, making them unfit to be deployed, rendering the faithfulness risk negligible.

## Keywords

Faithfulness, Explainability, Multi-Task Learning, Model Compression

## 1. Introduction

The deployment of artificial vision systems through Artificial Intelligence (AI) has transformed various domains such as autonomous driving and medical imaging, resulting in remarkable performance increases over the past decade. In Computer Vision (CV) tasks, Convolutional Neural Networks (CNNs) have emerged as one of the cornerstones of this technology. However, they are characterized by high computational costs and a lack of interpretability due to their complexity. Their computational footprint can however be tackled by Model Compression (MC) techniques. One of the main techniques for MC is pruning [1], which aims at removing (groups of) synapses in an Artificial Neural Network (ANN) according to specific criteria. Another viable technique consists in training models in a Multi-Task Learning (MTL) setting [2], where one ANN can generate predictions for multiple tasks given the same input, instead of having task-specific ANNs, each responsible for generating one prediction per task. MTL is useful in areas such as robotics and autonomous driving, where a machine learning model may be tasked to perform multiple tasks at the same time (e.g., Depth Estimation—DE—and semantic

*Late-breaking work, Demos and Doctoral Consortium, colocated with The 3rd World Conference on eXplainable Artificial Intelligence: July 09–11, 2025, Istanbul, Turkey*

\*Corresponding author.

<sup>†</sup>These authors contributed equally.

✉ x.demetriou3@gmail.com (X. Demetriou); sophie.sananikone12@gmail.com (S. Sananikone);  
vojo.westmoreland@gmail.com (V. T. Westmoreland); m.sabatelli@rug.nl (M. Sabatelli); marco.zullich@gmail.com  
(M. Zullich)

🌐 www.zullich.it (M. Zullich)

🆔 0009-0007-3677-3637 (X. Demetriou); 0009-0001-1509-6813 (S. Sananikone); 0009-0008-9252-9871 (V. T. Westmoreland);  
0009-0007-7540-8616 (M. Sabatelli); 0000-0002-9920-9095 (M. Zullich)



© 2025 Copyright for this paper by its authors. Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).

segmentation) in a resource-constrained environment. It can additionally be considered a MC tool, since MTL models are typically less computationally demanding than using separate models for each individual task [3, 4].

Explainable AI (XAI) is crucial for understanding complex model predictions, particularly in applications requiring transparency [5]. Feature attribution, a key XAI tool, assigns “saliency scores” to input features, highlighting their importance. The main model-agnostic methods are LIME [6], which relies on surrogate models, and SHAP [7], which has solid game theoretical foundations. Additionally, ANN-specific tools, such as Layerwise Relevance Propagation (LRP) [8], which tracks down components of an ANN responsible for a given activation, have been developed. In ANNs, feature attribution can be performed by backpropagating gradients of the output on the input features, like in the case of Guided Backpropagation [9]. However, these gradient-based methods have often been found to be unreliable, often acting similarly to edge detectors [10]. Grad-CAM [11], a feature attribution tool specific for CNNs, combines gradients and activations, being more class-discriminative and faithful to the underlying model than other gradient-based tools. If compared to LIME, SHAP, and LRP, gradient-based methods are often much more efficient [12], which is a reason why Grad-CAM is often the *de facto* choice for feature attributions in CNNs [13]. The present work utilizes and expands Grad-CAM to explain DE and Surface Normal (SN) prediction tasks. One of the paramount issues is that all these methods are *approximations* of the complex decision-making process of ANNs, and can hence produce unfaithful outputs, which highlight irrelevant features [14].

Analyzing the effect of pruning on the quality of explanations has been researched previously: Abbasi-Asl and Yu [15] studied the features learned by CNNs on image classification, noticing that redundancy in CNN filters coalesced when applying filter pruning. Weber et al. [16] provided a human-grounded evaluation of saliency maps generated on CNNs for image classification. Their findings suggest that explanations produced by moderately pruned models are evaluated as *superior* by humans with respect to their dense counterparts. Finally, Khakzar et al. [17] introduced a custom pruning method, termed PruneGrad, shown to quantitatively produce more faithful explanations; however, this is an input-specific method, thus reducing the real world practicality and scalability.

