

Efficient XAI Techniques: A Taxonomic Survey

Yu-Neng Chuang, Guanchu Wang, Fan Yang, Zirui Liu,
Xuanning Cai, Mengnan Du, and Xia Hu

Abstract—Recently, there has been a growing demand for the deployment of Explainable Artificial Intelligence (XAI) algorithms in real-world applications. However, traditional XAI methods typically suffer from a high computational complexity problem, which discourages the deployment of real-time systems to meet the time-demanding requirements of real-world scenarios. Although many approaches have been proposed to improve the efficiency of XAI methods, a comprehensive understanding of the achievements and challenges is still needed. To this end, in this paper we provide a review of efficient XAI. Specifically, we categorize existing techniques of XAI acceleration into efficient non-amortized and efficient amortized methods. The efficient non-amortized methods focus on data-centric or model-centric acceleration upon each individual instance. In contrast, amortized methods focus on learning a unified distribution of model explanations, following the predictive, generative, or reinforcement frameworks, to rapidly derive multiple model explanations. We also analyze the limitations of an efficient XAI pipeline from the perspectives of the training phase, the deployment phase, and the use scenarios. Finally, we summarize the challenges of deploying XAI acceleration methods to real-world scenarios, overcoming the trade-off between faithfulness and efficiency, and the selection of different acceleration methods.

Index Terms—Interpretability, Efficient Explainable Artificial Intelligence, Feature Attribution Explanation, Counterfactual Explanation.



1 INTRODUCTION

MACHINE learning (ML) has been successfully applied in a variety of domains, such as recommender systems [1], [2], machine translation [3], and voice recognition [4]. Despite the advancements in ML, providing transparency in the models, particularly in deep neural networks (DNNs), remains a substantial challenge. The lack of transparency can lead to mistrust and skepticism of ML model predictions, such as the block-box driving decisions made by autopilots. Towards this end, explainable artificial intelligence (XAI) has received increasing attention from the community. A multitude of XAI algorithms have been introduced and can be categorized into two main groups: local and global explanations [5]. Local explanations aim to explain the reasoning behind an individual model prediction, while global explanations aim to uncover the overall functioning of a complex model by examining its structure and parameters. XAI techniques can also be classified into post-hoc and intrinsic explainability [6]. These XAI methods provide valuable insight into the decision-making process of machine learning models.

In this article, we focus on the efficiency problems of post-hoc local explanation since it is one of the most commonly used methods in the field of XAI. As shown in Figure 1, local explanation methods derive feature attribution scores by locally examining the model prediction of each individual data instance one at a time. Traditional post-hoc local explanation methods are particularly inefficient because of the gradual process of creating a unique explainer for each instance. For example, SHAP [7] takes around 27 seconds to generate a local model explanation for a

32×32 CIFAR-10 image. Considering the efficiency of deriving a model explanation, traditional post-hoc local ones encounter other severe inefficiency issues since each instance requires a unique explainer during the derivation of the explanation. In addition, the local explanation suffers from extensive computational conditions due to the pending amounts of tested instances, where each instance requires massive permutation times to complete the importance score estimation. In this manner, two drawbacks push the post-hoc local explanation methods into a difficult situation, obtaining exceptionally high time complexity to fulfill the local explanation requirements. Despite the high time complexity, the reliable performance still leads to a vast exploitation of post-hoc local methods. Post-hoc local explanation reveals faithful and effective instance-based explanation with theoretical guarantees, such as Shapley-based estimation [7] and counterfactual examples [8]. However, the high computational complexity creates a barrier to deployment in real-world systems with sampling variability [9]. Providing a real-time explanation is still a remaining challenge in balancing the trade-off between efficacy and efficiency for a post-hoc local explanation method.

Based on the challenges above, we analyze efficiency issues from three different explanation perspectives: individual feature explanations, feature tuple explanations, and influential example explanations. The first class of explanation perspectives focuses on calculating feature scores among all input features, assuming that each feature is mutually independent. Specifically, we discuss the first perspective in the context of feature attribution with acceleration mechanisms. The second perspective involves generating scores for statistical feature interactions. In real-world problem settings, the features are not always independent, which means that the interaction between multiple features may be significantly related to the underlying prediction results. Unlike the previous two attribution tasks, the last explanation perspective, influential examples, aims to provide an instance

- Yu-Neng Chuang, Guanchu Wang, Fan Yang, Zirui Liu and Xia Hu are with the Department of Computer Science, Rice University. E-mail: {ynchuang, guanchu.wang, fyang, zl105, xia.hu}@rice.edu.
- Xuanning Cai is with Meta Platforms, Inc. Email: {caixuanning}@fb.com.
- Mengnan Du is with the Department of Data Science, New Jersey Institute of Technology. Email: {mengnan.du}@njit.edu.
- Correspondence to Yu-Neng Chuang and Xia Hu.

that explicitly helps users understand what is happening inside the model. The generated explanation instance serves as a proxy for transparency toward the prediction models, increasing the trustworthiness of the underlying prediction outcomes. We primarily target counterfactual examples (algorithmic recourse) in this taxonomic review as a representative. For simplicity and clarity, we use the terms "counterfactual example" and "algorithmic recourse" interchangeably throughout the rest of this survey.

We further categorize the efficient approaches based on their acceleration actions, including non-amortized acceleration and amortized acceleration. These two acceleration methods greatly help users to trust the model prediction in a real-time scenario. The first line of work locally accelerates the interpretation process for each instance, trying to decrease the computational complexity before or during the explanation generation. However, non-amortized ones may encounter a limitation while keeping the underlying explanation accuracy intact since original feature information is diminished. The second line of work is amortized acceleration, which utilizes a DNN model to capture the explanation distribution among all training instances globally. The derivation time of amortized explainers can be largely reduced since it provides the explanation via a single feedforward process. However, it may sacrifice explanation performance, since the explainers learn the general explanation distribution instead of focusing on a single data instance. Non-amortized acceleration can preserve higher explanation performance while speeding up individual derivation processes, but it is still slower than amortized ones. Amortized acceleration methods can provide a robust and real-time explanation by exploiting one DNN explainer, but they may sacrifice explanation performance.

The rest of the article is organized into sections as outlined below. In Section 2, we first introduce the context and formulation of efficiency issues to post-hoc local explanation methods. We then summarize the current state-of-the-art in efficient explanation methods, including non-amortized acceleration in Section 3 and amortized acceleration in Section 4. Following that, in Sections 5 and 6, we present challenges and future directions for improving the efficiency of explanation derivation. Finally, in Section 7, we summarize the discussions in this work.

2 BACKGROUND OF EFFICIENT EXPLAINABLE ARTIFICIAL INTELLIGENCE

This section depicts the efficiency issues that explainable artificial intelligence (XAI) faces, as well as the requirements for accelerating strategies to address efficiency concerns in explanation tasks.

2.1 Efficient Issues on XAI

Post-hoc local explanation is the most widely studied method that aims to provide instance-based explanations to locally examine model predictions (see Figure 1) [6]. This paradigm has been widely used in real-world systems due to its effectiveness and faithfulness. Moreover, some previous work [10], [11] focus on eliminating the uncertainty of post-hoc local explanations to derive more stable, consistent, and reliable model explanations, which can increase user

TABLE 1

Experiment of feature attribution task on Adult Dataset(13 features). The experiments compare the deriving time with Ground-truth Shapley Values(GT-Shapley), Traditional XAI (SHAP [7]), and Efficient XAI (SHEAR & CoRTX). SHAP-X represents SHAP with X-times of permutation. "Second" (sec.) is the measurement unit of deployment time per instance.

	GT-Shapley	SHAP-4000	SHEAR [13]	CoRTX [14]
ℓ_2 -error	–	0.0021	0.0019	0.007
Time/instance	5.8536	0.9869	0.0141	0.00015

trust in the decisions made by prediction models. However, some traditionally local explanation techniques, such as Shapley-based methods [7], [12], are seriously suffering from efficiency issues, although their generated model explanations are faithfulness and theoretical guarantee. As the experimental results compared to GT-Shapley (ground truth Shapley value) shown in Table 1, traditional XAI requires nearly 1.0 seconds to obtain a satisfactory explanation result under 4,000 permutation times. In contrast, efficient XAI only needs 0.00015 seconds to complete the explanation derivation process, which is much faster than GT-Shapley and traditional XAI. In the following sections, we will discuss the challenges and limitations of post-hoc local explanations in the context of three different tasks: feature attribution, statistical interaction detection, and counterfactual examples.

