

# Privacy Meets Explainability: A Comprehensive Impact Benchmark

SAIFULLAH SAIFULLAH<sup>1,2\*</sup>, DOMINIQUE MERCIER<sup>1,2\*</sup>, ADRIANO LUCIERI<sup>1,2\*</sup> ANDREAS DENGEL<sup>1,2</sup>, AND SHERAZ AHMED<sup>2</sup>

<sup>1</sup>Department of Computer Science, Technische Universität Kaiserslautern, Erwin-Schrödinger-Straße 52, 67663 Kaiserslautern, Germany

<sup>2</sup>Smart Data and Knowledge Services (SDS), DFKI GmbH, 67663 Kaiserslautern, Germany (e-mail: firstname.lastname@dfki.de)

Corresponding author: Saifullah Saifullah (E-mail: saifullah.saifullah@dfki.de).

\*Equal contribution

This work was supported by the BMBF projects SensAI (BMBF Grant 01IW20007) and ExplAINN (BMBF Grant 01IS19074). We thank all members of the Deep Learning Competence Center at the DFKI for their comments and support.

**ABSTRACT** Since the mid-10s, the era of Deep Learning (DL) has continued to this day, bringing forth new superlatives and innovations each year. Nevertheless, the speed with which these innovations translate into real applications lags behind this fast pace. Safety-critical applications, in particular, underlie strict regulatory and ethical requirements which need to be taken care of and are still active areas of debate. eXplainable AI (XAI) and privacy-preserving machine learning (PPML) are both crucial research fields, aiming at mitigating some of the drawbacks of prevailing data-hungry black-box models in DL. Despite brisk research activity in the respective fields, no attention has yet been paid to their interaction. This work is the first to investigate the impact of private learning techniques on generated explanations for DL-based models. In an extensive experimental analysis covering various image and time series datasets from multiple domains, as well as varying privacy techniques, XAI methods, and model architectures, the effects of private training on generated explanations are studied. The findings suggest non-negligible changes in explanations through the introduction of privacy. Apart from reporting individual effects of PPML on XAI, the paper gives clear recommendations for the choice of techniques in real applications. By unveiling the interdependencies of these pivotal technologies, this work is a first step towards overcoming the remaining hurdles for practically applicable AI in safety-critical domains.

**INDEX TERMS** Artificial Intelligence, Attribution, Deep Learning, Differential Privacy, Explainability, Federated Learning, Neural Networks, Privacy-preserving.

## I. INTRODUCTION

IN recent years, a wide variety of deep learning (DL) approaches have achieved outstanding performance in a wide range of application domains [1], [2]. The versatile applications of deep neural networks in areas such as image processing [3], object segmentation [4], document analysis [5], time series classification [6], time series prediction [7], layout classification [8], sensor analysis and other areas have contributed to immense growth. Intelligent and automated decision-making bears the potential to improve and transform a row of critical application domains, including finance, healthcare, transportation, and administration. However, DL-based systems rely on complex, data-driven black-box methods whose exact working mechanisms are still widely unexplained in the scientific community, while the secure application of algorithms in safety-critical domains requires transparency and traceability of decisions. Furthermore, data security is a serious concern in domains involving

critical and personal data. DL methods are data-driven, often requiring the transmission, processing, and storage of large amounts of data in multiple remote locations. A variety of works showed, that even after training, neural networks can leak sensitive information about training data [9], [10]. The lack of explainability and privacy of modern DL systems are some of the main challenges that prevent the practical use of these powerful methods in safety-critical domains.

The field of eXplainable AI (XAI) seeks to unveil the decision-making processes of black-box models and has been thoroughly researched in recent years [11]. Depending on the area of application, a variety of different methods have been developed that attempt to explain the prediction of networks as well as their underlying decision-making process. Especially in image processing, so-called attribution methods are often used [12]. These methods generate heatmaps that highlight the areas of the input that were significantly involved in the network prediction. Therefore, these methods may or may

arXiv:2211.04110v1 [cs.LG] 8 Nov 2022

not need access to the networks' internals. While this type of explainability is widely used, it is not always adequate. In some cases, an explanation for the complete process is required. This is beyond the capabilities of attribution methods, which provide local explanations. One way of approaching global explanations is the training of architectures that yield explanations by design. Well-known representatives of this group are prototype-based systems [13]. However, the special requirements regarding the model architecture introduce additional limitations. Although it is arguable whether any XAI method allows explaining a given process in its entirety, at least the representation learned by such a network is more interpretable and related to known concepts, as opposed to networks without an intrinsic interpretable design.

Some properties of DL-based models (i.e., gradients) are of great importance for decision explanation. However, these properties also provide interfaces for the targeted retrieval of sensitive information. It has been shown that only limited model access suffices to completely reconstruct models and steal their training data [9]. This constitutes a major risk for the deployment of AI in safety-critical applications. Moreover, it is a significant threat to individuals contributing to a model's training data and users alike [14]. To prevent this, many methods have been designed to protect neural networks from attacks during the training process and to reduce data leakage during later deployment. These methods also open new opportunities for collaboration between multiple parties, allowing for better predictions based on larger amounts of data.

While XAI methods increase the transparency and intelligibility of a model's decision-making behavior, privacy-protection techniques prevent the leakage of sensitive information. However, the safe deployment of data-driven systems in safety-critical areas is only possible if one can reconcile both goals. So far, there has been no work that investigates and describes the specific impact of different privacy-preserving methods on the quality of explanations in deep learning. Understanding the trade-off between both concurring objectives is crucial to improve XAI methods and assure their correct interpretation in a privacy-preserving setting, constituting an important step in the practical applicability of DL.

This work is the first to thoroughly analyze the influence of different privacy-preserving machine learning (PPML) techniques on the explanations generated by XAI methods. In an extensive, multivariate analysis, three different privacy techniques (*Differential Privacy (DP)* [15], *Federated Learning (FedAVG)* [16], and *Differential Private Federated Learning (FedAVG-DP)*) are combined with a series of XAI attribution methods, and applied to seven different model architectures trained on eleven different datasets from various domains including document image, natural image, medical image, as well as time series analysis. The evaluation qualitatively and quantitatively highlights the varying but non-negligible impact of PPML methods on the quality of explanations.

The analysis reveals important relationships between pri-

vate training and XAI.

- *Differential Privacy* hampers the *Interpretability* of explanations.
- *Federated Learning* often facilitates the interpretation of generated explanations.
- The *Fidelity* of explanations is potentially deteriorated when using *DP*.
- The negative effects introduced by *DP* can be moderated by combining it with *FedAVG*.
- Perturbation-based XAI methods are less affected by *DP*-based training procedures.

The remainder of this paper is structured as follows. Section II gives an overview of relevant XAI methods and PPML techniques. The various datasets used throughout this study are introduced in Section III. In Section IV the complete experimental setup is outlined, followed by the presentation of the respective results. Trends and findings of this analysis are discussed throughout Section V and the manuscript is concluded in Section VI.

## II. RELATED WORK

In the following section, the most common methods in the fields of XAI and PPML will be briefly outlined. A full review of methods is beyond the scope of this work. The interested reader can find extensive reviews of XAI and PPML in [17] and [18], respectively.

### A. XAI

In the field of XAI, attribution methods are very widely used due to their versatility and comprehensibility. Attribution maps approximate the relevance of input features or feature groups to the local model decision and belong to the group of so-called *post-hoc methods* [19], which are mainly characterized by their ability to explain models that have already been trained. In 2013, the *Saliency* [20] was published as one of the first methods in this field based on the backpropagation [21] algorithm used to train networks. An extension of this method is *InputXGradient* [22], in which the coherence of input features is additionally taken into account. Other methods that work similarly to the aforementioned methods include *GuidedBackpropagation* [23] and *IntegratedGradients* [24]. All these methods are so-called *gradient-based methods* and need access to the networks' internals to compute explanations.

In contrast, *perturbation-based* methods are usually *model-agnostic*, therefore not requiring any specific model architecture to work on. The *Occlusion* [25] method removes input areas sequentially and reevaluates each manipulated input to measure the influence of single regions on the network prediction. Another subgroup of *perturbation-based* algorithms derives surrogate models from the local model behavior. *Local Interpretable Model-Agnostic Explanations (LIME)* [26], for instance, applies perturbations to an input sample to obtain a local linear model from these inputs and the respective model predictions. *Shapley Additive Explanations (SHAP)* [27] has been proposed as a related method,

**TABLE 1.** Shows the datasets used to evaluate the impact of PPML on XAI methods. The datasets cover the image, document image, medical image and time series classification.

Modality & Dataset	Domain	Train	Test	Dimensions	Channels	Classes
<b>Time Series</b>						
Anomaly Detection	Synthetic	50,000	10,000	50	3	2
Character Trajectories	Communication	1,422	1,436	182	3	20
ECG5000	Medical	500	4,500	140	1	5
FordA	Manufacturing	3,601	1,320	500	1	2
Wafer	Information	1,000	6,164	152	1	2
<b>Images</b>						
RAF-Database	Facial Expressions Recognition	12,271	3,068	224 × 224	3	7
Caltech-256	Natural Image Classification	24,485	6,122	224 × 224	3	256
ISIC	Medical Image Analysis	26,521	2,947	224 × 224	3	8
SCDB	Synthetic	6,000	1,500	224 × 224	3	2
RVL-CDIP	Document Analysis	320,000	40,000	224 × 224	3	16
Tobacco3482	Document Analysis	2,782	700	224 × 224	3	10

with additional constraints based on game theory to provide certain mathematical guarantees.

For the sake of completeness, it should be mentioned that there are several other approaches besides the discussed *post-hoc* attribution methods. *Prototype-based* [13], *patch-based* [28], and *concept-based* methods [29], for instance, have also been used previously.

### B. PPML

Different techniques have been developed to protect data and model privacy during the training and in the subsequent inference phase. Anonymisation techniques (e. g., K-anonymity [30]) were among the first approaches developed to ensure privacy in model training. In the meantime, there have been outstanding breakthroughs in the area of privacy attacks. Membership [31] or model inversion attacks [9] allow reconstructing training data with extremely limited access to the models. Therefore, simplistic techniques such as anonymization are no longer sufficient.

A promising training technique that leads to a high degree of privacy for data and model is *Homomorphic Encryption* [32]. However, this method is rarely applicable to modern DL-based systems, due to its massive computational overhead. Another, more frequently used technique is *Differential Privacy (DP)* [15]. Here, a certain amount of noise is added to the training signal of deep networks to prevent its parameters to capture information held by specific training samples but instead focus on the general characteristics of the whole population. One advantage of this method is that it can be applied to a wide variety of architectures and requires only minimal changes to the training setup. Another prominent technique used to account for data-privacy is *Federated Learning (FedAVG)* [16]. In *FedAVG*, local models are trained on a data-owners subset of training samples, and only the locally computed gradients are sent to a centralized server. There, the average is calculated to obtain a global

model. This way, sensitive data does not need to leave the institution, but multiple institutions can collaborate to leverage a bigger training set for the global model. Moreover, *FedAVG* can be combined with *DP* to prevent the risk of data leakage from model gradients. Out of the many other privacy techniques [14], *DP* and *FedAVG* stand out as the most commonly used.