In the present work, we aim at answering the question: “**How is the faithfulness of multi-task model explanations affected by pruning?**”. We train three different CNNs for MTL on the NYUv2 and CityScapes datasets. We then proceed to apply different pruning methods, to obtain several pruned models at different sparsity levels. Subsequently, we generate task-level saliency maps using variations of the popular image-based XAI tool Grad-CAM [11]. Finally, we proceed to evaluate the faithfulness of these maps with respect to the models using the Iterative Removal of Features (IROF) technique [18]. Our findings highlight how sparsity up to a certain threshold (dependent upon the specific pruning technique) does not compromise the faithfulness of explanation; The main consequence is that, as far as explainability is concerned, well-performing, moderately sparse models will likely yield similarly faithful explanations.

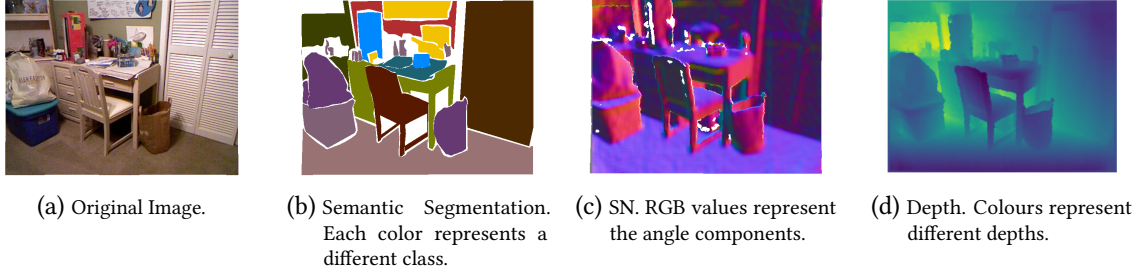
Our contributions are twofold: (i) We carry out an analysis on the effect of pruning on the faithfulness of models trained for vision-based MTL, and (ii) We provide an adaptation of Grad-CAM and its assessment for monocular DE and SN, a task which Grad-CAM is not originally designed for.

## 2. Methods

In this section, we provide the materials and methods used in the present work, alongside a summary of the experimental setup.

### 2.1. Datasets

We employed two datasets to train and evaluate our models; The NYU-Depth V2 (NYUv2) dataset [19], and the Cityscapes dataset [20], both commonly used benchmarks for MTL and scene understanding in indoor and urban settings. NYUv2 contains annotations for Semantic Segmentation, SN estimation, and DE, while Cityscapes only provides labels for Semantic Segmentation and DE. We provide an example



**Figure 1:** An example of an image in the NYUv2 dataset and its ground truth labels for each task.

of one image and the corresponding annotations for NYUv2 in Figure 1. Semantic Segmentation is a pixel-level classification task. We assess the quality according to the mean Intersection-over-Union (mIoU) metric. The SN task is a pixel-level three-dimensional regression task, where surface angles are described by their  $x$ ,  $y$ , and  $z$  coordinates in a 3D space. We evaluate it using the median angle error over all pixels in an image. DE is a pixel-level regression task where the distance of each pixel relative to the camera is measured. We assess it using the relative error between the ground truth and the predicted depth.

## 2.2. Model Architectures

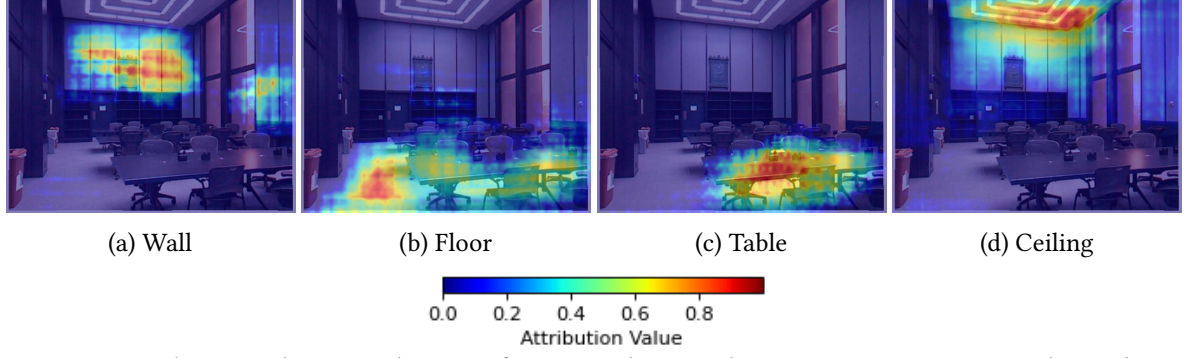
We made use of two CCN-based model architectures specifically crafted for MTL: DeepLab [21] and SegNetMTAN [22, 23]. Both techniques can be used in combination with the FAMO weighting method [24] to correctly balance the contribution of the various tasks.