- Feature Attribution Tasks:** In feature attribution tasks, the goal of efficient XAI is to accelerate the generation of feature importance scores for a model explanation. In this paper, we focus on Shapley value-based feature attribution tasks, which have been shown to be inefficient. Shapley values [15] originally aimed to estimate the feature contribution in the cooperative game theory. In feature attribution tasks, they are often used as the important scores of imputing feature set $\mathcal{U} = \{1, \dots, M\}$ to the black-box model behaviors. For any value function $f : 2^M \rightarrow \mathbb{R}$, the Shapley values $\phi_i(f, \mathcal{U}) \in \mathbb{R}$ of feature i can be formalized as follows:

$$\phi_i(f, \mathcal{U}) = \sum_{S \subseteq \mathcal{U} \setminus \{i\}} \binom{M-1}{|S|}^{-1} [f(\{i\} \cup S) - f(S)]. \quad (1)$$

where $S \subseteq \mathcal{U}$ is a feature coalition set. In other words, the average preceding difference considering all possible feature coalitions indicates the contribution of feature i . Nevertheless, the calculation of Shapley values relies on 2^M times of model evaluation to estimate the contribution of feature i , where the computational complexity is $T[\phi_i(f, \mathcal{U})] = O(2^M)$. One of the inefficient challenges is how to rapidly estimate the Shapley values but without traversing through all feature coalitions.

- Statistical Interaction Detection Tasks:** Unlike most feature attribution tasks, which assume that features are independent, this task assumes that different features interact with one another. Following the definition of statistical interaction, $I \subseteq \mathcal{U}$ is said to be a feature interaction if and only if there does not exist $f_{\setminus i}$ for $i \in I$

satisfying the following:

$$f(\mathcal{U}) = \sum_{i \in I} f_{\setminus i}(\mathcal{U} \setminus \{i\}) \quad (2)$$

where $f_{\setminus i}$ denotes a function that does not depend on the feature i . To analyze the interaction of features that contributed to the prediction results of the model, the measurement of interaction scores $\phi_I(f, \mathcal{U})$ can be formally defined as follows:

$$\phi_I(f, \mathcal{U}) = \sum_{L \subset I} (-1)^{|L|-1} f(L), \quad (3)$$

where L represents a n -variable sub-interaction from I . The high-level idea of the definition is to accurately consider all possible sub-interactions to form the actual importance score of feature interaction I . However, the computational process is extremely time-consuming and needs to go through all the possible sub-interactions and feature coalitions. Specifically, traditional interaction detection methods require calculating all statistical feature interactions to generate scores among feature interactions, where the number of interaction candidates is extremely large. The number of interaction candidates is 2^N if there exist N features, which makes it not applicable to real-world systems. The goal here is to alleviate the inefficiency of estimating scores for detecting feature interactions.

- **Counterfactual Example Tasks:** Another goal is to improve the efficiency challenges in counterfactual explanation tasks. Counterfactual explanations formalize the exploration of “what-if” scenarios, which are an instance of example-based reasoning using a set of hypothetical data samples. Considering a classification model $f : \mathbb{R}^M \rightarrow \{-1, 1\}$ as an example, well-established counterfactuals x^* are required to flip the prediction outcomes of the original queried instance q_0 . Given $f(q_0) = -1$, the explanation problem can be formally illustrated as follows:

$$x^* = \arg \min_{x \sim \mathcal{X}} \mathcal{L}(x, q_0) \quad \text{s.t.} \quad f(q_0) = -1, f(x^*) = 1 \quad (4)$$

where \mathcal{X} denotes the potential data distribution to the counterfactual universe of a given queried instance q_0 , and \mathcal{L} is the distance measurement in the input space. To solve Equation 4, traditional post-hoc local ones generate counterfactual examples by gradually adjusting the features to achieve good validity and sparsity [16]. Under good conditions of validity and sparsity, actionability and closeness are two more essential factors that lead the counterfactual examples to be user accessible. To meet these four requirements, additional operations such as mix-integer programming are necessary to generate counterfactual examples, which results in slow progress and hinders the ability to obtain efficient results.

In this paper, we concentrate on addressing the efficiency issues related to the three aforementioned explanation tasks. The goal of efficient XAI is to speed up the generating process while maintaining the effective model explanation and running it in a short period of time. To cope with the challenges, existing work proposes two kinds of framework for efficient XAI: non-amortized acceleration methods and amortized acceleration methods, which are shown in Figure 2.

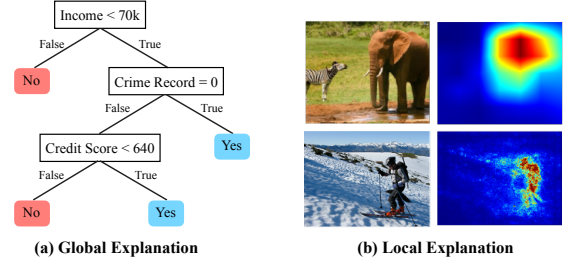


Fig. 1. Comparison between global (model-wise) explanation and local (instance-wise) explanation. (a) The decision tree intrinsically illustrates the model-wise important features of loan decisions from banks; (b) while heatmaps produced the instance-wise feature attributions toward each individual on the image classification task. In this survey, we focus on the efficient issues in local explanation.

2.2 Non-amortized Acceleration Methodology

Efficient non-amortized methodologies aim to address the efficiency issues of post-hoc local explanation methods on feature attribution, interaction detection, and counterfactual explanation tasks. Non-amortized post-hoc local methodologies require a unique explainer for each instance in order to derive a model explanation, which makes the time complexity grow with the number of tested instances. Despite the high latency in deriving explanations, non-amortized methods are effective at providing satisfactory model explanations because each explainer is only trained on a specific tested instance. Existing approaches can be divided into two groups based on their explanation angles: attribution-oriented and counterfactual-oriented methods.

2.2.1 Non-amortized Attribution-oriented Methods

Non-amortized attribution-oriented methods focus on generating instance-wise model explanations, which provides the model explanation of underlying model prediction among each tested feature or feature interaction. According to the training paradigms, attribution-oriented methods can be divided into three aspects to accelerate inefficient challenges in Equation (1) and Equation (3). The first line of methods [17], [18] utilizes proxy models (e.g., linear regression) to fit the distribution of feature importance scores and treats the coefficient of proxy models as estimated Shapley values. The second group of work adopts the preceding difference of the value function to generate the explanation, which averagely calculates the difference in model prediction as the feature importance scores. For example, RISE [19], Permutation Sampling [20], and Antithetical Permutation Sampling [21] are three representative works. The third group of methods exploits the model gradient to yield a model explanation, such as Integrated Gradient [22], GradCAM [23] and SmoothGrad [24]. Generally, the third group obtains a relatively faster explanation process than the first and second groups. However, extra time is still required to perform the sampling process on each instance while generating the model explanation. The time complexity of these three groups of methods is still highly dependent on the number of sampling and tested instances. Despite their effectiveness and efficacy, attribution-based explanation methods often face the challenge of long computation times, which presents a significant obstacle in deploying non-amortized methods to real-time systems.

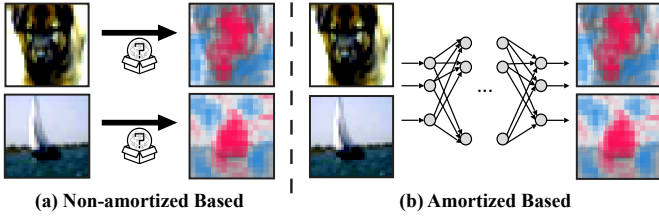


Fig. 2. A non-amortized explanation method using four explainers for four instances, and an amortized explanation method exploiting only one explainer among four instances.

2.2.2 Non-amortized Counterfactual-oriented Methods

Generally, Equation (4) can be merely solved by algorithm-based counterfactual approaches [25], [26], [27], which employ different optimization strategies to generate counterfactual examples for each queried instance q_0 . However, this approach suffers from efficiency issues since the counterfactual instances are decided instance by instance. The optimization process of each new query requires one specific optimization problem of Equation (4) at one time. The explanation progress among a group of queried instances could be exceptionally time consuming. In this case, existing work provides data adjustment on imputing instances before optimizing the counterfactual distribution, accelerating the convergence of the training process. Another pack of prior achievements incorporates counterfactual-oriented constraints to speed up the training of yielding counterfactual examples. Despite the efficacy of algorithm-based ones, the high computational cost eventually creates a barrier to deployment on real-world systems that require low latency to provide real-time service.

2.3 Amortized Acceleration Methodology

Unlike non-amortized methods, amortized methodologies utilize a unified explainer to learn the distribution of model explanation. The goal of amortized ones is to replace the local training pattern in post-hoc local explanation with a deep neural network (DNN). The alternative replacement accelerates the process of generating an instance-based model explanation. The unified explainer only requires conducting a single feed-forward process in the inference phase. In this way, the time complexity of amortized methods is restricted to constant, which explicitly solves the efficient issues discussed in Section 2.1. Existing methods are designed to deal with two discussed explanation tasks: the feature attribution task and the counterfactual example task.