To this date, both XAI and PPML have mostly been considered in separation by the research community. Few works [33], [34] have made the first attempts at combining explanation and privacy preservation in a single framework. However, there is still a lack of insight regarding the exact influence private training has on the generated explanations in deep learning.

### III. DATASETS

To comprehensively analyze the impact of privacy-preserving methods on explanations, a variety of different datasets from different domains in time series and image analysis were utilized, as listed in Table 1.

#### A. TIME SERIES DATASETS

The first modality that is evaluated uses time series data. Time series data is usually acquired using different sensors and differs from image data concerning various characteristics such as the locality constraints and the dependence on a sequential order.

With the exception of the *Anomaly Detection* dataset [35], the datasets for the time series analysis come from the UEA & UCR repository [36]. This selection includes both univariate and multivariate time series with different numbers of classes. The *Anomaly Detection* dataset and the *FordA* dataset consider the task of anomaly detection. The *Anomaly Detection* dataset deals with point anomalies and the *FordA* dataset with sequence anomalies. The point anomalies are very interpretable for humans as in their case the data is more or less noise and contains a large peak that indicates

the anomaly. Even without the annotation, it is possible to understand whether the explanation for such a sample is correct or not. This is not the case for the *FordA* data as the sequences are very long and there is no annotation. In this dataset, the anomaly can be a long part of the sequence that varies from the expected behavior. The *Character Trajectories* dataset was selected as it is possible to transform it back to the 2d input space to understand the explanation. It consists of three channels covering the acceleration within the x and y direction and the pen force. Therefore, it is a real-world dataset that enables precise identification of whether an explanation is good or not. In addition, it is important to mention that the dataset size of the time series datasets differs significantly, to properly represent the influence of data volume.

## B. IMAGE DATASETS

### 1) Natural Image Datasets

Image datasets range from specialized domains such as facial recognition, medical image analysis, and document analysis to toy datasets covering a varying number of classes, channels, and dataset sizes. The Real-world Affective Faces Database (*RAF-Database*) [37] is a collection of 15,339 face images with crowd-sourced annotations. It classifies images in one of the seven facial expressions (Surprise, Fear, Disgust, Happiness, Sadness, Anger, Neutral). *Caltech-256* [38] is a natural images classification dataset comprising of a total of 30,607 labeled images with 256 unique object classes.

### 2) Medical Image Datasets

The International Skin Imaging Collaboration (ISIC) provides a large publicly accessible library of digital skin images<sup>1</sup> and hosts annual challenges. The *ISIC* dataset used in this work is a cleaned combination of all ISIC challenge datasets. The datasets have been merged and freed from duplicates according to the recommendations in [39]. All images are labeled as either Melanoma (*MEL*), Nevus (*NV*), Basal Cell Carcinoma (*BCC*), Actinic Keratosis (*AK*), Benign Keratotic Lesion (*BKL*), Dermatofibroma (*DF*), Vascular Lesion (*VASC*) or Squamous Cell Carcinoma (*SCC*). The classification is based on complex combinations of distributed and overlapping biomarkers, posing particular challenges for the explanation of automated decisions. The seven-point checklist criteria dataset (*Derm7pt*) proposed in [40] consists of clinical and dermoscopic images of 1,011 skin lesions with extensive annotation. In this work, only the subset of dermoscopic images along with the respective diagnosis annotations for pre-training of ISIC classifiers is considered in experiments involving *DP*. The *SCDB* [41] dataset is a synthetic toy dataset inspired by the problems of skin lesion analysis. Images are classified into one of two classes based on the combinations of shapes present in a base shape, depicting the skin lesion. The shapes can be overlapping and redundant, but classification evidence is sparse and localized.

Along with the class label, each image is supplemented by shape annotation maps, serving as ground truth explanations.

### 3) Document Image Datasets

Business documents are a fundamental component of modern industry. Recent advances in deep learning have sparked a growing interest in automating document processing tasks such as document search, and extraction of document information. However, business documents often contain highly personal user data and sensitive information pertaining to a company's intellectual property, which makes the secure application of Deep Learning in this area a major concern. On the other hand, deep learning-based decision-making processes have been shown to be susceptible to learning biases in the data [42]. An example of such a system involves automatically analyzing resumes to make hiring decisions, which may lead to discrimination against women or members of minority groups. Explainability of such systems is therefore of paramount importance for their safe and practical deployment.

In order to analyze the interdependence between PPML and XAI for document domain, two popular document benchmark datasets are utilized in this study. *RVL-CDIP* [43] is a large-scale document dataset that has been extensively used as a benchmark for document analysis tasks. The dataset contains a total of 400,000 labeled document images with 16 different categories and consists of training, testing, and validation split of 320,000, 40,000, and 40,000 images respectively. *Tobacco3482*<sup>2</sup>, is another popular but small-scale dataset with 3482 labeled document images. Since there is no split defined for this dataset, the training and test, and validation splits of sizes 2,504, 700, and 278 images were defined. Since both of these datasets are subsets of a bigger dataset, there exists some overlap between them. Therefore, for all the experiments, the overlapping images were removed from the *RVL-CDIP* dataset. The two datasets were used in combination to analyze the effects of transfer learning on the privacy and explainability aspects of the models.

## IV. EXPERIMENTS & RESULTS

A broad experimental basis, covering various domains, applications, and configurations is necessary, to make general statements about the impact of privacy techniques on explanations. Therefore, a selection of state-of-the-art classifiers is trained on a range of different datasets and applications covering both time series and image domains. Each combination of model and dataset is trained in four different settings, including training without privacy (*Baseline*), with differential privacy (*DP*), federated training (*FedAVG*) and federated training with client-side differential privacy (*FedAVG-DP*). Different explanation methods are finally applied to every model instance to compare their generated explanations.

Evaluating explanations and judging their quality is a common problem not only in XAI research [44], but also in

<sup>1</sup>The ISIC Archive is accessible at <https://www.isic-archive.com>.

<sup>2</sup><https://www.kaggle.com/patrickaudriaz/tobacco3482jpg>.

the social sciences [45]. Multiple evaluation dimensions have to be considered to make clear statements about the impact of privacy-preserving model training on the explainability of DL-based models. Human-centered evaluation is laborious and requires domain experts. Instead, functionality-grounded methods are best suited for the domain- and dataset-wide fair comparison and quality assessment of XAI and are therefore utilized throughout this study.

In the experiments, the focus lies on the two main properties of explanations as defined in [44], namely their **Fidelity** and **Interpretability**. *Fidelity* measures soundness and completeness to ensure that explanations accurately reflect a model's decision-making behavior. *Interpretability* refers to the clarity, parsimony, and broadness of explanations, and therefore describes factors related to the ease of communication on the interface of machines and humans. Functionality-grounded methods make use of formal mathematical definitions as proxies of perceived interpretability.

In this section, details of the experimental setup are described, and evaluation metrics are introduced. Afterward, the results and analysis of all setups are thoroughly presented.

#### A. EXPERIMENT SETUP

Baseline networks were trained using the standard SGD or ADAM optimizers with varying numbers of epochs per dataset, to ensure convergence. For all other settings, training and privacy hyperparameters have been manually tuned to find a good trade-off between privacy and model performance matching the baseline. This is important to guarantee a sufficiently fair comparison between the methods since a significantly worse network would also show worse attribution results. However, all models were trained with overall comparable settings. Moreover, fixed seeds were used to ensure reproducibility. The reported performances correspond to the accuracies of the model performing best on the test datasets. In time series analysis, *InceptionTime* [46] and *ResNet-50* [47] were used as representative networks, since they achieve state-of-the-art performances for the utilized datasets. *ResNet-50*, *NFNet* [48] and *ConvNeXt* [49] have been used for classification of natural and medical images. Since document images differ significantly from natural images, a different set of models has been used for this domain, including *AlexNet* [50], *VGG-16* [51], *ResNet-50*, *EfficientNet* [52], and *ConvNext*, which have shown the best performance in the past. The training data was split between training and validation with a factor of 0.9, wherever no validation dataset had been provided.

Some XAI methods pose specific requirements on the model architecture or training procedure, complicating the application of privacy-protection techniques. Therefore, this work solely focus on the commonly used *post-hoc* explanations. Different attribution methods vary considerably in their realization and their associated underlying assumptions. Therefore, it was decided to apply a broad range of diverse methods differing in their implementations and theoretical foundations. The work covers a total of nine XAI

methods, including gradient-based Saliency [20], InputX-Gradient [22], GuidedBackpropagation [23], IntegratedGradients [24], DeepSHAP [27], and DeepLift [53], but also gradient-free methods such as Occlusion [25], LIME [26] and KernelSHAP [27].

#### B. EVALUATION METRICS

The *Fidelity* of the explanations is quantified using *Sensitivity* [54], *Infidelity* [54], *Area Over the Perturbation Curve* [55], and *Ground Truth Concordance* evaluation metrics. When measuring the *Sensitivity*, insignificant perturbations are applied to the input and the change in the attribution map is measured. Small changes in the input should not result in large changes in the attribution map. Thus, a smaller value is better. *Infidelity* applies significant perturbations to the input and the corresponding attribution map. The value measures the mean-squared error between the perturbed heatmap and the difference in the predictions of perturbed and unperturbed input. The *Area Over the Perturbation Curve (AOPC)* measures the alignment between an attribution map's relevance values and the effect of perturbing the corresponding input regions on the model prediction. Intuitively, removing features with lower importance should affect the prediction less than the deletion of important features. The *AOPC* is a scalar value computed by integrating the curve over the impact of consecutive perturbations with decreasing attribution importance, relative to random perturbations. Therefore, higher values indicate an attribution map's meaningfulness. Most fidelity measures evaluate the degree to which an explanation is faithful to the local model behavior. Whether the explanation is human-aligned, on the other hand, can only be evaluated with ground truth explanations available. The *Ground Truth Concordance* was measured by computing the overlap between the highest attributed input region with the segmentation maps of the ground truth explanations. The real-valued attribution maps are binarized to allow comparison with the binary segmentations. The overlap is computed for increasing binarization thresholds from zero to one. The overlap is quantified through the *Intersection over Union (IoU)*.