## 2.3. Pruning

Pruning compresses ANNs by removing synapses based on saliency. *Unstructured* pruning removes individual connections, while *structured* pruning removes groups (neurons, filters, etc.). Structured pruning offers hardware-agnostic speed-ups via tensor dimension reduction, whereas unstructured pruning requires specialized hardware. Pruning can be *dynamic* (allowing synapse regrowth) or *static* (no regrowth). This work employs DiSparse [3] (unstructured, gradient-based MTL), Network Slimming (NetSlim) [25] (structured, channel-level L1-norm), and HRank [26] (structured, filter-rank-based). DiSparse avoids task-conflicting pruning, NetSlim uses L1-norm for channel importance, while HRank prunes low-rank filters.

## 2.4. Explanations

For an RGB image  $x$ , we generate a saliency map  $s$  indicating pixel importance. Grad-CAM, using CNN layer activations and output gradients, produces a matrix  $s'$ , normalized and rescaled to the dimensions of  $x$ , so each pixel is assigned a saliency from 0 to 1. Following Selvaraju et al. [11], we apply it to the last convolutional layer. For semantic segmentation, we use SegGrad-CAM [27], explaining per-class pixel groups. An example of the explanations generated using SegGrad-CAM is present in Figure 2. In addition, we operate the following methodology to adapt Grad-CAM to DE and SN. We group predictions according to their values into specific categories: for DE, we divide the  $[0 - 1]$  range of depth prediction into quintiles and compute the gradients according to each of these groups; in this way, we aim at explaining the feature importance in determining a given quintile of the depth predictions. Similarly, for SNs, we group the predictions in the eight *octants* of the 3d plane and generate explanations accordingly. Despite the problems of DE and SN estimation being *pixel-level regression* tasks, the proposed methodology restructures the output into a classification problem, thus allowing us to use Grad-CAM.



**Figure 2:** Example SegGrad-CAM explanations for various classes in the Semantic Segmentation task. Attribution, representing feature importance, is shown on the scale at the bottom.

## 2.5. Assessing XAI

As highlighted by Nauta et al. [28], faithfulness is paramount in XAI evaluation. While other explanation properties exist, unfaithful explanations highlight irrelevant features, potentially causing loss of trust in AI systems by stakeholders. The assessment of faithfulness boils down to determining whether the features marked as *salient* are effectively important for the model. One way of evaluating this is by progressively *masking out* features according to their importance and measuring the change in the model prediction. The expectation is that highly important features should cause a large change in the prediction. This is the rationale behind IROF: after the process of iterative masking, the change in prediction is plotted against the proportion of features removed, and faithfulness is computed as the area-over-the-curve (AOC). A detailed visualization of the IROF procedure is shown in Figure 3.

Formally, every input  $X$  is segmented in superpixels using SLIC [29]. Each segment is given a saliency score according to the mean importance provided by Grad-CAM. Next, we compute the output of the model on the perturbed input  $F(X^l)_y$ , where  $l$  is the number of segments removed, and divide it by the output on the unperturbed input  $F(X^0)_y$ . By varying this procedure over all segments  $l \in \{1, \dots, L\}$  in order of saliency, the aforementioned curve can be plotted and AOC computed as a faithfulness score. This metric is then averaged over a whole set of  $N$  datapoints to produce the final estimate:

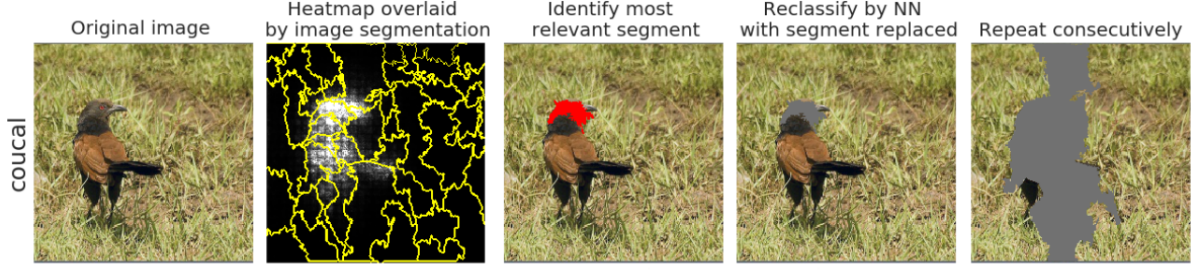
$$\text{IROF}(X, F) = \frac{1}{N} \sum_{n=1}^N \text{AOC} \left( \frac{F(X_n^l)_y}{F(X_n^0)_y} \right)_{l=0}^L.$$

In the present study, we used the IROF implementation from the Quantus library by Hedström et al. [30]. Since IROF requires scalar outputs to be compared, for our three tasks, we adapted the procedure as follows: (a) for Semantic Segmentation, we considered the mean of the predicted logits for each specific class of interest; (b) for DE, we considered the relative error between unperturbed and perturbed outputs, while (c) for SN, we considered the cosine similarity between the unperturbed and perturbed outputs.

## 2.6. Experimental Setup

The present study investigates the effect of two independent variables (sparsity and pruning method) on two dependent variables (model performance and IROF Scores). We tested this on three model architectures to improve the robustness of the conclusions. Different models addressed different combinations of tasks. We applied a different pruning method for each model architecture. Table 1 shows an overview of the model architectures used, with their respective pruning method used, sparsity percentages, tasks, and datasets. The specific optimizers, training hyperparameters, and data augmentation strategy are instead presented in Appendix A. We point out that the aim of our experiments, as far as



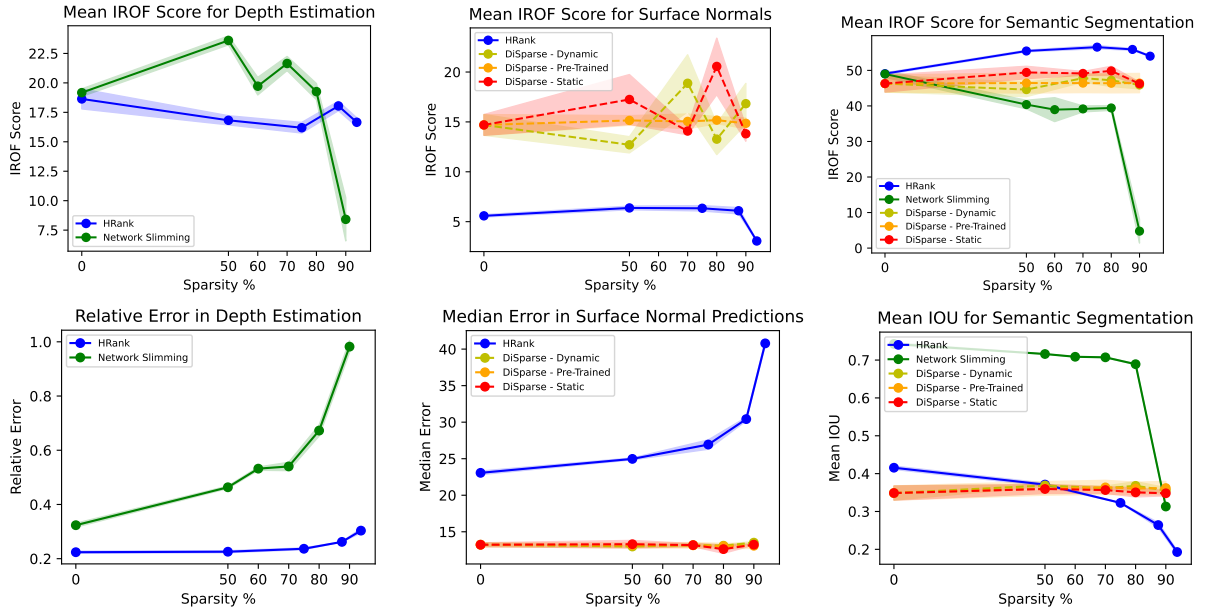


**Figure 3:** A visual example of the IROF procedure for the computation of faithfulness.

**Table 1**

Overview of models, pruning methods, sparsity levels, datasets, and tasks used in this study. Each model undergoes a specific pruning technique—either structured or unstructured—while exploring different sparsity levels. They are evaluated on tasks such as semantic segmentation, SN, and DE. Additional information on the training hyperparameters can be found in Appendix A.