2.3.1 Amortized Attribution-oriented Methods

Amortized attribution-oriented methods employ predictive DNN models to simultaneously approximate the unified distribution of feature attributions for each instance. The approximation of data distribution has been widely exploited in various domains, such as matrix factorization [28], [29] in recommendation systems, to learn the distribution of task-specific information. The deductive output of attribution-oriented models can help humans understand what prediction models consider when making predictions. Amortized methods fundamentally address the efficiency issues of non-amortized methods, which reduce the number of explainers

during local explanation derivation. However, finding a balance between explanation effectiveness and efficiency poses a challenge in establishing an optimal training paradigm.

2.3.2 Amortized Counterfactual-oriented Methods

Amortized counterfactual-oriented methods use generative-based models or reinforcement agents to learn the general rules of counterfactual examples by solving Equation (4). Different from attribution-oriented ways, which are plausible to rely on ground-truth explanation information, counterfactual cases usually do not obtain ground truth as the label reference when updating the explanation models. The example generation process in the inference phase thereby requires only a single forward passing, making it a real-time explanation derivation.

- **Generative-based Methods:** To cope with the challenges mentioned above, counterfactual-oriented methods utilize generative models to synthesize general rules for producing counterfactual examples in the latent code space. Specifically, the generative-based model builds an adversarial learning framework, including a generator to produce counterfactual examples and a discriminator to prevent deviated examples from the generator. The framework typically employs a model based on generative adversarial networks (GANs) [30], [31] and incorporates counterfactual-oriented objectives in the training phase. In this way, generative-based frameworks can accelerate the derivation process by producing multiple counterfactual examples for a given instance simultaneously.
- **Reinforcement-based Methods:** Counterfactual-oriented methods take advantage of deep reinforcement agents to efficiently formulate the decision policy to generate counterfactual examples. The decision process for counterfactual examples is typically a discrete action, which is not differentiable and cannot be directly used in model training. Instead, deep reinforcement learning (Deep RL) [32], [33], [34] employs an action space module to transform discrete actions into corresponding continuous parameters. By taking advantage of the properties of Deep RL, reinforcement agents in counterfactual-oriented methods can synthesize the discrete rules for generating counterfactual examples. The Deep RL framework for model explanation can then accelerate the explanation-generation process with the aid of pre-trained agents during the inference phase.

3 NON-AMORTIZED ACCELERATION

In this section, we divide the non-amortized acceleration into two broad families: data-centric acceleration and model-centric acceleration.

3.1 Data-centric Acceleration

The data-centric mechanism adjusts the data to fundamentally improve the performance of the ML model. Superior data quality usually leads to fast convergence of the learning process and competitive prediction results since it eliminates the noise from the original data. The goal of data-centric acceleration is to dismiss the unimportant features or feature coalitions before the head, which reduces the computational complexity of explanation derivation. In Section 3.1.1, we

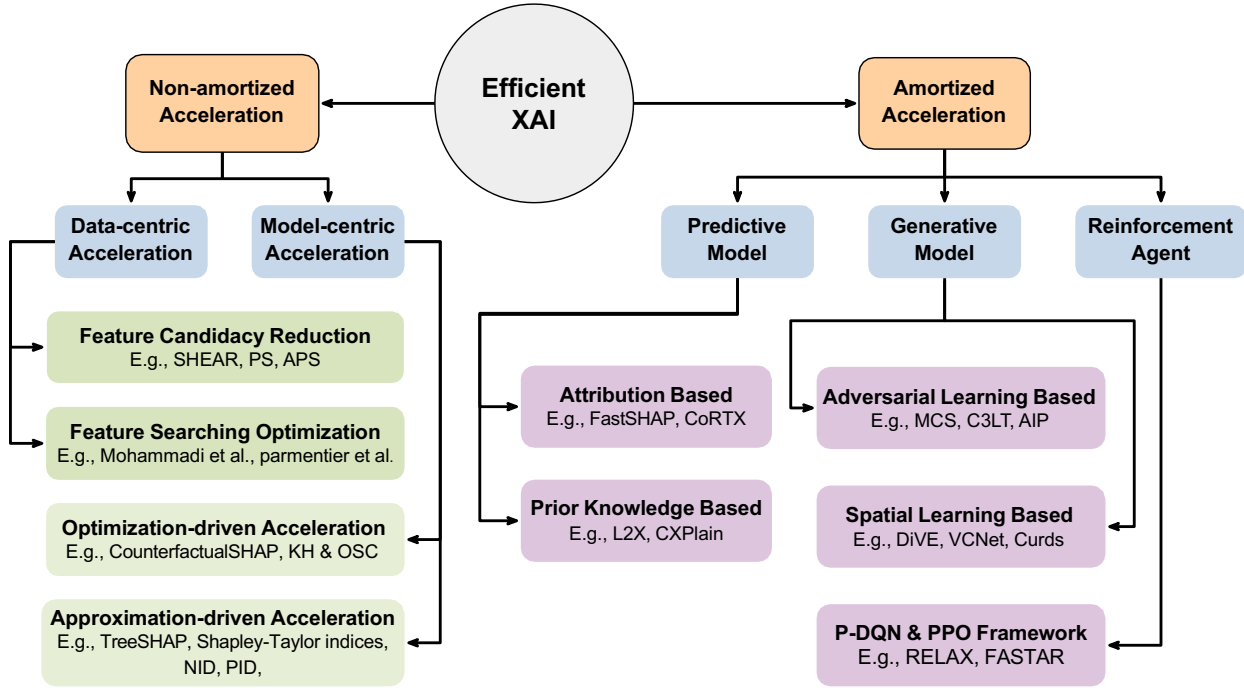


Fig. 3. Structure of the survey. Two main categorizations are discussed in this paper: Non-amortized Acceleration and Amortized Acceleration.

summarize acceleration methods for feature attribution, and in Section 3.1.2, we summarize acceleration methods for counterfactual example generation.

3.1.1 Feature Candidacy Reduction

Feature reduction accelerates the estimation of model explanation by eliminating the feature candidates, making the precise computation feasible and practical. The chosen features are typically decided by the causative observations from the prediction models and feature properties, such as pre-calculating contributions of feature that are independent with prediction model [35].

One of the existing work SHEAR [13] accelerates the estimation of Shapley value by contributive cooperative feature selection. The cooperative features are selected based on the well-contributed feature interactions from the prediction model. The Shapley values are estimated by only considering the selected coalitions of contributive features. Different from SHEAR, another work groupShapley [36] estimates the Shapley values under feature groups instead of individual features. The feature groups are decided in terms of feature type or dependence. Calculating the Shapley values of given features typically requires iterating all the features in combination with the given features, which inherently increases the execution time. One previous work [37] also proposes two measures for instance-wise feature importance scoring: L-Shapley and C-Shapley on the graph-structured dataset, which selects the feature candidacy by neighborhood information of the given graph structural dataset. The two measurements compute the Shapley values by exploiting the underlying graph structure, which reduces the time complexity from exponential to linear level. To cope with the challenge, the target features are replaced with a group of features on deriving the Shapley values, where the time complexity can be decreased once the group number is small.

Considering the intractable situation of lacking prior information, some previous works [38], [39], [40] focus on reducing the number of features by randomly sampling partial candidate coalitions before obtaining feature attribution scores. The deduction on coalitions usage can directly speed up the calculation process. One of the representative works is permutation sampling [20], which is an unbiased sampling method to speed up the estimation of the Shapley values. As the random sampling is still a bottleneck of efficiency, the antithetic sampling [41] is combined with the permutation sampling to alleviate the efficient drawback from random sampling. Specifically, the antithetic sampling method takes the correlated pairs sampling instead of standard i.i.d. random instance sampling. For simplicity in practice, the correlated pairs are considered the complementary pairs of the feature sets. This enables the pairs of permutations to be negatively correlated, theoretically reducing the estimation variance. Moreover, antithetic sampling reduces half of the sampling time since it only needs to sample one element in each complementary pair. As an efficient and effective method, antithetic permutation sampling has been widely deployed in real-world systems.

3.1.2 Feature Searching Optimization

According to counterfactual optimization problems formulated in Equation (4), the derivation process of counterfactual examples can be formulated as the contrasted optimization problems. Generally, traditional approaches such as integer programming are able to solve the problems and get the optimal solutions by adjusting the values of features. Nevertheless, the conventional optimization approaches are typically time-consuming due to their complex search steps. Some existing work even uses some kind of random or approximate search [65], [66], which are also computationally slow, particularly with high-dimensional instances. The

TABLE 2

Assessment of the collected efficient XAI papers on the key properties. "Non-amtz." and "Amtz." represent non-amortized acceleration and amortized acceleration, respectively. Details about the methods in the full table are given in Sections 3 and 4. "NLP" represents natural language data, "CV" denotes image data, and "Tab." means tabular data.