The *Interpretability* of explanations was measured using the *Continuity* metric. In general, humans have difficulties interpreting information that is both high dimensional and scattered. *Continuity* is defined as the sum of the absolute changes between two consecutive importance scores in an attribution map. For time series, the continuity is the absolute change between each subsequent point in a sequence whereas in the image domain, the absolute changes are measured and aggregated separately in X and Y directions. Better *Interpretability* is indicated by lower continuity scores.

All evaluation metrics were computed on the respective test datasets. The attribution maps were normalized to have zero mean and unit standard deviation for a fair comparison across methods and different privacy types. Due to computational and time restrictions, the influence of PPML on attribution methods was quantified using a subset of the

respective test sets, limited to a maximum of 1,000 examples. This work assumes that 1,000 randomly selected examples represent a sufficient quantity to generalize the findings to the complete test datasets.

### C. CRITICAL DIFFERENCE DIAGRAMS

Intuitive visualization of high-dimensional data is particularly challenging when the data origins from multiple distinct configurations as in this case. Critical Difference (CD) diagrams, proposed by Demšar in 2006 [56], allow the high-level visualization of complex experimental data intuitively and were therefore chosen to present most quantitative results. Their ability to condense ordinal information across different datasets, models and attribution methods extracts relevant information and helps to pick up universal trends in a benchmark study. Moreover, the method includes statistical tests, indicating the data's significance.

CD diagrams report the average rank of a given item, in a series of different settings. If the statistical significance for two distinct items is not guaranteed, these items are connected by a horizontal line and referred to as a "clique". In this study, the Friedman test is used to decide the statistical significance of a group of different observations. The Holm-adjusted Wilcoxon's signed rank test is then applied for post-hoc analysis as suggested in [57]. For all statistical analyses, an  $\alpha$  of 0.05 is assumed.

### D. IMPACT ON MODEL PERFORMANCE

Applying privacy-preserving training techniques for DL-based models can have a very diverse impact on their test performances. The severity depends on multiple factors including model architecture, type of dataset, as well as the various hyperparameters for model and private training. Table 2 shows the results on the respective test sets for all experiment configurations when trained with different private training techniques. Results are sorted by domains and datasets to provide a better overview.

Even in privacy-preserving training settings, all models converged and demonstrated acceptable accuracies. However, the best accuracies were usually achieved in *Baseline* or *FedAVG* settings. Across all domains, it can be observed that *DP* has a considerable impact on the models' test performances. For some configurations, a higher  $\epsilon$ -value was required to achieve comparable results (e. g., *NFNet* and *ConvNeXt-B* for *ISIC*). However, no consistent pattern indicating higher robustness of one model architecture over another, against noise introduced by *DP*, is obvious. In contrast to *DP*, *FedAVG* always resulted in significantly lower performance losses. The combination of *FedAVG* and *DP* almost exclusively resulted in a lower performance considering a comparable  $\epsilon$ -value.

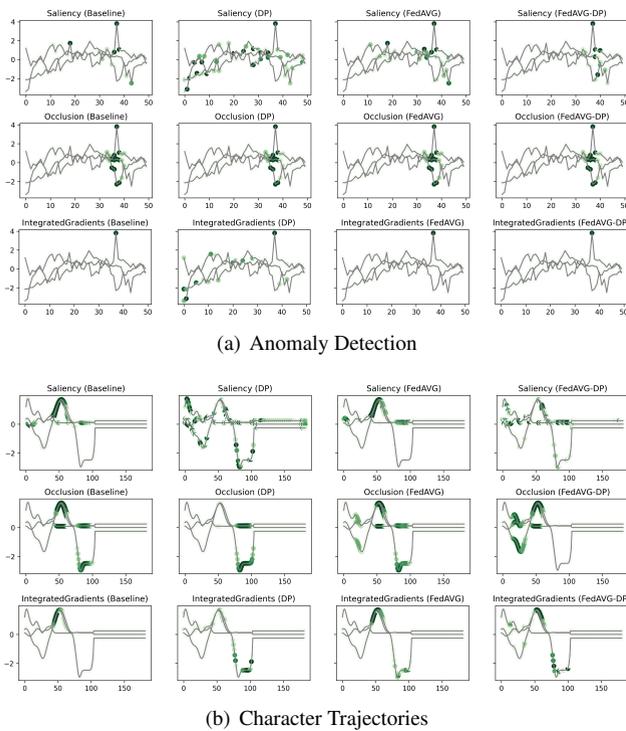
The results from the time series domain for most datasets indicate that *InceptionTime* is usually affected slightly less by private training, in direct comparison with *ResNet-50*. The only exception is the *Anomaly* dataset, which experienced almost no performance loss with *ResNet-50*. One possible

**TABLE 2.** Test accuracies on all datasets for different architectures and privacy-preserving settings, divided by application domain. For configurations containing *DP* or *FedAVG*, the  $\epsilon$  and  $n_c$  values are provided, respectively.

Datasets & Models	Acc <sub>Baseline</sub>	Acc <sub>DP</sub> / $\epsilon$	Acc <sub>FedAVG</sub>	Acc <sub>FedAVG-DP</sub> / $\epsilon$	
Time Series Analysis	<b>Anomaly</b>				
	InceptionTime	98.74	92.87 / 5.0	98.77	$n_c = 4$ 89.50 / 5.0
	ResNet-50	98.70	97.02 / 5.0	98.60	97.36 / 5.0
	<b>Character Trajectories</b>				
	InceptionTime	99.44	91.85 / 5.0	98.82	$n_c = 4$ 87.26 / 50.0 68.73 / 5.0
	ResNet-50	99.44	85.03 / 5.0	98.19	82.10 / 50.0 59.19 / 5.0
	<b>ECG5000</b>				
	InceptionTime	94.38	89.07 / 5.0	93.36	$n_c = 4$ 89.29 / 5.0
	ResNet-50	94.16	89.64 / 5.0	92.78	88.87 / 5.0
	<b>FordA</b>				
	InceptionTime	95.61	92.88 / 5.0	97.70	$n_c = 4$ 94.17 / 50.0 91.43 / 5.0
	ResNet-50	94.32	86.14 / 5.0	93.94	87.12 / 50.0 76.44 / 5.0
Document Analysis	<b>Wafer</b>				
	InceptionTime	99.22	89.21 / 5.0	97.81	$n_c = 4$ 89.21 / 5.0
	ResNet-50	98.75	89.21 / 5.0	89.21	89.21 / 5.0
	<b>RVL-CDIP</b>				
	AlexNet	87.90	70.30 / 4.5	85.54	$n_c = 8$ 61.35 / 5.3
	VGG-16	91.00	69.67 / 4.4	89.41	62.38 / 5.4
	ResNet-50	90.50	72.55 / 5.0	88.25	68.85 / 8.8
	Efficientnet-B4	92.60	60.20 / 4.2	90.59	45.09 / 6.5
	ConvNeXt-B	93.64	75.60 / 3.7	92.60	73.23 / 7.7
	<b>Tobacco3482</b>				
	AlexNet	89.57	86.14 / 3.9	91.85	$n_c = 4$ 85.71 / 8.0
	VGG-16	94.14	85.14 / 4.9	93.99	87.00 / 7.5
ResNet-50	92.57	75.42 / 2.7	92.14	78.43 / 7.3	
Efficientnet-B4	94.42	89.42 / 4.4	93.99	88.57 / 8.0	
ConvNeXt-B	94.71	87.14 / 4.8	94.85	85.42 / 6.0	
Natural Images	<b>Caltech-256</b>				
	ResNet-50	87.30	59.00 / 5.0	87.97	$n_c = 4$ 61.82 / 35.94
	NFNet	88.50	60.17 / 5.0	91.39	75.28 / 35.13
	ConvNeXt-B	91.57	78.32 / 5.0	93.74	79.86 / 14.85
	<b>RAF-Database</b>				
	ResNet-50	81.91	67.42 / 5.0	80.82	$n_c = 4$ 64.43 / 13.23
NFNet	83.96	69.68 / 5.0	82.82	69.75 / 13.33	
ConvNeXt-B	86.63	69.32 / 4.79	88.23	71.97 / 13.23	
Medical	<b>ISIC</b>				
	ResNet-50	86.08	71.09 / 4.66	82.15	$n_c = 4$ 70.89 / 8.09
	NFNet	90.16	77.23 / 14.61	86.90	71.39 / 18.28
ConvNeXt-B	87.20	69.63 / 14.02	81.87	68.78 / 30.69	
Synthetic	<b>SCDB</b>				
	ResNet-50	90.20	86.20 / 4.60	92.19	$n_c = 4$ 85.13 / 60.00
	NFNet	94.40	88.33 / 4.46	94.87	85.19 / 18.27
ConvNeXt-B	92.46	85.80 / 13.19	92.40	87.07 / 30.08	

explanation for this is the advanced architecture of *InceptionTime* including residual connections and inception modules. This enables the *InceptionTime* to be more robust against noise and outliers. All datasets in the time series domain, with the exception of *Anomaly*, *ECG5000*, and *Wafer*, considerably suffered from the combination of *DP* with *FedAVG*. For *Character Trajectories* and *FordA*, the  $\epsilon$ -value had to be increased to achieve adequate results.

The image domain indicates similar findings. Since the models for the *Tobacco3482* dataset were trained after being pretrained on the *RVL-CDIP* dataset, the models were able to



**FIGURE 1.** Shows the change in the attribution for *ResNet-50* for three selected samples of the *Anomaly Detection* and *Character Trajectories* dataset, respectively. *DP*-based training techniques tend to add additional noise and alter the explanation. *FedAVG*, by contrast, is closer to the original attribution of the baseline.

achieve higher performance even with *DP* and *FedAVG-DP*. However, despite pretraining on *Derm7pt*, the impact of *DP* on *ISIC*-trained models is still considerable. For all datasets in the image domain, a moderately higher  $\epsilon$ -value was also experimented with, to improve the performance of the models for *DP* and *FedAVG-DP* cases. However, it did not seem to provide any significant improvement in most of the cases. It is worth noting that for larger datasets (i. e., *RVL-CDIP*, *Caltech-256*, *RAF-Database* and *ISIC*), *DP* and *FedAVG-DP* severely degraded the performance of the models whereas the performance for *FedAVG* is still comparable to the *Baseline* setting.