Model	Pruning Method	Sparsity (%)	Tasks	Dataset
DeepLab	Unstructured (DiSparse)	0, 50, 70, 80, 90	Semantic Segmentation, SN	NYUv2
SegNetMTAN w/ FAMO	Structured (NetSlim)	0, 50, 60, 70, 80, 90	Semantic Segmentation, DE	Cityscapes
SegNetMTAN w/o FAMO	Structured (HRank)	0, 50, 75, 87.5, 93.75	Semantic Segmentation, SN, DE	NYUv2



**Figure 4:** Results for faithfulness (first row) and model performance (second row) as a function of sparsity for the different tasks (one task per column).

model performance is concerned, is to reproduce the results achieved in the original implementations [21, 22, 23], not to beat the state-of-the-art on the specific tasks. The goal of the present paper is indeed the faithfulness of the explanations.

### 3. Results

Results, as shown in Figure 4 following the setup in Table 1, reveal diverse trends across tasks and pruning methods. For **DE**, Relative Error rises exponentially with sparsity for both model/dataset

combinations. IROF scores remain stable until 75% sparsity; HRank maintains better IROF at higher sparsity compared to NetSlim. For **SN**, DiSparse (DeepLab) achieves significantly higher performance and IROF than HRank (SegNetMTAN), even with substantial pruning due to DiSparse’s unstructured nature allowing for aggressive pruning. IROF scores show a similar difference between DiSparse and HRank, with sparsity minimally affecting the score until 90%. In **Semantic Segmentation**, NetSlim (SegNetMTAN) achieves stable, high performance (around 70% mIoU), with a drop at 90% sparsity. HRank (SegNetMTAN) and DiSparse (DeepLab) perform worse (35-45% mIoU). IROF scores remain largely stable across pruning rates for all methods, with declines at high sparsity ( $\sim 90\%$ ) for structured pruning.

## 4. Discussion and Conclusions

In the present work, we investigated the ties between pruning and *faithfulness* of the explanations in the case of MTL. We trained three different model architectures on the datasets CityScapes and NYUv2. We then pruned these models using iterative structured and unstructured pruning techniques, reaching sparsity levels above 90%. We then proceeded to generate per-task input attribution explanations using variations of GradCAM. Finally, we assessed the faithfulness of these explanations, faithfulness being the degree to which the explanations capture the complex predictive process of the underlying black-box model. We evaluated the faithfulness using the popular Iterative Removal of Features (IROF) metric, which proceeds to iteratively delete progressively more important (as measured by the explanation) information from the input, expecting a variation in the model’s output in the process. The results indicate how the DiSparse unstructured pruning technique seems not to particularly affect the faithfulness of the explanations; however, the structured pruning techniques generally show a decrement in faithfulness as the pruning rate increases—a trend which seems to be loosely connected to the decrease in performance—which comes with excessive pruning—measured across the various tasks.

Our findings highlight several positive impacts. Firstly, it encourages the implementation of DNNs with unstructured pruning in embedded systems and resource-constrained environments without compromising explanation faithfulness. In turn, this also helps foster user trust and regulatory compliance in such systems. We suggest positive findings in terms of AI Democratization through the decreased computational demands that come as a result of compression, making these models more accessible to a broader audience. For structured pruning, the results show that sparsity up to a certain threshold maintains explanation faithfulness. Beyond this point, both faithfulness and performance decline significantly, making such highly pruned models unlikely to be deployed. Thus, the overall risk from reduced faithfulness is minimal.