Paper	Development Process			Derivation Process		
	Modeling Paradigm	Acceleration Methods	Datatype	Shapley Based	Explanation Perspectives	
Non-amtz.	[13]	Coalition	Sample by Prior Information	Tab.	Yes	Feature Attribution
	[36]	Feature Cluster	Sample by Prior Information	Tab.	Yes	Feature Attribution
	[37]	Graph Neighborhood	Sampling by Prior Information	CV, NLP	Yes	Feature Attribution
	[20]	Coalition	Random Sampling	Tab.	Yes	Feature Attribution
	[21]	Coalition	Antithetic Sampling	CV, Tab.	Yes	Feature Attribution
	[21]	Coalition	KH & OSC	CV, Tab.	Yes	Feature Attribution
	[42]	MIP	Constraint Search	NLP	No	Counterfactuals
	[43]	MIP	Nearest Search	Tab.	No	Counterfactuals
	[44]	MIP	Constraint Search	Tab.	No	Counterfactuals
	[45]	QP	Constraint Search	CV, Tab.	No	Counterfactuals
	[46]	MIP	Task-oriented Constraint	Tab.	No	Counterfactuals
	[47]	Bayesian Opt.	Task-oriented Constraint	Tab.	Yes/No	Counterfactuals
	[48]	Spatial Opt.	Task-oriented Constraint	CV, Tab.	No	Counterfactuals
	[49]	Spatial Opt.	Shapley-based Constraint	Tab.	Yes	Counterfactuals
	[18]	Linear Model	Variance Reduction	Tab.	Yes	Feature Attribution
	[9]	Tree Model	Marginal Path Approximation	Tab.	Yes	Feature Attribution
	[50]	Tree Model	Banzhaf Values	Tab.	No	Feature Attribution
	[51]	MLP Model	Examine from Model Weights	Tab.	No	Feature Interaction
	[52]	MLP Model	Examine from Model Weights	CV, Tab.	No	Feature Interaction
	[12]	Multilinear Extension	Weighted Least Squares Solver	NLP, Tab.	Yes	Feature Interaction
[53]	Linear Model	Weighted Least Squares Solver	NLP	Yes	Feature Interaction	
Amtz.	[14]	Self-supervised Learning	Contrastive Framework	CV, Tab.	Yes	Feature Attribution
	[54]	Supervised Learning	Sampling-generated Labels	CV, Tab.	Yes	Feature Attribution
	[55]	Supervised Learning	Unified Feature Selection	CV, Tab.	No	Feature Attribution
	[56]	Supervised Learning	Unified Causality Relationship	CV, Tab., NLP	No	Feature Attribution
	[57]	Adversarial Learning	GAN	CV, NLP	No	Counterfactuals
	[58]	Adversarial Learning	CGAN	Tab.	No	Counterfactuals
	[59]	Adversarial Learning	GAN	CV	No	Counterfactuals
	[60]	Spatial Learning	Conditional Subspace VAE	Tab.	No	Counterfactuals
	[61]	Spatial Learning	VAE	CV	No	Counterfactuals
	[62]	Spatial Learning	VAE	CV, Tab.	No	Counterfactuals
	[63]	Reinforcement Learning	P-DQN Framework	Tab.	No	Counterfactuals
	[64]	Reinforcement Learning	PPO Algorithm	Tab.	No	Counterfactuals

goal of feature searching optimization is to provide an efficient searching algorithm with linear or convex quadratic programming (QP) solvers to speed up the derivation of counterfactuals by changing the feature values.

Prior works [43], [44], [45], [42] focus on seeking efficient programming algorithms as the simulation of distance measurement function $\mathcal{L}(\cdot)$ in Equation (4), providing the validity and proximity guarantees on counterfactual examples. One of the existing works [42] proposes a novel set of constraints, named mixed polytope, to incorporate with an integer programming solver to efficiently find coherent counterfactual explanations. The proposed mixed polytope can guarantee that the derived counterfactuals obtain a set of desiderata on diversity and closeness. Another work [43] proposes efficient approaches to search for the nearest counterfactual explanation within a given interval in the input feature space, where the proposed searching ones can formulate the optimization problem with mixed-integer programming (MIP) solvers. The experimental results also prove a significant improvement in runtime efficiency in yielding MIP-based approaches for counterfactual generation. Regarding the application of feature searching optimization, a work [67] proposes several searching strategies that aim to efficiently extract concise explanations under constraints without access to the internal recommendation model. The

key idea is to provide counterfactual explanations that are actionable and sparse, defined as minor changes to the user's interaction logs for explaining recommendation outputs.

3.2 Model-centric Acceleration

Model design is one of the essential factors that directly impact the performance of ML models, including speed and performance. Providing fitted task-oriented objectives and well-designed model structures can accelerate the training progress of ML models. In this manner, model-centric acceleration boosts the derivation speed of model explanation based on the model-wise adjustments, such as optimizing the objective functions or proposing new approximation models, which better estimate the explanation of the underlying prediction models with a faster speed.

3.2.1 Optimization-driven Acceleration

Optimization-driven methods accelerate the explanation-deriving process by replacing it with relatively optimal training strategies. This line of work follows the philosophy that task-oriented training can solve efficient issues from the learning perspective, allowing it to converge faster. In this direction, we consider feature attribution tasks and counterfactual example tasks. Existing achievements focus

on addressing the efficiency issues of counterfactual and feature attribution tasks.

- Feature Attribution Tasks:** Besides optimizing the training paradigm with the same objectives for faster derivation, some previous works reformulate the training objectives to accelerate the convergence speed of generating stable model explanations. Existing work [21] proposes two simple and effective methods to accelerate the convergence of permutation sampling. The first method is Kernel Herding (KH) which adopts a dynamic process to select the buffer of permutations. In each update step, a permutation is selected to minimize the Kendall correlation coefficient (KCC) [68] with the permutations in the buffer and maximize the KCC with the proposed efficient components. In this manner, one more permutation is added to the buffer until the selected permutations reach the number limitation. The second one is Orthogonal Spherical Codes (OSC), which reformulates the selection of permutations into orthogonal transformation in the hyper-sphere. To be concrete, OSC selects the correlated samples on the hyper-sphere from a basis of orthogonal vectors. OSC can significantly accelerate the convergence since the KCC of the selected permutations can be explicitly bounded in the range of $[-\frac{1}{2}, +\frac{1}{2}]$, which is better than random sampling with $[-1, +1]$. However, the computational overhead of OSC is $O(n^2)$, where n is the number of features.
- Counterfactual Example Tasks:** As for counterfactual tasks, various of prior works [46], [69], [47], [70] proposes new loss functions to extract actions towards an interpretable counterfactual. [48] offers a fast and model-agnostic method by incorporating the class prototypes in the objective function to guide the perturbations quickly to generate counterfactual explanations. The prototypes are beneficial to remove the computational bottlenecks caused by black-box prediction models. Another work, CounterfactualSHAP [49], exploits "background information" with SHAP [7] to guide on using feature attributions scores for rapidly generating counterfactual instances. The proposed framework comprises two new loss functions, the cost function and the action function, to enrich the counterfactual-ability of an approximated feature attribution from SHAP.

3.2.2 Approximation-driven Acceleration

In this subsection, we introduce the approximation-driven methods according to feature attribution tasks and interaction detection tasks, respectively.

- Feature Attribution Tasks:** Approximation-driven methods propose white-box models to estimate the feature importance scores according to the ML predictions around the neighborhoods of the given inputs. Generally, the white-box models are self-interpretable and less complex, ensuring the explanation deviating process is fast and explainable. The white-box model does not have to work well globally, but it has to approximate the black-box model well in a small sampling neighborhood near the original input. Then the contribution score for each feature can be obtained by examining the parameters of the white-box model. One existing work [18] aims to rapidly estimate the unbiased feature importance scoring by approximating the

optimal coefficients over the weighted linear regression. Considering the prohibitively large ensembling times to achieve an unbiased explanation performance, this work proposes a variance reduction technique to speed up the approximation process of acquiring optimal coefficients as estimated unbiased Shapley values. Unlike linear-based model approximation, some studies claim that the distribution of sampling neighborhood near the original input may be extremely non-linear, and linear-based model could lead to biased and ineffective estimations for model explanation [6]. In this case, another group of works [50], [9] takes advantage of non-linear models, such as tree-based models, to reduce the time complexity of estimating Shapley values. For instance, TreeSHAP [9] decreases the time complexity from $O(TL2^M)$ to $O(TLD^2)$ by reducing the number of trees (T) and maximum depth (D) of trees with L leaves.