### E. GENERAL IMPACT ON EXPLAINABILITY (QUALITATIVE)

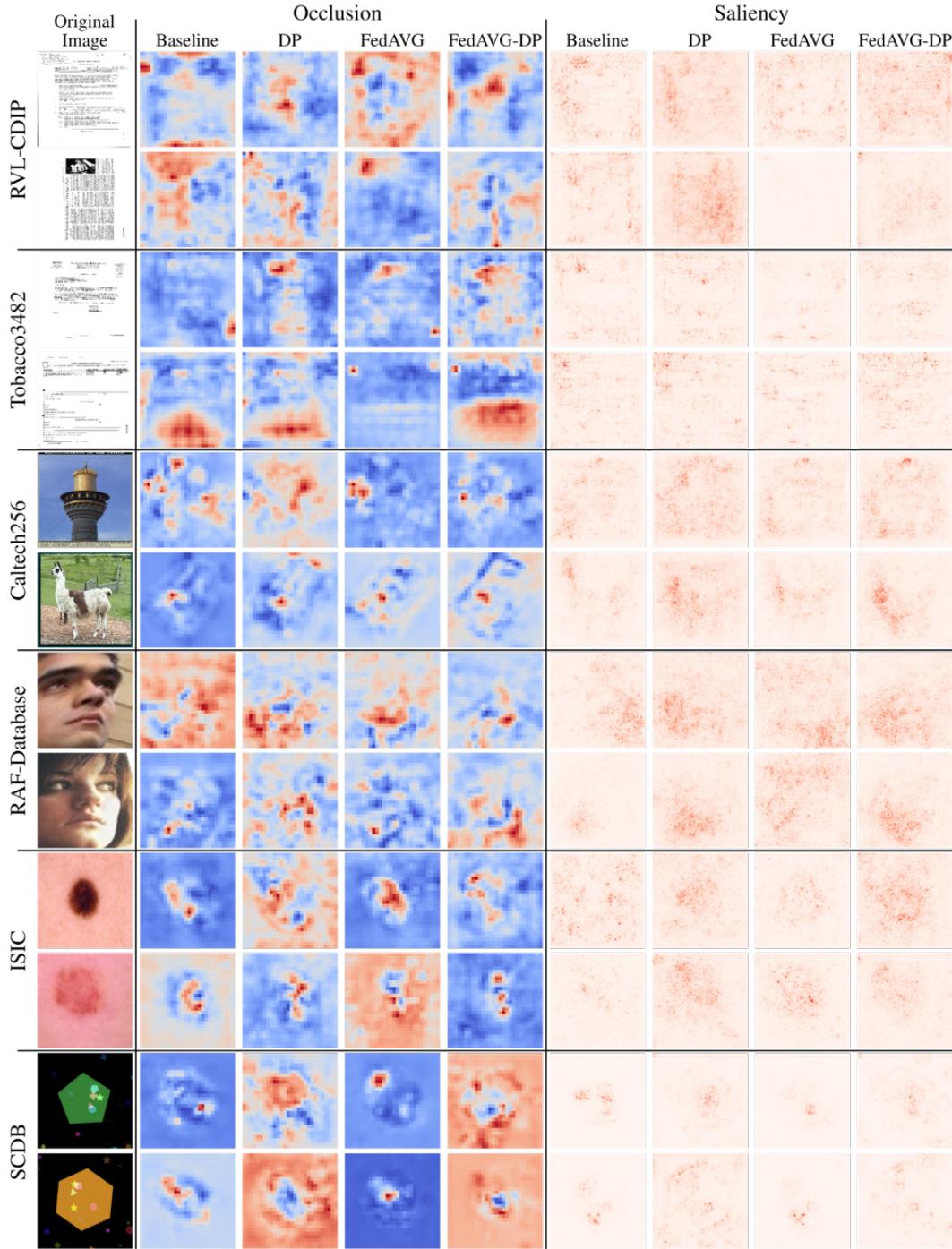
A visual inspection of individual explanations gives a first impression of the influence privacy-preserving techniques can have on the trained models. These local impressions are then further validated on the dataset level through the qualitative analysis of summary statistics.

#### 1) Individual Analysis

Figure 1 shows explanations generated for the *Anomaly Detection* and *Character Trajectories* datasets in the time series domain. For the *Anomaly Detection* dataset, it can be observed that there is a general overlap between the explanations from different training settings, always highlighting the

anomaly. In some cases, *DP* increases the amount of noise in the signal's relevance around the anomaly, yielding unclear and misleading explanations by highlighting distant points which do not correspond to the anomaly at all. However, this is not the case when additionally adding *FedAVG* in the *FedAVG-DP* setting. By contrast, *DP*-trained models show remarkable deviations from the original *Baseline* explanation when trained on the *Character Trajectories* dataset. This observations holds not only true for *DP*, but also *FedAVG-DP* settings. *FedAVG*, on the other hand, shows explanations close to the *Baseline* setting, with only minor deviations.

Figure 2 shows samples from all image datasets along with the generated *Occlusion* and *Saliency* explanations from *ResNet-50* models trained with and without privacy techniques. It can be seen that different training settings yield heatmaps that visibly differ. However, the areas of highest relevance roughly overlap for most samples. Particularly for *Saliency* attributions, it can be observed that *DP* and *FedAVG-DP* often add additional noise to the generated explanations. In some cases, this can also be observed in *Occlusion*. This is particularly striking in the *SCDB* samples in the last two rows, where both *DP*-based methods highlight regions outside the decision-relevant area. Another striking example is the second sample from *RVL-CDIP*. It can be seen that both *FedAVG* attribution maps present smoother and more focused heatmaps pointing to a specific location on the image. On closer inspection of the samples from document datasets, it was found that *FedAVG* prominently focused on specific class-relevant cues such as dates, titles, figures, etc. Moreover, comparing *DP* with *FedAVG-DP* attribution maps, it can be observed that the addition of federated training leads to less noise in the attribution for some samples. When observing the *SCDB* sample in the last row, it can be observed that for *Occlusion*, *Baseline* highlights both rectangle and star shapes, whereas *FedAVG* only focuses on the rectangle shapes. In *SCDB*, rectangles are exclusive markers for class two. However, both star and star markers are also part of the decision-relevant shape combination. *FedAVG* appears to have focused only on the single relevant marker, whereas in the *Baseline* setting, multiple relevant markers were highlighted. This is also evident from the *SCDB* sample in the second last row, where *FedAVG* successfully focused on the ellipse shape, which is exclusive for class one. It has to be noted that the *Occlusion* attribution map shows two relevant regions, where the most relevant region highlights the edge of the big, green pentagon. This could be attributed to a limitation of *Occlusion*, corresponding to distraction due to the generation of out-of-distribution samples during perturbation. Interestingly, *Occlusion* attribution maps for *Baseline* and *FedAVG* show different decision-relevant cues for the last row's sample, as compared to the corresponding *Saliency* maps. The former highlight rectangle and star shapes, whereas the latter most prominently highlight the star marker on the lower part of the image. As already mentioned, star markers are also decision-relevant for that sample. However, they are no exclusive markers and could also indicate



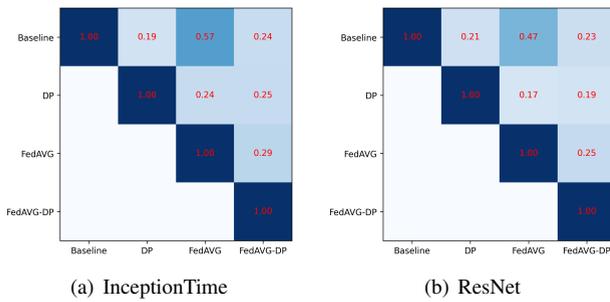
**FIGURE 2.** Examples of perturbation- and gradient-based attribution maps computed on models trained in varying privacy configurations. Two random, correctly classified samples are provided per dataset. *Occlusion* and *Saliency* attribution maps are computed on *ResNet-50* models. Red and blue regions indicate positive and negative relevance, respectively.

shape combinations appearing in class one. Nevertheless, it has to be mentioned that for some individual samples, these observations do not apply. For the second sample of *ISIC* and the first sample of *RAF-Database*, for instance, *Baseline* and *FedAVG* did not yield clearer heatmaps as compared to the *DP*-based approaches. To capture overall trends and characteristics in attribution maps of different configurations,

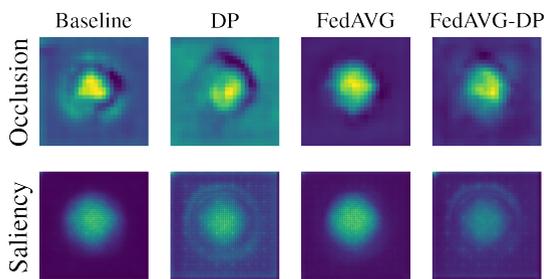
further analysis on dataset-level statistics of the generated attribution maps was performed.

## 2) Dataset-wide analysis

Figure 3 shows the Pearson correlation of the explanations generated by different training settings for the *Anomaly Detection* dataset. Therefore, the correlation across the different



**FIGURE 3.** Shows the average Pearson correlation of the attribution maps compared between the different privacy approaches for the Anomaly Detection dataset. *FedAVG* shows a higher similarity to the *Baseline* setting, as compared to the *DP*-based approaches.



**FIGURE 4.** Average attribution heatmaps for samples from *SCDB* computed for each of the non-private / private approaches.

training approaches was computed using all available attribution maps. Precisely speaking, the attribution between the corresponding attribution maps was calculated and the average over the number of samples was taken. The final correlation shows the score averaged over the attribution methods and the samples. For both architectures, it is evident that the privacy methods significantly change the produced attribution maps. However, *FedAVG* yields significantly higher correlation to the *Baseline* setting as compared to the *DP*-based approaches. Moreover, it is surprising that the correlation between *DP* and *FedAVG-DP* is rather low. The remaining correlation matrices are excluded as they showed similar results to the presented and do not provide any additional information.

*SCDB* is a synthetic dataset that, by design, carries all relevant information in the center of the image. Figure 4 shows a visual comparison of the overall impact of privacy on the distribution of attribution in the explanations. For each setting, the attribution maps are averaged across all test samples to highlight the frequency with which regions were attributed as relevant. It is evident that involving *DP* during training drastically increased the diffusion of attribution values, also including image areas that are not related to the actual classification task. Combining *DP* with *FedAVG*, on

the other hand, led to a moderation of the noise added by *DP*. Interestingly, the results show that *FedAVG* alone usually results in the cleanest and most focused heatmaps, even improving the baseline. As the remaining datasets possess less spatial standardization, it is not trivial to interpret their results in the same way.

This qualitative analysis already drew an interesting initial picture, proving that to some degree, any privacy-preserving training technique has an impact on the generated explanations. Furthermore, the results suggest that *DP*-trained models generate explanations that tend to be noisier and cover potentially unimportant regions, harboring the danger of misleading the explainer. However, it is unclear whether noise added by *DP* only concerns the explanations, or whether this reflects the model decision (Fidelity). First results also indicate that the *FedAVG* approach can even improve explanations, leading to more focused, and meaningful explanations in some instances.