However, our work comes with several limitations: first of all, it is strictly tied to vision-based MTL. In addition, we assess faithfulness only according to IROF, while other methodologies exist for the evaluation. Moreover, our analysis lacks considerations on other facets of evaluation of explanations, such as *robustness* and *coherence*, which still should heavily be subject to considerations on faithfulness. Our analysis is also strictly tied to functional assessment, while several XAI studies operate human-grounded evaluation as a step for *validating* the explanations. These limitations also guide our next steps in extending the project: (i) extend the study to non-vision tasks, (ii) include other metrics for faithfulness and other facets of explanations, (iii) and carry out a human-grounded evaluation step. Considering the paramount importance of having faithful feature importance explanations, we still believe our findings to be of use for practitioners willing to implement MC on MTL.

## 5. Acknowledgements

We thank the Center for Information Technology of the University of Groningen for their support and for providing access to the Hábrók high performance computing cluster.

# Declaration on Generative AI

The author has not employed any Generative AI tools.

## References

- [1] T. Hoefler, D. Alistarh, T. Ben-Nun, N. Dryden, A. Peste, Sparsity in deep learning: Pruning and growth for efficient inference and training in neural networks, *J. Mach. Learn. Res.* 22 (2021).
- [2] R. Caruana, Multitask learning, *Machine learning* 28 (1997) 41–75.
- [3] X. Sun, A. Hassani, Z. Wang, G. Huang, H. Shi, Disparse: Disentangled sparsification for multitask model compression, 2022 IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR) (2022) 12372–12382.
- [4] Y. Zhang, Q. Yang, A survey on multi-task learning, *IEEE transactions on knowledge and data engineering* 34 (2021) 5586–5609.
- [5] M. E. Kaminski, The right to explanation, explained, in: *Research handbook on information law and governance*, Edward Elgar Publishing, 2021, pp. 278–299.
- [6] M. T. Ribeiro, S. Singh, C. Guestrin, “Why Should I Trust You?”: Explaining the Predictions of Any Classifier, in: *Proceedings of the 22nd ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ACM, 2016, pp. 1135–1144.
- [7] S. M. Lundberg, S.-I. Lee, A Unified Approach to Interpreting Model Predictions, in: *Advances in Neural Information Processing Systems*, volume 30, 2017, pp. 4765–4774.
- [8] S. Bach, A. Binder, G. Montavon, F. Klauschen, K.-R. Müller, W. Samek, On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation, *PloS one* 10 (2015) e0130140.
- [9] J. T. Springenberg, A. Dosovitskiy, T. Brox, M. Riedmiller, Striving for simplicity: The all convolutional net, *arXiv preprint arXiv:1412.6806* (2014).
- [10] W. Nie, Y. Zhang, A. Patel, A theoretical explanation for perplexing behaviors of backpropagation-based visualizations, in: J. Dy, A. Krause (Eds.), *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, PMLR, 2018, pp. 3809–3818.
- [11] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, D. Batra, Grad-CAM: Visual Explanations from Deep Networks via Gradient-based Localization, *International Journal of Computer Vision* 128 (2020) 336–359.
- [12] M. Miró-Nicolau, A. J. i Capó, G. Moyà-Alcover, A comprehensive study on fidelity metrics for xai, *Information Processing & Management* 62 (2025) 103900.
- [13] L. Arrighi, I. A. de Moraes, M. Zulich, M. Simonato, D. F. Barbin, S. B. Junior, Explainable Artificial Intelligence Techniques for Interpretation of Food Datasets: a Review, *arXiv preprint arXiv:2504.10527* (2025).
- [14] L. Arrighi, S. Barbon Junior, F. A. Pellegrino, M. Simonato, M. Zullich, Explainable Automated Anomaly Recognition in Failure Analysis: is Deep Learning Doing it Correctly?, in: *World Conference on Explainable Artificial Intelligence*, Springer, 2023, pp. 420–432.
- [15] R. Abbasi-Asl, B. Yu, Structural compression of convolutional neural networks with applications in interpretability, *Frontiers in big Data* 4 (2021) 704182.
- [16] D. Weber, F. Merkle, P. Schöttle, S. Schlögl, Less is more: The influence of pruning on the explainability of CNNs, *arXiv preprint arXiv:2302.08878* (2023).
- [17] A. Khakzar, S. Baselizadeh, S. Khanduja, C. Rupprecht, S. T. Kim, N. Navab, Improving feature attribution through input-specific network pruning, *arXiv preprint arXiv:1911.11081* (2019).
- [18] L. Rieger, L. Hansen, Irof: a low resource evaluation metric for explanation methods, in: *Proceedings of the Workshop AI for Affordable Healthcare at ICLR 2020*, 2020.
- [19] S. Gupta, P. Arbelaez, J. Malik, Perceptual Organization and Recognition of Indoor Scenes from RGB-D Images, in: *2013 IEEE Conference on Computer Vision and Pattern Recognition*, IEEE, Portland, OR, USA, 2013, p. 564–571.