- Feature Interaction Tasks:** Besides the single feature attribution task that assumes each feature to be independent, some existing works focus on analyzing the attributions of feature interactions. However, the interaction candidacy is extremely huge. The number of interaction candidates is 2^N if there exist N features. In the problem formulation in Equation 3, traditional statistical post-hoc tests, such as ANOVA with Fisher's Least Significant Difference (LSD), need to conduct comprehensive tests through all interaction candidates, which is impractical for applying to large-scale problems. In this case, the goal of neural network-based interaction detection methods is to efficiently analyze the statistical interactions of Equation 3, which are highly related to the results of prediction models. One of the existing works, NIT [71], first demonstrates that any interacting features must follow strongly weighted connections to a common hidden unit before reaching the final output layer. Based on this observation, NID [51] detects interactions from weights of learned neural networks by examining the interacting paths between the weight of the first layer and those from subsequent layers. PID [52] further explored the idea of interacting paths by extending the theory of persistent homology to the interaction detection problem. Although these methods claim to generate a post-hoc global-wise explanation, it is still available to generate locally instance-wise explanation [52] with interaction importance scores. Regarding the application of interaction detection methods, GLIDER [72] utilizes gradient interaction detection methods to efficiently build synthetic crossing features for each detected group of interacting features to improve the recommendation performance.

Another aspect is the acceleration of Shapley-based interaction scores. Following the definition of Shapley interaction index [73], the Shapley interaction scores $\phi_I(f, \mathcal{U})$ of targeted feature interaction $I \subseteq \mathcal{U}$ can be defined as:

$$\phi_I(f, \mathcal{U}) = \sum_{S \subseteq \mathcal{U} \setminus I} \binom{M-1}{|S| - |I| + 1}^{-1} (-1)^{|I| - |L|} \sum_{L \subseteq I} f(L \cup S),$$

where L represents the potential of feature sub-interactions happening in targeted feature interaction I . However, it is extremely time-consuming to calculate the exact Shapley

interaction scores. The time complexity $O(2^{M+|I|})$ is even higher than the exact calculation of Shapley values for individual features, which is implausible to apply to any real-world system. In this manner, some prior works focus on approximating the interaction scores on top of Shapley interaction scores. One of the works, Faith-SHAP [53], aims to extend Shapley values from individual feature importance scores to interaction importance scores. The key idea is to interpret Shapley values as coefficients of the most faithful linear approximation to the pseudo-Boolean coalition game value function, where the process of leveraging the weighted linear approximation can lead to efficient explanation derivation of feature interaction. Another work, Shapley-Taylor indices [12], detects feature interactions by expanding the Taylor series of the multi-linear extension with the set-theoretic model behavior. Under the samples of interaction permutations by random process over orderings of features, the derivation process can be accelerated through the approximation. By exploiting the concept of second-order Shapley-Taylor indices, [74] extend traditional non-amortized methods, such as GradCAM, LIME, and SHAP, to extract pairwise correspondences between images from trained opaque-box models. The other work [17] detects feature interaction based on Taylor expansion by interpreting the importance of interactions as mixed partial derivatives, which achieve a runtime-efficiency process.

4 AMORTIZED ACCELERATION

After understanding the non-amortized ways to accelerate post-hoc local explanation, we turn to introduce how amortized acceleration affects the efficiency of deriving model explanation. Unlike non-amortized ones, amortized methods require only one explanation model to generate the model explanation among all instances, which speeds up the progress with its fast inference phase of a unified explainer. In this section, we introduce three different mechanisms of amortized acceleration methods on both feature attribution and counterfactual tasks.

4.1 Predictive-driven Method

The predictive-driven methods maintain a unified DNN explainer to generate a fast model explanation among each data instance. The explanation can be generated via a single feed-forward process of the DNN explainer, providing real-time estimation on feature attribution tasks. Compared to the existing traditional XAI methods, the advantages of the predictive model mainly lie in two folds: (1) faster explanation generation; and (2) more robust explanation derivation. Generally, existing works attempt to learn the overall explanation distribution using two lines of methodologies, which are attribution-based approaches [75], [54], [76], [76], [14] and prior-knowledge-based approaches [55], [77], [78], [56], [79], [80], [81], [82], [83].

The first line of work employs the DNN explainers to simulate a given approximated Shapley distribution for generating explanation results. One of the representative works, FastSHAP [54], exploits a DNN model to capture the Shapley distribution among training instances for efficient

real-time explanations, which is a supervised paradigm learned with approximated Shapley value labels. Although the approximate attribution labels are rapid, degradation has been shown to affect the performance of explanation broadly [14]. Another work, CoRTX [14], proposes an unsupervised learning paradigm, which can significantly reduce the dependency on the Shapley labels and accelerate the derivation progress. The unsupervised CoRTX benefits the explanation tasks by exploiting a contrastive framework for generating latent explanations. After that, CoRTX fine-tunes the latent explanations with extremely few-shot labels to get the final model explanation.

The second line of work assumes the specific pre-defined causal relationship or feature distributions, as well as formulates the explainer learning process based on the given assumptions. One of the works, L2X [55], provides a feature masking generator for real-time feature selection. The training process of the mask generator is under the constraint from the given predicted label distribution of masked imputing instances. In addition to the previous definition of data distribution, CXPlain [56] train an explanation model by using a causal objective function that follows the definition of Granger causality [84]. To guarantee efficacy, CXPlain provides the uncertainty estimations for feature importance scores that are strongly correlated with the efficacy of the provided importance scores on previously unseen test data.

In this manner, both the first line and second line of work exploit DNN models to learn the unified distribution of model explanation for providing the explanation derivation with pre-trained explainers. Due to the inference phase of well-trained explainers, predictive models can provide real-time model explanations.

4.2 Generative-driven Method

The goal of generative methods is to learn the unified counterfactual rules by utilizing generative models from raw data instances such as textual and image data. The derivation process can be effective and efficient by exploiting generative models conditioned with a counterfactual universe. Unlike traditional local-wise counterfactual methods that modify instances in the data space, generative models focus on constructing the attribute-informed latent space, guiding the generative models to fit the counterfactual distribution of multiple query data instances. Existing work proposes their DNN frameworks based on two training paradigms: adversarial-learning-based [58], [57], [85], [86], [62], [59] and spatial-learning-based [87], [88], [61], [89], [60], [90], [91], [87], [92], which allows counterfactual examples to be rapidly generated by reusing same models among multiple queried data instances.

The first aspect is adversarial-learning-based ones, which utilize adversarial learning to guarantee the effectiveness of derived counterfactuals. A representative adversarial-based work, MCS [58], constructs a model-based synthesizer by using a conditional generative adversarial network (CGAN) [31] to capture counterfactual information faithfully. By incorporating model inductive bias, MCS can accurately exploit the causal dependence of attributes, which attempts to ensure the design correctness of the generative models through the causality identification process. The second

aspect is spatial-learning-based approaches, which adopt generative models to build up the algorithmic recourse frameworks by reformulating the counterfactual latent space. One of the existing works, DiVE [61], is built upon the VAE [93]-based structure with counterfactual constraints, allowing the deviation process to be efficiently completed via single forward passing. This work aims to learn the disentangled latent space by leveraging the Fisher information matrix of the underlying prediction ML models. With the learned representation, the generated counterfactual explanation is enforced to be proximity, sparsity, and diversity. However, compared to non-amortized counterfactual methods, there is a trade-off between explanation efficacy and efficiency in generative-driven strategies, making it a remaining challenge to balance these two goals in the training phase.

4.3 Reinforcement Agent

Reinforcement agents aim to reformulate the counterfactual explanation problem into sequential decision-making progress, where they utilize reinforcement learning to find the optimal counterfactual instances. Owing to the fast inference phase of reinforcement agents, the explanation deviating process is extremely fast and can generate multiple instances simultaneously. Prior work, RELAX [63], has produced model-agnostic counterfactual examples using the P-DQN framework. The derivation process is reformulated as a Markov decision process (MDP) with hybrid discrete-continuous actions to ensure the sparsity and proximity of the generated counterfactuals. Experimental results have shown that these agents lead to significant improvements in both efficacy and efficiency. Another work [64] propose a stochastic-control-based approach that uses the proximal policy optimization (PPO) algorithm to generate sequential algorithmic recourse. The derivation process allows the data instance to move stochastically and sequentially across the intermediate states to reach the final generated examples, which are treated as the generated counterfactuals. Following the practical guidance of algorithmic recourse, the framework translates the algorithmic recourse problem into an MDP with the constraints of discrete state space and discrete action space. In particular, the design of state space ensures actionability and action space fulfills the sparsity. The evaluation results indicate the successful generation of sequential counterfactual instances that meet other recourse desiderata.