## F. GENERAL IMPACT ON EXPLAINABILITY (QUANTITATIVE)

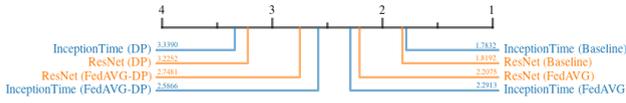
The quantitative analysis serves as a means to further verify the findings from the previous section and allows the investigation of whether privacy-preserving techniques impact only the explanations, or also the underlying model behavior.

### 1) Continuity

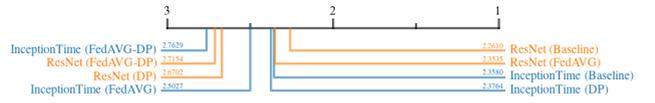
Measuring the continuity of an explanation helps to understand how difficult the interpretation of an explanation might be for an explainer. Humans usually struggle when confronted with high-dimensional, diffuse data.

The continuity for time series data is defined as the sum of the absolute changes between each pair of subsequent points within the attribution map. For image attributions, continuity is computed as the sum of absolute gradients in both spatial directions of the attribution map. A smoother map results in a lower continuity. Figure 5 shows the CD diagrams for the *Continuity* across all domains. For each domain, the ranked results are averaged over all datasets and attribution methods.

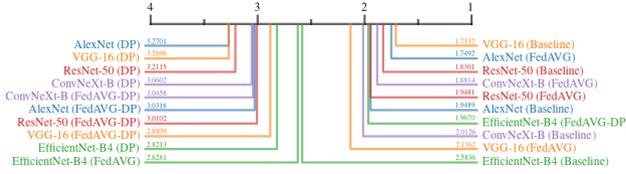
The results for all domains clearly show that *Baseline* and *FedAVG* settings yield better continuity scores as compared to the *DP*-based approaches. This confirms that *DP*-based private models generate significantly more discontinuous attribution maps compared to *Non-DP* training techniques. The only outlier to this observation is *EfficientNet-B4* trained on document images in Figure 5(b), where surprisingly *FedAVG-DP* achieved the highest rank. Comparing only *Baseline* and *FedAVG* models, it can not be clearly stated whether one is better than the other, as this seems to be highly dependent on the exact model architecture and domain combination. Moreover, *FedAVG-DP* achieved better ranks as compared to *DP* in most configurations across all domains. *DP* and *Non-DP* approaches even show a clear visual separation in the CD diagram in most cases (i. e., document, natural, and medical). For document image datasets, there are



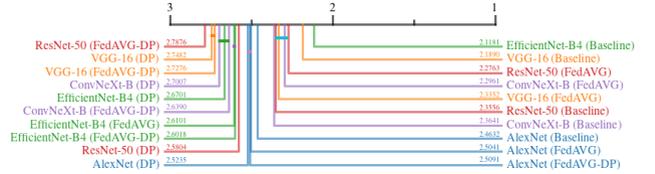
(a) Time series datasets



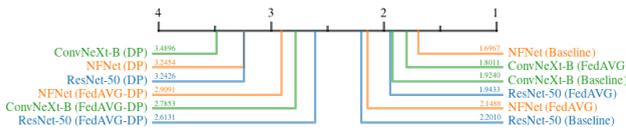
(a) Time series datasets



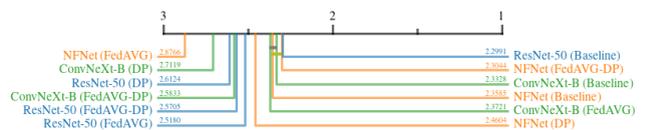
(b) Document image datasets



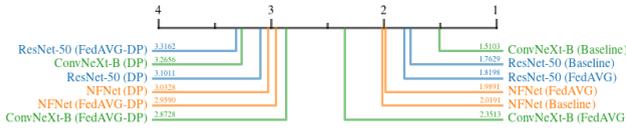
(b) Document image datasets



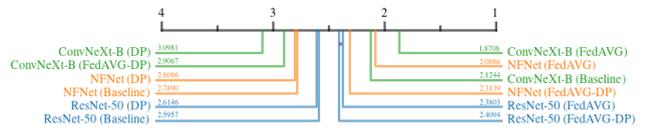
(c) Natural image datasets



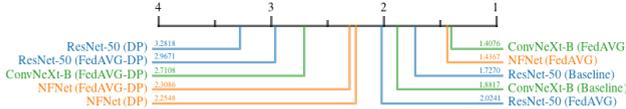
(c) Natural image datasets



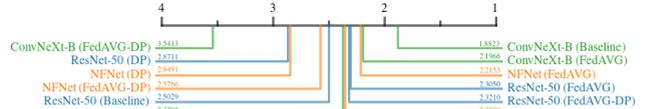
(d) Medical image datasets



(d) Medical image datasets



(e) Synthetic image datasets



(e) Synthetic image datasets

**FIGURE 5.** Critical difference diagrams for the Continuity of models trained on datasets from different domains. Privacy results in less continuity and therefore noisier explanations.

**FIGURE 6.** Critical difference diagrams for the AOPC of models trained on datasets from different domains.

only minor differences within the ranks of *DP* and *Non-DP* regions, making it very hard to draw clear conclusions.

### 2) Area over the Perturbation Curve

The *AOPC* measures how removing features deemed relevant by the explanation affects local model predictions. This provides important insights into the Fidelity of the explanations. Intuitively, removing features with lower importance should affect the prediction less, whereas the deletion or perturbation of important features should result in significant prediction changes. In this experiment, features were removed sequentially starting with the most important, as per the attribution map.

Figure 6 shows all critical difference diagrams for the *AOPC* measure. Over all domains, the most prevalent pattern is that of *Non-DP*-based settings occupying the higher ranks. In the time series domain, *ResNet-50* shows the clear superiority of *Baseline* and *FedAVG* compared to *DP* and *FedAVG-DP*. However, the results from the image domain

indicate that a clear superiority of *FedAVG* over *Baseline* can not be reported. In contrast to other domains, *FedAVG*-based approaches on average achieved higher ranks on medical and synthetic images. For *InceptionTime*, *DP* surprisingly achieved almost similar performance as compared to *Baseline*. Apart from this outlier, *DP* almost exclusively ranked last in direct comparison with all other training settings. The presented results suggest that adding *Differential Privacy* during training decreases the explanation’s fidelity.

### 3) Infidelity

The *Infidelity* measure provides information about an explanation’s fidelity by evaluating a model’s adversarial robustness in regions of varying explanation relevance. Perturbations are both applied to the attribution map and the input image while comparing the predictions of the unperturbed and noisy input. It is expected that the perturbation of a more important feature leads to a larger change in prediction.

Figure 7 shows the *Infidelity* ranks for all configurations.

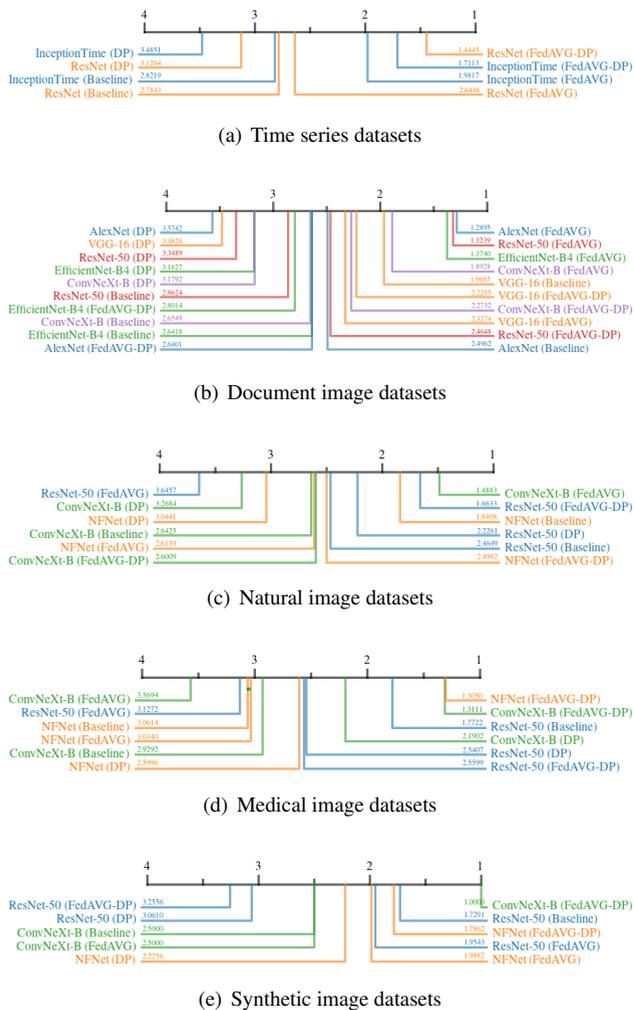


FIGURE 7. Critical difference diagrams for the Infidelity of models trained on datasets from different domains.

In the time series domain, approaches involving *FedAVG* achieved the highest scores. Interestingly, the addition of *DP* to *FedAVG* settings resulted in higher scores, whereas the sole use of *DP* during training led to the worst outcomes. For all image domains, again, results depend on the architecture and the domain. However, there is an overall tendency similar to the results of the time series domain with *DP* being the worst. Furthermore, *FedAVG* approaches being the best performing approach for time series datasets. The *Baseline* setting ranked highest in several configurations, such as *ResNet-50* in medical and synthetic images, and *NFNet* in natural images. This indicates that *FedAVG* increases adversarial robustness and Fidelity while *DP* alone leads to lower Fidelity.

#### 4) Sensitivity

In contrast to the *Infidelity*, *Sensitivity* quantifies the fidelity by perturbing the input directly. The change in the generated explanation is measured before and after the input is insignificantly perturbed. Small changes in the input should not result in large changes in the attribution map.