- [20] M. Cordts, M. Omran, S. Ramos, T. Rehfeld, M. Enzweiler, R. Benenson, U. Franke, S. Roth, B. Schiele, The Cityscapes dataset for semantic urban scene understanding, in: Proceedings of the IEEE conference on computer vision and pattern recognition, 2016, pp. 3213–3223.
- [21] L.-C. Chen, G. Papandreou, I. Kokkinos, K. Murphy, A. L. Yuille, Deeplab: Semantic image segmentation with deep convolutional nets, atrous convolution, and fully connected crfs, IEEE transactions on pattern analysis and machine intelligence 40 (2017) 834–848.
- [22] V. Badrinarayanan, A. Kendall, R. Cipolla, Segnet: A deep convolutional encoder-decoder architecture for image segmentation, IEEE transactions on pattern analysis and machine intelligence 39 (2017) 2481–2495.
- [23] S. Liu, E. Johns, A. J. Davison, End-to-end multi-task learning with attention, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2019, pp. 1871–1880.
- [24] B. Liu, Y. Feng, P. Stone, Q. Liu, FAMO: Fast adaptive multitask optimization, Advances in Neural Information Processing Systems 36 (2024).
- [25] Z. Liu, J. Li, Z. Shen, G. Huang, S. Yan, C. Zhang, Learning efficient convolutional networks through network slimming, in: Proceedings of the IEEE international conference on computer vision, 2017, pp. 2736–2744.
- [26] M. Lin, R. Ji, Y. Wang, Y. Zhang, B. Zhang, Y. Tian, L. Shao, Hrank: Filter pruning using high-rank feature map, in: Proceedings of the IEEE/CVF conference on computer vision and pattern recognition, 2020, pp. 1529–1538.
- [27] K. Vinogradova, A. Dibrov, G. Myers, Towards Interpretable Semantic Segmentation via Gradient-weighted Class Activation Mapping, Proceedings of the AAAI Conference on Artificial Intelligence 34 (2020) 13943–13944.
- [28] M. Nauta, J. Trienes, S. Pathak, E. Nguyen, M. Peters, Y. Schmitt, J. Schlötterer, M. van Keulen, C. Seifert, From anecdotal evidence to quantitative evaluation methods: A systematic review on evaluating explainable ai, ACM Computing Surveys 55 (2023) 1–42.
- [29] R. Achanta, A. Shaji, K. Smith, A. Lucchi, P. Fua, S. Süsstrunk, SLIC Superpixels Compared to State-of-the-Art Superpixel Methods, IEEE Transactions on Pattern Analysis and Machine Intelligence 34 (2012) 2274–2282.
- [30] A. Hedström, L. Weber, D. Krakowczyk, D. Bareeva, F. Motzkus, W. Samek, S. Lapuschkin, M. M.-C. Höhne, Quantus: An explainable ai toolkit for responsible evaluation of neural network explanations and beyond, Journal of Machine Learning Research 24 (2023) 1–11.

## A. Training hyperparameters

We report here the optimizer and hyperparameters we used for training our models. The settings are as similar as possible to the original implementations found in the source references [21, 22, 23].

Model + Pruning	Optimizer	Learning Rate & Annealing	Batch size	Iterations	Augmentation
DeepLab+DiSparse	Adam	0.001, w/ Anneal $\times 0.5$ every 4k iterations	32	20k	None
SegNetMTAN+NetSlim		0.0001	8	36.4k	Random Scaling [1.0, 1.2, 1.5 $\times$ ] and Random horizontal flip
SegNetMTAN+HRank		0.0001, w/ Anneal $\times 0.5$ at half training	2	145k	