5 DISCUSSION OF THE LIMITATIONS

We briefly introduce the limitations of efficient XAI we have surveyed from the training phase, deployment phase, and using scenarios. Then we present suggestions for choosing efficient XAI methods that might be more suitable to match the using scenarios for users.

5.1 Limitation on Training Phase

The learning strategies of non-amortized and amortized methods have been a bottleneck in making XAI local methods efficient. Current efficient non-amortized explanations are usually given in reducing queried features and providing advanced efficient-oriented algorithms. Without the training-testing process, non-amortized methods can neglect to

prepare the pre-trained explainers. In contrast, amortized methods need to train a unified model before deriving a model explanation, causing additional computational resources to an explainer model in advance. In addition, the training speed of amortized ones highly depends on the explanation data scale, which can be referred to as the usage amount of explanation label in the training phase of supervised learning paradigm [54], [57], and the fine-tuning stage of self-supervised learning paradigm [14]. This means that the amortized method moves the heavy computational complexity to the training phase, leading to a fast derivation via its single-forward inference process. There exists a trade-off between training speed and explanation performance in amortized methods on feature attribution and counterfactual tasks. In general, the trade-off between training speed and explanation performance does not have a significant impact in real-time online service, as amortized explanation models are typically pre-trained. However, if a production system requires online learning [95], the long training time required for amortized methods could be a barrier to deployment.

5.2 Limitation on Deployment Phase

There are several metrics that can be used to evaluate the efficiency of non-amortized and amortized acceleration methods in the deployment phase. One of the commonly used metrics is algorithmic throughput [13], [14], [96]. Specifically, the throughput is calculated by $\frac{N_{\text{test}}}{t_{\text{total}}}$, where N_{test} and t_{total} denote the testing instance number and the overall time consumption of explanation derivation, respectively. In this case, higher throughput indicates higher efficiency of the explanation process. Table 4 indicates the execution time of the deployment phase. The experimental results show that the efficient amortized methods (e.g., CoRTX and FastSHAP) obtain significantly larger algorithmic throughput than efficient non-amortized methods (e.g., SHEAR). In addition, efficient non-amortized methods significantly improve the deployment time compared to traditional non-amortized methods. Instead of fixing the time to count the executed instances (e.g., algorithm throughput), some other works [54] use the metric by fixing the tested instances and reckoning the execution time.

Despite the significant improvement in the efficiency of non-amortized ones, there is still a considerable gap in the deployment of real-time online services. Compared to amortized methods that only require an extremely short time to derive explanations, traditional non-amortized ones, such as SHAP [7], can take up to 200x longer to complete the explanation process [54], [14]. While non-amortized methods are more accurate than amortized ones, their slower execution time makes them impractical for real-time applications. Some gradient-based methods, such as [22], [24], are important explanation methods that yield a relatively faster explanation than other traditional non-amortized methods. However, gradient-based methods still need extra time to conduct the sampling process on each data instance, slowing the execution time while generating the instance-wise model explanation. The explanation derivation time highly depends on the scalability of sampling and testing instances. As a result, gradient-based methods are inadequate for deployment in real-time systems. In Table 3,

TABLE 3

Experiment results of feature attribution task on CIFAR-10 dataset from two previous works [54], [14]. "Millisecond" (ms) is the measurement unit of deployment time.

CIFAR-10 Dataset	CoRTX [14]	FastSHAP [54]	SHAP [7]	IG [22]	SmoothGrad [24]	GradCAM [23]	DeepSHAP [7]
Deployment Time (ms)	0.4	0.4	27221.4	54.6	60.0	22.8	323.4

TABLE 4

Experiment of feature attribution tasks on Adult dataset from two previous works [13], [14]. We utilize throughput to measure the deployment time. Higher throughput indicates higher efficiency of the explanation process.

Adult Dataset	CoRTX [14]	FastSHAP [54]	SHEAR [13]	SHAP [7]	Permutation Sampling [94]
Deployment Time (Throughput)	6202.85	6202.85	71.064	1.0940	22.077

we conduct an experiment on CIFAR-10 dataset to observe the efficiency of explanation methods. The experimental results reveal the execution time on the inference phase per image. The results show that amortized-based methods are much faster than other efficient non-amortized methods (e.g., gradient-based methods), indicating that amortized ones are still the best candidate to deploy on real-time systems.

5.3 Limitation on Using Scenario

The trade-off between efficacy and efficiency brings the limitation to choosing efficient XAI methods for certain using scenarios appropriately. One reason is that selecting the best algorithm for a given scenario necessitates balancing the two factors. This is because a highly effective but inefficient algorithm may be impractical for large-scale applications, whereas a highly efficient but ineffective algorithm may not provide useful solutions. One prior work [14] designs a flexible data proportional parameter on the amortized method, providing an extra degree of freedom to coordinate the balance of efficacy and efficiency based on the users' specific requirements of the problem at hand. Another example is interaction detection tasks. Users may choose different interaction detection techniques based on the prediction models. While explaining the impact of feature interaction toward multi-layer perceptron (MLP), NID and PID can provide fast and accurate explanations via detecting interactions of variable order, which are extremely time-consuming by using traditional statistical tests (e.g., four-way ANOVA). When it comes to analyzing the prediction models that are not in the zone of MLP-based models (e.g., probit models and logistic models), NID and PID are not well-suited, whereas traditional statistical tests, such as [97], are able to fit in this situation. Therefore, the decision to use amortized or non-amortized acceleration methods should be based on the specific needs of the user, including the type of prediction model being targeted and the tolerance for computation latency and performance.

6 RESEARCH CHALLENGES

Despite recent advances in efficient XAI, there are still several urgent challenges, particularly in the deployment of efficient explanation methods and the issues of their trustworthiness raised by stakeholders and regulations. These challenges

need to be addressed in order to successfully deploy XAI in real-world applications.

6.1 Deployment on Efficient XAI

The first research challenge lies on the deployment side. We present the deployment challenges for non-amortized acceleration and amortized acceleration methods in two parts: software and hardware. On the software side, we explore centralized and decentralized training. On the hardware side, we discuss the limitations of different devices and the impact on the deployment of acceleration methods. However, only limited previous works [98] discuss the topic.

6.1.1 Non-amortized Acceleration Methods

The efficiency improvement from non-amortized acceleration methods makes it possible to deploy on real-world applications. Non-amortized ones obtain the advantage of deploying with decentralized training since each explainer is independent, which opens a big map on real-world deployment. Despite its natural advantage, the decentralized training paradigm may face a performance drop, since we cannot manage the training resource individually to meet Pareto efficiency. It is also difficult to coordinate conflicts consistently made by different local explainers when encountering similar data instances but with a different model explanation. One of the possible ways is to use federated learning [99] with a centralized working node. On the hardware side, efficient non-amortized methods are possibly applicable to the edge-computing paradigm. One of the main reasons is that the requirement for computational resources is relatively lower than traditional XAI local methods. Different from cloud computing platforms, edge computing has the advantages of low latency of network congestion but limited data storage and computing resources, which efficient non-amortized ones can partially handle. Additionally, it is essential to establish a standardized experimental setup for analyzing both i.i.d. and non-i.i.d. data distributions among edge clients. Finding the right balance between the software and hardware is still an open problem for deploying efficient non-amortized methods to real-world applications.

6.1.2 Amortized Acceleration Methods

Unlike non-amortized acceleration explainers that are separable and independent, amortized acceleration methods

aim to generate a unified explainer for real-time derivation, which requires a sufficient explanation dataset for centralized training. This may cause risks in collecting private user data from different devices, including data leakage or attack. One of the challenges in amortized ones is how to efficiently and effectively train them under the decentralized learning paradigm, which prevents the data gathering requests to the clients. In terms of the hardware side, amortized acceleration explainers can provide real-time services on both edge devices and online service platforms. Based on the current state-of-the-art efficient amortized ones, the computational requirements are much less than non-amortized ones, making it applicable to deploy on edge-computing devices with limited computational resources. However, a pre-trained amortized explainer may be large and complex, which requires substantial storage and computing capacity and obtains a long training process. In this way, the distributed training paradigm [100] can be one of the possible ways to cope with this difficulty. Generally, data parallelism and model parallelism are the two main types of distributed training. From a data parallelism perspective, a single worker device needs without large memory usage and disk storage, contributing to faster training time. However, there are times when the explanation models are too large to fit in a single worker device, and thus model parallelism is here to accelerate the training process from another perspective. As for the model parallelism, the explanation model itself is split into several parts that are trained simultaneously across different worker devices. Throughout sharing the model weights from different devices, the final explainer can be provided with a shorter training time with parallel architectures. Nevertheless, the design of an efficient training paradigm on amortized methods still remains an open problem in generating faithful and accurate model explanations.