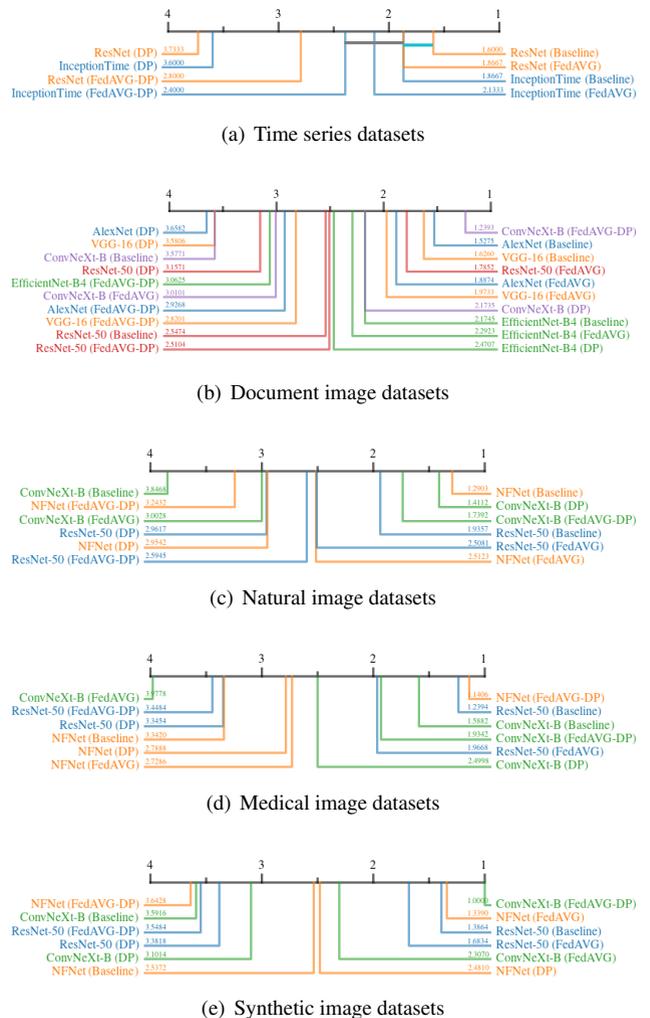


FIGURE 8. Critical difference diagrams for the Sensitivity of models trained on datasets from different domains.

Figure 8 shows the *Sensitivity* ranks for all configurations. For the time series domain, Figure 8(a) shows a very clear ranking with *Baseline* and *FedAVG* being superior to *FedAVG-DP*, followed by *DP*. However, no clear statistical distinction can be made between *Baseline*, *FedAVG*, and *FedAVG-DP* for *InceptionTime*, and for *Baseline* and *FedAVG* for *ResNet-50*. The superiority of *Non-DP*-based over *DP*-based approaches is further confirmed by seven more configurations within the different image domains, including *ResNet-50* in natural, medical, and synthetic images. Interestingly, both *DP*-based methods ranked highest in combination with *ConvNeXt* in all image-domains except for medical imaging.

#### 5) Ground Truth Concordance

Measuring the concordance of attribution maps and ground truth explanations is the best way to ensure the truthfulness of explanations, but comes with some limitations. *Ground Truth Concordance* can only be computed on synthetically constructed datasets without ambiguous decision paths. There-

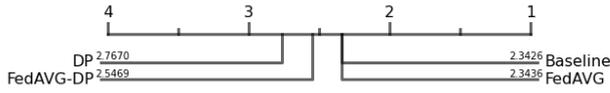


FIGURE 9. Critical difference diagram for *Ground Truth Concordance* on SCDB data for different privacy-preserving settings. *FedAVG* and *Baseline* settings clearly outperform *DP*-based techniques.

fore, the *SCDB* dataset was utilized which provides segmentation maps for the different visible shapes to construct ground truth explanation maps containing only decision-relevant shapes for each image. All attribution maps are blurred before computing the concordance, to moderate the drawback of gradient-based methods, which generate noisier explanations by design.

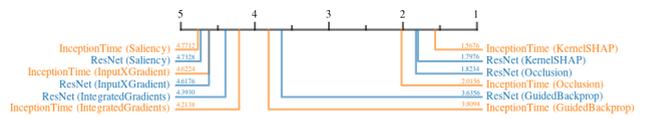
Figure 9 shows the critical difference diagram for the *Ground Truth Concordance* computed over all attribution methods for the *SCDB* dataset. The ranking clearly indicates the superiority of *Baseline* and *FedAVG* settings over *DP*-based training techniques. Moreover, it can be observed that the addition of *FedAVG* to the *DP*-trained setting alleviates the divergence from the ground truth explanations. Overall it can be noted that *DP*-based methods indeed reduce the fidelity of models, while there is promising evidence that a *DP*-based training in a federated constellation can mitigate its effects to a certain degree.

G. IMPACT OF NOISE ON DIFFERENT SETTINGS

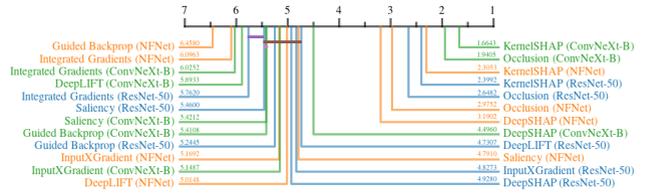
The results so far suggested that the introduction of *DP* during the training process has a considerable impact on the generated explanations. Moreover, it was found that using the combination of *FedAVG* and *DP* can sometimes mitigate the negative effects of the added noise during the training process. This section covers an investigation whether the degree to which the quality of explanations is affected, differs for different attribution methods and datasets. Therefore, the relative increase in continuity score was measured when comparing the *Baseline* with the *DP* training setting. A higher relative increase indicates a bigger impact, resulting in a lower rank.

1) Attribution Methods

Figure 10 shows the ranks of different attribution methods when applied to different architectures before and after adding *DP* to the training, for the time series and image datasets. For both modalities, a prominent separation of two distinct groups can be noticed. In time series datasets, both *KernelSHAP* and *Occlusion* are affected significantly less by differential privacy as compared to the remaining, gradient-based methods. Similarly, *KernelSHAP* and *Occlusion* clearly outperformed most other methods. *DeepSHAP* is the only exception, which even achieved a similar score to *Occlusion* when applied to *NFNet*.



(a) Time series datasets



(b) Image datasets

FIGURE 10. Critical difference diagrams showing the impact of adding *Differential Privacy* during training, on the quality of explanations generated by different attribution methods.

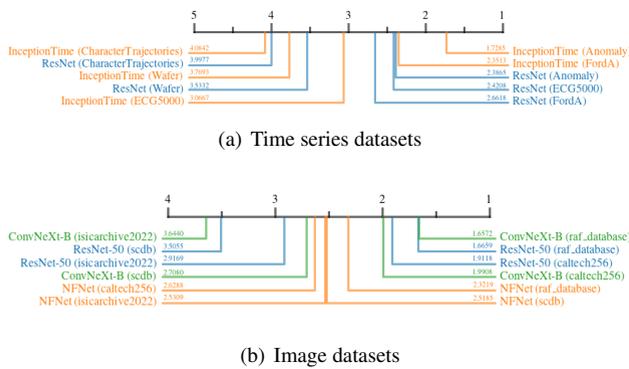
2) Datasets

Figure 11 shows the impact of *DP* on the quality of explanations for different datasets. For the time series domain, it can be seen that noise has the least impact on the *Anomaly* dataset, as the decision-relevant anomaly is not affected much by the added noise. On the other hand, *Character Trajectories* dataset is highly affected by noise. This can be explained by the fact that the dataset consists of raw sensor values that describe drawn letters. Slight noise distributed over the time series can have a devastating influence on the meaning of a given sample, as the error adds up over time. In the image-domain, *RAF-Database* and *Caltech-256* are influenced less by noise, whereas *ISIC*, on average, shows a higher susceptibility. This is understandable, as *ISIC* heavily relies on fine-grained patterns and complex features which might be more susceptible to added noise as compared to coarse-grained features used for emotion recognition and object detection. Surprisingly, the results suggest a rather high impact on *SCDB* as well. At first, this might seem unexpected due to the relevance of clean and uniform shapes for classification. However, considering the low resolution of input images, these shapes might be particularly fragile under the influence of noise, as small perturbations can easily shift the resemblance of one shape to another.

V. DISCUSSION

In the last few years, explainability and data privacy are drastically gaining importance in the field of Deep Learning. It is therefore all the more important to take a closer look at their interaction. The presented results revealed a significant impact of privacy-preserving training techniques on generated explanations. However, the influence on XAI strongly depends on the privacy technique used, as well as further factors.

First of all, it has been shown that not every PPML method has the same impact on model performance. *DP*-based models were shown to almost always deteriorating test



**FIGURE 11.** Critical difference diagrams showing the impact of adding *Differential Privacy* during training, on the quality of explanations when applied to different datasets.

accuracy. Moreover, experience showed that they drastically complicate model convergence and hyperparameter search. *FedAVG*, on the other hand, yielded accuracies similar to the *Baseline* setting, sometimes even improving the results. It has to be mentioned, though, that both *DP* and *FedAVG* follow different goals in the domain of privacy. Whereas *DP* aims at preventing models to capture individual sample information, which could be used for reconstruction, *FedAVG* mainly aims at minimising the exposure of sensitive information by keeping the training data local. Although *FedAVG* also generates an aggregated model which might have less vulnerability to reconstruction attacks due to averaging effects, it still needs to transfer information about the local models to the orchestration server. Therefore, the combination of *FedAVG* and *DP* provide the highest privacy, yielding in many cases similar performance compared to only *DP*.

The qualitative and quantitative analysis revealed various interesting findings regarding the impact of different privacy-preserving techniques on explanations. *Differential Privacy*, for example, stood out in almost all configurations for its property to add noise to the attribution maps. This has been reported in many individual samples and could be confirmed by dataset level analysis, as well as quantitative analysis, where *DP*-based methods stood out for increased *Continuity* values. One possible reason for this phenomenon is the addition of noise during the training process with *DP*. The introduction of noise in the parameter update most likely leads to contortions in the parameter space, which are never completely compensated, and translate into the prediction process. This effect might be counteracted by slightly tweaking the optimization, such as fine-tuning public datasets, or by increasing the batch sizes during training.

The degree to which noise is added has been investigated in Section IV-G. The results suggest, that perturbation-based methods are a lot less prone to changing their explanation's *Continuity* under influence of noise. Inspecting the individual examples, as well as the average heatmaps in Figure 4, this finding can again be verified. The difference between *Occlusion* and *Saliency* is particularly notable in the contrast be-

tween highly relevant areas and low relevant areas. Whereas *Saliency* produces monotonous heatmaps, peaks and areas of interest are much more prominently highlighted in the average *Occlusion* maps. The main reason for perturbation-based methods being less affected by noise in terms of *Continuity* is their higher resolution which neglects fine nuances in relevance, and the fact that randomly introduced noise is prone to cancel out within a patch. However, *Continuity* is only a mathematical approximation of an explanation's interpretability. Figures 1(b) and 2 illustrate that *Occlusion*-based explanations are often significantly changed when introducing *DP* during training. Furthermore, the high interpretability of heatmaps is worthless, if their fidelity is not ensured. As reported in Section IV-F, *DP* exclusively led to the deterioration of metrics indicating an explanation method's fidelity. Therefore, even when applying *Occlusion*, it needs to be clarified how truthful the generated explanations remain to be.

In contrast to *DP*, *Federated Learning* often resulted in smoother attribution. Interestingly, combining *FedAVG* with *DP* often times even led to more continuous attribution maps compared to the *Baseline* setting, reducing the negative effects introduced by *DP* alone. However, *FedAVG-DP* has also been reported to decrease the fidelity of explanations in many cases. Therefore, whenever XAI is required and *Differential Privacy* is applied, it might be worth considering a combination of *DP* and *Federated Learning*. This will also be possible in cases where *Federated Learning* is not required, as the federated setting can easily be simulated by dividing the dataset into chunks. Although some outlier experiments report a better *Continuity* score for *Baseline* settings, the fact that *FedAVG* leads to better *Continuity* scores has a strong theoretical basis. Averaging models during training inevitably prevents the final model from overemphasizing granular features or noise.

The present study also showed that the influence of PPML on XAI is not really dependent on the application domain, but rather on the choice and feature scales of the dataset at hand. The noise introduced by *DP* has, above all, a detrimental impact on classification tasks that rely on fine-grained and nuanced features or patterns. For simpler anomaly detection tasks or tasks focusing on the detection of overall, coherent structures seem to be less affected by privacy-preserving training techniques.

Besides the different influences PPML has on XAI, there is another fact that needs to be considered when combining both techniques. No matter how private a system has been made, exposing an explanation is in itself always a potential point of attack for a system, revealing sensitive information about the decision-making process. This is, for instance, particularly evident with *Saliency*, which provides the raw gradients of a single input instance. For truly critical applications one should ask the question of who, in the end, should be authorized to request explanations, and under which circumstances. Moreover, it might even be required to further obfuscate the exact generation process for explanations, or

rely exclusively on global explanations for applications with extremely high privacy requirements.

This study revealed several general trends which will affect explanations on a global scale when applying private training strategies to DL-based models. However, one major limitation of such studies is the examined basis of comparison. When comparing explanations of separate model instances, there is always the risk of obtaining different local minima, i.e., different classification strategies. Previous research [58] suggests that one dataset can have multiple, redundant, but fundamentally different features. Therefore, even models with identical test performance could have, in theory, picked up completely different cues to solve the same problem, hence yielding deviant explanations per model. When training models using different training strategies, it cannot be avoided to obtain models with deviating classification strategies. This is also clearly reflected in the naturally lower model performance of *DP*-based models.

Further limitations are related to the evaluation of the explanation's quality through quantitative metrics. As already mentioned, quantitative quality metrics for XAI are simply mathematical approximations of factors that could account for human interpretability or test assumptions of fidelity that should be satisfied by good explanations. Many such metrics still have inherent limitations like *AOPC*, *Sensitivity*, and *Fidelity*, introducing out-of-distribution samples through the perturbation of samples. *Ground Truth Concordance* assumes that, for each sample, there is exclusively one single decision path, and therefore a ground truth explanation. To approach this assumption, a synthetic dataset was utilized that allowed the construction of ground truth segmentation maps, highlighting all decision-relevant shapes. For somewhat complicated problems, explanations are always redundant as are the corresponding human problem definitions. The human-made logic behind the dataset postulates that a set of pre-defined shapes (i. e., star or triangle) need to be present to associate a sample with a class. However, instead of selecting the entirety of a shape, networks could also simply define triangles by the angle of their apices. This way, the explanation would not need to cover the complete shape, but only an arbitrarily small area around a single apex.

## VI. CONCLUSION

Both eXplainable AI and privacy-preserving machine learning constitute pivotal technologies for the safe translation of state-of-the-art AI algorithms into everyday applications. It is particularly important to get an early understanding of the effect of private training on XAI, to actively develop countermeasures, and avoid blind interpretations of explanations. This work showed that, although the exact effect on explanations depends on a multitude of factors including the privacy technique, dataset, model architecture, and XAI method, some overall trends can be identified. It has been found that *Differential Privacy*, on average, decreases both the *Interpretability* and *Fidelity* of heatmaps. However, *Federated Learning* was found to moderate both effects when used in

combination. When used, alone, *FedAVG* was even sometimes found to improve the *Interpretability* of attribution maps by generating more continuous heatmaps. The results suggest to consider *Federated Learning* before *Differential Privacy*, where appropriate. Moreover, it is recommended always to choose *Differential Private Federated Learning* as well as perturbation-based XAI methods, if an application requires both privacy and explainability. As the first work to investigate the impact of privacy on XAI, this study opens up a series of interesting follow-up questions, including the in-depth analysis of the trade-off between privacy and interpretability under different privacy constraints, and the impact of using privacy techniques beyond *DP* and *FedAVG*. Moreover, the field would benefit from a deeper analysis of the effect of privacy on the interpretation by human users through application-grounded and human-grounded evaluation methods. The union of PPML and XAI solves one of the remaining regulatory and safety-critical hurdles, freeing the way for innovative and high-performing AI-based applications that can bring significant advancements in crucial domains of our everyday lives.

## REFERENCES

- [1] X. Liu, L. Faes, A. U. Kale, S. K. Wagner, D. J. Fu, A. Bruynseels, T. Mahendiran, G. Moraes, M. Shamdas, C. Kern et al., "A comparison of deep learning performance against health-care professionals in detecting diseases from medical imaging: a systematic review and meta-analysis," *The lancet digital health*, vol. 1, no. 6, pp. e271–e297, 2019.
- [2] R. Sujatha, J. M. Chatterjee, N. Jhanjhi, and S. N. Brohi, "Performance of deep learning vs machine learning in plant leaf disease detection," *Microprocessors and Microsystems*, vol. 80, p. 103615, 2021.
- [3] D. J. Hemanth and V. V. Estrela, *Deep learning for image processing applications*. IOS Press, 2017, vol. 31.
- [4] F. Chen, H. Yu, R. Hu, and X. Zeng, "Deep learning shape priors for object segmentation," in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*, 2013, pp. 1870–1877.
- [5] A. Gilani, S. R. Qasim, I. Malik, and F. Shafait, "Table detection using deep learning," in *2017 14th IAPR international conference on document analysis and recognition (ICDAR)*, vol. 1. IEEE, 2017, pp. 771–776.
- [6] H. Ismail Fawaz, G. Forestier, J. Weber, L. Idoumghar, and P.-A. Muller, "Deep learning for time series classification: a review," *Data mining and knowledge discovery*, vol. 33, no. 4, pp. 917–963, 2019.
- [7] B. Lim and S. Zohren, "Time-series forecasting with deep learning: a survey," *Philosophical Transactions of the Royal Society A*, vol. 379, no. 2194, p. 20200209, 2021.
- [8] G. M. Binmakhshen and S. A. Mahmoud, "Document layout analysis: a comprehensive survey," *ACM Computing Surveys (CSUR)*, vol. 52, no. 6, pp. 1–36, 2019.
- [9] M. Fredrikson, S. Jha, and T. Ristenpart, "Model inversion attacks that exploit confidence information and basic countermeasures," in *Proceedings of the 22nd ACM SIGSAC conference on computer and communications security*, 2015, pp. 1322–1333.
- [10] C. Chen and N. D. Campbell, "Understanding training-data leakage from gradients in neural networks for image classification," *arXiv preprint arXiv:2111.10178*, 2021.
- [11] A. Das and P. Rad, "Opportunities and challenges in explainable artificial intelligence (xai): A survey," *arXiv preprint arXiv:2006.11371*, 2020.
- [12] I. E. Nielsen, D. Dera, G. Rasool, N. Bouaynaya, and R. P. Ramachandran, "Robust explainability: A tutorial on gradient-based attribution methods for deep neural networks," *arXiv preprint arXiv:2107.11400*, 2021.
- [13] O. Li, H. Liu, C. Chen, and C. Rudin, "Deep learning for case-based reasoning through prototypes: A neural network that explains its predictions," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 32, no. 1, 2018.
- [14] X. Liu, L. Xie, Y. Wang, J. Zou, J. Xiong, Z. Ying, and A. V. Vasilakos,