6.2 Trustworthy Issues in Efficient XAI

The second research challenge involves issues of three parts: privacy, fairness, and security. Careful consideration must be given to these potential issues in order to ensure the trustworthy and responsible deployment of amortized acceleration methods of efficient XAI. We introduce the human-centric challenges below for non-amortized acceleration methods and amortized acceleration methods.

6.2.1 Issues on Non-amortized Acceleration Methods

Non-amortized acceleration methods in XAI could be subject to fairness and security issues. One of the issues is the fairness problem caused by sensitive attributes in the original prediction model. No matter in attribution or influential example task, a fair situation of model explanation should remain unchanged when sensitive attributes, such as race and gender in the criminal dataset, are considered [101]. Typically, sensitive attributions are the essential features of the biased prediction model, and data-centric methods, for example, can then easily pick up those sensitive attributions as essential features to further explainer learning. This can result in the inclusion of biased information from the prediction model, which will be further emphasized after the explainer learning, leading to fairness issues. Another one is the security issue in non-amortized ones. In non-amortized

settings, each individual user obtains their own explainer, which heavily overfits the data from a specific user. Model-centric acceleration methods can easily suffer from the data poison attack [102] since efficient optimization is sensitive to the adjustment of data values. In addition, data-centric methods are also vulnerable to data poison attacks, where the search optimization process can be significantly affected if the feature values are not of good quality. Thus, it is essential to carefully consider and address these issues in developing non-amortized acceleration methods for XAI.

6.2.2 Issues on Amortized Acceleration Methods

The use of amortized acceleration methods also poses an increased risk of encountering privacy, fairness, and security issues, as these methods rely on additional DNN models to learn the explanation distribution from individual users. One risk is the potential for privacy violations. Amortized acceleration methods require a large amount of private data to establish, since the targeted features to explain are usually the personal features of individual users. In an amortized setting, all users rely on the same explainer, which means that private information can be easily shared through the public access of an explainer without their consent or knowledge. This raises additional concerns about the security and privacy of sensitive information. Another issue faced by amortized ones is fairness issues led by the biased training patterns in the DNN-based explainer. This is the challenge for most DNN-based models trained with sensitive attributes [103]. The learning paradigm may tend to learn the spurious relevance between specific sensitive attributes and explanation results, causing algorithmic discrimination. The fairness issues should be carefully emphasized in the future direction of efficient XAI. Last but not least, one of the important issues that need to be awarded for amortized ones is security. Security measures are designed to protect user data from being hacked or stolen. However, these measures can be easily circumvented if the explainer is encoded with whole batches of user data, such as using the techniques of model inversion attack [104]. The attack can recover personal information from API access to amortized explanation models. This emphasizes the importance of security issues in amortized explanation models.

6.3 Towards Human-aware Selection

Efficient XAI provides fast and accurate model explanations to users, enabling them to trust the decision making of prediction models. An open question can be raised here: How can we effectively choose a type of explainer in order to provide a satisfactory explanation for certain scenarios? It is necessary to carefully choose the desired types of explainers according to the given tasks. For example, autopilot systems require the explanation derivation time to be extremely short for real-time judgments, whereas medical diagnosis cases are relatively flexible in the speed of generating model explanations. One criterion to consider when choosing an explainer is the tolerance of latency in deriving a model explanation. For example, in real-time bidding systems, users require real-time model explanations to support the real-time decisions made by bidding systems. However, efficient non-amortized methods may not be fast enough to meet the

tight deviation time constraints. Efficient amortized methods are able to meet the real-time requirements, as they only obtain a single forward pass during the inference phase. Another example is online streaming services, where users are expected to obtain fast and personalized explanations. Efficient non-amortized methods are particularly useful in this context because they typically provide personalized explanations. In this case, efficient non-amortized methods are more suitable to fit the requirements than amortized methods, as they provide each user with its own unique explainer. As a result, the usage scenarios of the efficient approaches are one of the essential criteria leading to an effective derivation process.

7 CONCLUSION

Efficient XAI is a rapidly evolving research area with significant demand from practical applications. This work provides a clear taxonomy and a comprehensive overview of existing techniques to aid practitioners and researchers in selecting the most suitable efficient algorithms for their specific needs. Examining the perspective of efficiency may be beneficial to the community in better understanding the limitations of existing XAI methods. Despite the progress in efficient XAI, there are still significant challenges that require future solutions to further emphasize the importance of efficiency. We hope that this work will serve as a valuable resource for both newcomers and professionals who are interested in the broad field of efficient XAI.

REFERENCES

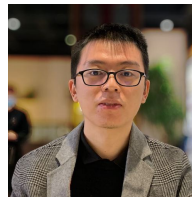
- [1] X. He, K. Deng, X. Wang, Y. Li, Y. Zhang, and M. Wang, "Lightgcn: Simplifying and powering graph convolution network for recommendation," in *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval*, 2020, pp. 639–648.
- [2] Y.-N. Chuang, C.-M. Chen, C.-J. Wang, M.-F. Tsai, Y. Fang, and E.-P. Lim, "Tpr: Text-aware preference ranking for recommender systems," in *Proceedings of the 29th ACM International Conference on Information & Knowledge Management*, 2020, pp. 215–224.
- [3] R. Dabre, C. Chu, and A. Kunchukuttan, "A survey of multilingual neural machine translation," *ACM Computing Surveys (CSUR)*, vol. 53, no. 5, pp. 1–38, 2020.
- [4] N. Chandollikar, C. Joshi, P. Roy, A. Gawas, and M. Vishwakarma, "Voice recognition: A comprehensive survey," in *2022 International Mobile and Embedded Technology Conference (MECON)*. IEEE, 2022, pp. 45–51.
- [5] C. Molnar, "Interpretable machine learning: A guide for making black box models explainable," 2022.
- [6] M. Du, N. Liu, and X. Hu, "Techniques for interpretable machine learning," *Communications of the ACM*, vol. 63, no. 1, pp. 68–77, 2019.
- [7] S. M. Lundberg and S.-I. Lee, "A unified approach to interpreting model predictions," in *Proceedings of the 31st international conference on neural information processing systems*, 2017, pp. 4768–4777.
- [8] S. Wachter, B. Mittelstadt, and C. Russell, "Counterfactual explanations without opening the black box: Automated decisions and the gdpr," *Harv. JL & Tech.*, vol. 31, p. 841, 2017.
- [9] S. M. Lundberg, G. Erion, H. Chen, A. DeGrave, J. M. Prutkin, B. Nair, R. Katz, J. Himmelfarb, N. Bansal, and S.-I. Lee, "From local explanations to global understanding with explainable ai for trees," *Nature machine intelligence*, vol. 2, no. 1, pp. 56–67, 2020.
- [10] D. Slack, A. Hilgard, S. Singh, and H. Lakkaraju, "Reliable post hoc explanations: Modeling uncertainty in explainability," *Advances in Neural Information Processing Systems*, vol. 34, pp. 9391–9404, 2021.
- [11] C. Agarwal, N. Johnson, M. Pawelczyk, S. Krishna, E. Saxena, M. Zitnik, and H. Lakkaraju, "Rethinking stability for attribution-based explanations," *arXiv preprint arXiv:2203.06877*, 2022.
- [12] M. Sundararajan, K. Dhamdhere, and A. Agarwal, "The shapley taylor interaction index," in *International conference on machine learning*. PMLR, 2020, pp. 9259–9268.
- [13] G. Wang, Y.-N. Chuang, M. Du, F. Yang, Q. Zhou, P. Tripathi, X. Cai, and X. Hu, "Accelerating shapley explanation via contributive cooperator selection," *arXiv preprint arXiv:2206.08529*, 2022.
- [14] Y.-N. Chuang, G. Wang, F. Yang, Z. Quan, P. Tripathi, X. Cai, and X. Hu, "CoRTX: Contrastive framework for real-time explanation," in *Openreview*, 2023. [Online]. Available: <https://openreview.net/forum?id=L2MUOU0beo>
- [15] H. W. Kuhn and A. W. Tucker, *Contributions to the Theory of Games*. Princeton University Press, 1953, vol. 2.
- [16] S. Verma, J. Dickerson, and K. Hines, "Counterfactual explanations for machine learning: A review," *arXiv preprint arXiv:2010.10596*, 2020.
- [17] M. Tsang, S. Rambhatla, and Y. Liu, "How does this interaction affect me? interpretable attribution for feature interactions," *Advances in neural information processing systems*, vol. 33, pp. 6147–6159, 2020.
- [18] I. Covert and S.-I. Lee, "Improving kernelshap: Practical shapley value estimation using linear regression," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 3457–3465.
- [19] V. Petsiuk, A. Das, and K. Saenko, "Rise: Randomized input sampling for explanation of black-box models," *arXiv preprint arXiv:1806.07421*, 2018.
- [20] J. Castro, D. Gómez, and J. Tejada, "Polynomial calculation of the shapley value based on sampling," *Computers & Operations Research*, vol. 36, no. 5, pp. 1726–1730, 2009.
- [21] R. Mitchell, J. Cooper, E. Frank, and G. Holmes, "Sampling permutations for shapley value estimation," 2022.
- [22] M. Sundararajan, A. Taly, and Q. Yan, "Axiomatic attribution for deep networks," in *International Conference on Machine Learning*. PMLR, 2017, pp. 3319–3328.
- [23] R. R. Selvaraju, M. Cogswell, A. Das, R. Vedantam, D. Parikh, and D. Batra, "Grad-cam: Visual explanations from deep networks via gradient-based localization," in *Proceedings of the IEEE international conference on computer vision*, 2017, pp. 618–626.
- [24] D. Smilkov, N. Thorat, B. Kim, F. Viégas, and M. Wattenberg, "Smoothgrad: removing noise by adding noise," *arXiv preprint arXiv:1706.03825*, 2017.
- [25] B. Ustun, A. Spangher, and Y. Liu, "Actionable recourse in linear classification," in *Proceedings of the conference on fairness, accountability, and transparency*, 2019, pp. 10–19.
- [26] A.-H. Karimi, J. Von Kügelgen, B. Schölkopf, and I. Valera, "Algorithmic recourse under imperfect causal knowledge: a probabilistic approach," *Advances in neural information processing systems*, vol. 33, pp. 265–277, 2020.
- [27] A.-H. Karimi, B. Schölkopf, and I. Valera, "Algorithmic recourse: from counterfactual explanations to interventions," in *Proceedings of the 2021 ACM conference on fairness, accountability, and transparency*, 2021, pp. 353–362.
- [28] Y. Koren, R. Bell, and C. Volinsky, "Matrix factorization techniques for recommender systems," *Computer*, vol. 42, no. 8, p. 30–37.
- [29] S. Kabbur, X. Ning, and G. Karypis, "Fism: Factored item similarity models for top-n recommender systems," in *Proceedings of the 19th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining*, ser. KDD '13. Association for Computing Machinery, p. 659–667.
- [30] I. Goodfellow, J. Pouget-Abadie, M. Mirza, B. Xu, D. Warde-Farley, S. Ozair, A. Courville, and Y. Bengio, "Generative adversarial networks," *Communications of the ACM*, vol. 63, no. 11, pp. 139–144, 2020.
- [31] M. Mirza and S. Osindero, "Conditional generative adversarial nets," *arXiv preprint arXiv:1411.1784*, 2014.
- [32] M. Hausknecht, P. Mupparaju, S. Subramanian, S. Kalyanakrishnan, and P. Stone, "Half field offense: An environment for multiagent learning and ad hoc teamwork," in *AAMAS Adaptive Learning Agents (ALA) Workshop*, vol. 3. sn, 2016.
- [33] W. Masson, P. Ranchod, and G. Konidaris, "Reinforcement learning with parameterized actions," in *Thirtieth AAAI Conference on Artificial Intelligence*, 2016.
- [34] E. Wei, D. Wicke, and S. Luke, "Hierarchical approaches for reinforcement learning in parameterized action space," *arXiv preprint arXiv:1810.09656*, 2018.
- [35] J. Yang, "Fast treeshap: Accelerating shap value computation for trees," *arXiv preprint arXiv:2109.09847*, 2021.