- “Privacy and security issues in deep learning: A survey,” *IEEE Access*, vol. 9, pp. 4566–4593, 2020.
- [15] M. Abadi, A. Chu, I. Goodfellow, H. B. McMahan, I. Mironov, K. Talwar, and L. Zhang, “Deep learning with differential privacy,” in *Proceedings of the 2016 ACM SIGSAC conference on computer and communications security*, 2016, pp. 308–318.
- [16] J. Konečný, H. B. McMahan, F. X. Yu, P. Richtárik, A. T. Suresh, and D. Bacon, “Federated learning: Strategies for improving communication efficiency,” *arXiv preprint arXiv:1610.05492*, 2016.
- [17] G. Vilone and L. Longo, “Explainable artificial intelligence: a systematic review,” *arXiv preprint arXiv:2006.00093*, 2020.
- [18] A. Boulemtafes, A. Derhab, and Y. Challal, “A review of privacy-preserving techniques for deep learning,” *Neurocomputing*, vol. 384, pp. 21–45, 2020.
- [19] F.-L. Fan, J. Xiong, M. Li, and G. Wang, “On interpretability of artificial neural networks: A survey,” *IEEE Transactions on Radiation and Plasma Medical Sciences*, vol. 5, no. 6, pp. 741–760, 2021.
- [20] K. Simonyan, A. Vedaldi, and A. Zisserman, “Deep inside convolutional networks: Visualising image classification models and saliency maps,” *arXiv preprint arXiv:1312.6034*, 2013.
- [21] Y. LeCun, Y. Bengio, and G. Hinton, “Deep learning,” *nature*, vol. 521, no. 7553, pp. 436–444, 2015.
- [22] A. Shrikumar, P. Greenside, A. Shcherbina, and A. Kundaje, “Not just a black box: Learning important features through propagating activation differences,” *arXiv preprint arXiv:1605.07173*, 2016.
- [23] J. T. Springenberg, A. Dosovitskiy, T. Brox, and M. Riedmiller, “Striving for simplicity: The all convolutional net,” *arXiv preprint arXiv:1412.6806*, 2014.
- [24] M. Sundararajan, A. Taly, and Q. Yan, “Axiomatic attribution for deep networks,” in *International conference on machine learning*. PMLR, 2017, pp. 3319–3328.
- [25] M. D. Zeiler and R. Fergus, “Visualizing and understanding convolutional networks (2013),” *arXiv preprint arXiv:1311.2901*, 2013.
- [26] M. T. Ribeiro, S. Singh, and 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*, 2016, pp. 1135–1144.
- [27] S. M. Lundberg and S.-I. Lee, “A unified approach to interpreting model predictions,” *Advances in neural information processing systems*, vol. 30, 2017.
- [28] C. Chen, O. Li, D. Tao, A. Barnett, C. Rudin, and J. K. Su, “This looks like that: deep learning for interpretable image recognition,” *Advances in neural information processing systems*, vol. 32, 2019.
- [29] B. Kim, M. Wattenberg, J. Gilmer, C. Cai, J. Wexler, F. Viegas et al., “Interpretability beyond feature attribution: Quantitative testing with concept activation vectors (tcav),” in *International conference on machine learning*. PMLR, 2018, pp. 2668–2677.
- [30] H. Hellani, R. Kilany, and M. Sokhn, “Towards internal privacy and flexible k-anonymity,” in *2015 International Conference on Applied Research in Computer Science and Engineering (ICAR)*. IEEE, 2015, pp. 1–2.
- [31] M. A. Rahman, T. Rahman, R. Laganière, N. Mohammed, and Y. Wang, “Membership inference attack against differentially private deep learning model,” *Trans. Data Priv.*, vol. 11, no. 1, pp. 61–79, 2018.
- [32] Y. Aono, T. Hayashi, L. Wang, S. Moriai et al., “Privacy-preserving deep learning via additively homomorphic encryption,” *IEEE Transactions on Information Forensics and Security*, vol. 13, no. 5, pp. 1333–1345, 2017.
- [33] M. A. Rahman, M. S. Hossain, A. J. Showail, N. A. Alrajeh, and M. F. Alhamid, “A secure, private, and explainable ioh framework to support sustainable health monitoring in a smart city,” *Sustainable Cities and Society*, vol. 72, p. 103083, 2021.
- [34] D. Franco, L. Oneto, N. Navarin, and D. Anguita, “Toward learning trustworthily from data combining privacy, fairness, and explainability: an application to face recognition,” *Entropy*, vol. 23, no. 8, p. 1047, 2021.
- [35] S. A. Siddiqui, D. Mercier, M. Munir, A. Dengel, and S. Ahmed, “Tsviz: Demystification of deep learning models for time-series analysis,” *IEEE Access*, vol. 7, pp. 67 027–67 040, 2019.
- [36] A. Bagnall, J. Lines, W. Vickers, and E. Keogh, “The uea & ucr time series classification repository,” 2021. [Online]. Available: [www.timeseriesclassification.com](http://www.timeseriesclassification.com)
- [37] S. Li and W. Deng, “Reliable crowdsourcing and deep locality-preserving learning for unconstrained facial expression recognition,” *IEEE Transactions on Image Processing*, vol. 28, no. 1, pp. 356–370, 2019.
- [38] G. Griffin, A. Holub, and P. Perona, “Caltech-256 object category dataset,” 2007.
- [39] B. Cassidy, C. Kendrick, A. Brodzicki, J. Jaworek-Korjakowska, and M. H. Yap, “Analysis of the isic image datasets: usage, benchmarks and recommendations,” *Medical Image Analysis*, vol. 75, p. 102305, 2022.
- [40] J. Kawahara, S. Daneshvar, G. Argenziano, and G. Hamarneh, “Seven-point checklist and skin lesion classification using multitask multimodal neural nets,” *IEEE Journal of Biomedical and Health Informatics*, vol. 23, no. 2, pp. 538–546, mar 2019.
- [41] A. Lucieri, M. N. Bajwa, A. Dengel, and S. Ahmed, “Explaining ai-based decision support systems using concept localization maps,” in *International Conference on Neural Information Processing*. Springer, 2020, pp. 185–193.
- [42] E. Ntoutsi, P. Fafalios, U. Gadiraju, V. Iosifidis, W. Nejdl, M.-E. Vidal, S. Ruggieri, F. Turini, S. Papadopoulos, E. Krasanakis, I. Kompatsiaris, K. Kinder-Kurlanda, C. Wagner, F. Karimi, M. Fernandez, H. Alani, B. Berendt, T. Kruegel, C. Heinze, K. Broelemann, G. Kasneci, T. Tiropanis, and S. Staab, “Bias in data-driven ai systems – an introductory survey,” 2020. [Online]. Available: <https://arxiv.org/abs/2001.09762>
- [43] A. W. Harley, A. Ufkes, and K. G. Derpanis, “Evaluation of deep convolutional nets for document image classification and retrieval,” in *International Conference on Document Analysis and Recognition (ICDAR)*.
- [44] J. Zhou, A. H. Gandomi, F. Chen, and A. Holzinger, “Evaluating the quality of machine learning explanations: A survey on methods and metrics,” *Electronics*, vol. 10, no. 5, p. 593, 2021.
- [45] T. Miller, “Explanation in artificial intelligence: Insights from the social sciences,” *Artificial intelligence*, vol. 267, pp. 1–38, 2019.
- [46] H. Ismail Fawaz, B. Lucas, G. Forestier, C. Pelletier, D. F. Schmidt, J. Weber, G. I. Webb, L. Idoumghar, P.-A. Muller, and F. Petitjean, “Inceptiontime: Finding alexnet for time series classification,” *Data Mining and Knowledge Discovery*, vol. 34, no. 6, pp. 1936–1962, 2020.
- [47] V. Feng, “An overview of resnet and its variants,” *Towards data science*, 2017.
- [48] A. Brock, S. De, S. L. Smith, and K. Simonyan, “High-performance large-scale image recognition without normalization,” in *International Conference on Machine Learning*. PMLR, 2021, pp. 1059–1071.
- [49] Z. Liu, H. Mao, C.-Y. Wu, C. Feichtenhofer, T. Darrell, and S. Xie, “A convnet for the 2020s,” in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 11 976–11 986.
- [50] A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” *Communications of the ACM*, vol. 60, no. 6, pp. 84–90, 2017.
- [51] K. Simonyan and A. Zisserman, “Very deep convolutional networks for large-scale image recognition,” *arXiv preprint arXiv:1409.1556*, 2014.
- [52] M. Tan and Q. Le, “Efficientnet: Rethinking model scaling for convolutional neural networks,” in *International conference on machine learning*. PMLR, 2019, pp. 6105–6114.
- [53] A. Shrikumar, P. Greenside, and A. Kundaje, “Learning important features through propagating activation differences,” in *International conference on machine learning*. PMLR, 2017, pp. 3145–3153.
- [54] C.-K. Yeh, C.-Y. Hsieh, A. Suggala, D. I. Inouye, and P. K. Ravikumar, “On the (in) fidelity and sensitivity of explanations,” *Advances in Neural Information Processing Systems*, vol. 32, 2019.
- [55] W. Samek, A. Binder, G. Montavon, S. Lapuschkin, and K.-R. Müller, “Evaluating the visualization of what a deep neural network has learned,” *IEEE transactions on neural networks and learning systems*, vol. 28, no. 11, pp. 2660–2673, 2016.
- [56] J. Demšar, “Statistical comparisons of classifiers over multiple data sets,” *The Journal of Machine Learning Research*, vol. 7, no. 1, pp. 1–30, 2006.
- [57] A. Benavoli, G. Corani, and F. Mangili, “Should we really use post-hoc tests based on mean-ranks?” *The Journal of Machine Learning Research*, vol. 17, no. 1, pp. 152–161, 2016.
- [58] A. Ilyas, S. Santurkar, D. Tsipras, L. Engstrom, B. Tran, and A. Madry, “Adversarial examples are not bugs, they are features,” *Advances in neural information processing systems*, vol. 32, 2019.



SAIFULLAH received the B.S. degree in mechanical engineering and the M.S. degree in robotics and intelligent machine engineering from the National University of Sciences and Technology (NUST), Pakistan. He is currently pursuing his Ph.D. at the University of Kaiserslautern and is working as a researcher at the German Research Center for Artificial Intelligence (DFKI GmbH) under the supervision of Prof. Dr. Prof. H. C. Andreas Dengel. His research interests include

document understanding and analysis, explainability and robustness of deep learning models, and privacy preservation in deep learning.



ANDREAS DENGEL is Scientific Director at DFKI GmbH in Kaiserslautern. In 1993, he became Professor in Computer Science at TUK where he holds the chair Knowledge-Based Systems. Since 2009 he is appointed Professor (Kyakuin) in the Department of Computer Science and Information Systems at Osaka Prefecture University. He received his Diploma in CS from TUK and his Ph.D. from the University of Stuttgart. He also worked at IBM, Siemens, and Xerox Parc.

Andreas is a member of several international advisory boards, has chaired major international conferences, and founded several successful start-up companies. He is a co-editor of international computer science journals and has written or edited 12 books. He is the author of more than 300 peer-reviewed scientific publications and supervised more than 170 Ph.D. and master theses. Andreas is an IAPR Fellow and received many prominent international awards. His main scientific emphasis is in the areas of Pattern Recognition, Document Understanding, Information Retrieval, Multimedia Mining, Semantic Technologies, and Social Media.



DOMINIQUE MERCIER received his Master degree in computer science from the Technische Universitaet Kaiserslautern, Germany in 2018. The topic of his Master's thesis was 'Towards Understanding Deep Networks for Time Series Analysis'. Currently, he is pursuing his Ph.D. at the German Research Center for Artificial Intelligence (DFKI GmbH) under the supervision of Prof. Dr. Prof. h.c. Andreas Dengel. His areas of interest include the interpretability of deep learning methods,

time series analysis, and document analysis. His work includes the development of novel interpretability methods for deep neural networks for time series analysis. Furthermore, he actively working in the NLP domain with a focus on citation and community management.



SHERAZ AHMED is Senior Researcher at DFKI GmbH in Kaiserslautern, where he is leading the area of Time Series Analysis. He received his MS and Ph.D. degrees in Computer Science from TUK, Germany under the supervision of Prof. Dr. Prof. h.c. Andreas Dengel and Prof. Dr. habil. Marcus Liwicki. His Ph.D. topic is Generic Methods for Information Segmentation in Document Images. Over the last few years, he has primarily worked on the development of various systems

for information segmentation in document images. His research interests include document understanding, generic segmentation framework for documents, gesture recognition, pattern recognition, data mining, anomaly detection, and natural language processing. He has more than 30 publications on the said and related topics including three journal papers and two book chapters. He is a frequent reviewer of various journals and conferences including Pattern Recognition Letters, Neural Computing and Applications, IJDAR, ICDAR, ICFHR, and DAS.



ADRIANO LUCIERI completed his BE in Mechatronic Engineering from Duale Hochschule Baden-Württemberg (DHBW) Mannheim and MS in Mechatronic Systems Engineering from Hochschule Pforzheim in Germany. He is presently pursuing a Ph.D. from Technische Universität Kaiserslautern (TUK), Germany, and is also working as Research Assistant at Deutsches Forschungszentrum für Künstliche Intelligenz GmbH (DFKI). His research focus lies on improving the explainability and transparency of Computer-Aided Diagnosis

(CAD) systems based on Deep Learning for medical image analysis. His work includes a concept-based explanation of skin lesion classifiers as well as the localization of concept regions in input images.

...