- [36] M. Jullum, A. Redelmeier, and K. Aas, "groupshapley: efficient prediction explanation with shapley values for feature groups," *arXiv preprint arXiv:2106.12228*, 2021.
- [37] J. Chen, L. Song, M. J. Wainwright, and M. I. Jordan, "L-shapley and c-shapley: Efficient model interpretation for structured data," *arXiv preprint arXiv:1808.02610*, 2018.
- [38] M. Ancona, C. Oztireli, and M. Gross, "Explaining deep neural networks with a polynomial time algorithm for shapley value approximation," in *International Conference on Machine Learning*. PMLR, 2019, pp. 272–281.
- [39] A. Ghorbani, M. Kim, and J. Zou, "A distributional framework for data valuation," in *International Conference on Machine Learning*. PMLR, 2020, pp. 3535–3544.
- [40] U. Bhatt, A. Weller, and J. M. Moura, "Evaluating and aggregating feature-based model explanations," *arXiv preprint arXiv:2005.00631*, 2020.
- [41] M. Lomeli, M. Rowland, A. Gretton, and Z. Ghahramani, "Antithetic and monte carlo kernel estimators for partial rankings," *arXiv preprint arXiv:1807.00400*, 2018.
- [42] C. Russell, "Efficient search for diverse coherent explanations," in *Proceedings of the Conference on Fairness, Accountability, and Transparency*, 2019, pp. 20–28.
- [43] K. Mohammadi, A.-H. Karimi, G. Barthe, and I. Valera, "Scaling guarantees for nearest counterfactual explanations," in *Proceedings of the 2021 AAAI/ACM Conference on AI, Ethics, and Society*, 2021, pp. 177–187.
- [44] A. Parmentier and T. Vidal, "Optimal counterfactual explanations in tree ensembles," in *International Conference on Machine Learning*. PMLR, 2021, pp. 8422–8431.
- [45] M. Á. Carreira-Perpiñán and S. S. Hada, "Counterfactual explanations for oblique decision trees: Exact, efficient algorithms," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 35, no. 8, 2021, pp. 6903–6911.
- [46] K. Kanamori, T. Takagi, K. Kobayashi, and H. Arimura, "Dace: Distribution-aware counterfactual explanation by mixed-integer linear optimization." in *IJCAI*, 2020, pp. 2855–2862.
- [47] E. Albin, A. Rago, P. Baroni, and F. Toni, "Influence-driven explanations for bayesian network classifiers," in *Pacific Rim International Conference on Artificial Intelligence*. Springer, 2021, pp. 88–100.
- [48] A. V. Looveren and J. Klaise, "Interpretable counterfactual explanations guided by prototypes," in *Joint European Conference on Machine Learning and Knowledge Discovery in Databases*. Springer, 2021, pp. 650–665.
- [49] E. Albin, J. Long, D. Dervovic, and D. Magazzeni, "Counterfactual shapley additive explanations," in *2022 ACM Conference on Fairness, Accountability, and Transparency*, 2022, pp. 1054–1070.
- [50] A. Karczmarz, A. Mukherjee, P. Sankowski, and P. Wygocki, "Improved feature importance computations for tree models: Shapley vs. banzhaf," *arXiv preprint arXiv:2108.04126*, 2021.
- [51] M. Tsang, D. Cheng, and Y. Liu, "Detecting statistical interactions from neural network weights," in *International Conference on Learning Representations*, 2018. [Online]. Available: <https://openreview.net/forum?id=ByOfBggRZ>
- [52] Z. Liu, Q. Song, K. Zhou, T.-H. Wang, Y. Shan, and X. Hu, "Detecting interactions from neural networks via topological analysis," *Advances in Neural Information Processing Systems*, vol. 33, pp. 6390–6401, 2020.
- [53] C.-P. Tsai, C.-K. Yeh, and P. Ravikumar, "Faith-shap: The faithful shapley shapley interaction index," *arXiv preprint arXiv:2203.00870*, 2022.
- [54] N. Jethani, M. Sudarshan, I. C. Covert, S.-I. Lee, and R. Ranganath, "Fastshap: Real-time shapley value estimation," in *International Conference on Learning Representations*, 2021.
- [55] J. Chen, L. Song, M. Wainwright, and M. Jordan, "Learning to explain: An information-theoretic perspective on model interpretation," in *International Conference on Machine Learning*. PMLR, 2018, pp. 883–892.
- [56] P. Schwab and W. Karlen, "Cxpain: Causal explanations for model interpretation under uncertainty," *arXiv preprint arXiv:1910.12336*, 2019.
- [57] F. Yang, N. Liu, M. Du, and X. Hu, "Generative counterfactuals for neural networks via attribute-informed perturbation," *ACM SIGKDD Explorations Newsletter*, vol. 23, no. 1, pp. 59–68, 2021.
- [58] F. Yang, S. S. Alva, J. Chen, and X. Hu, "Model-based counterfactual synthesizer for interpretation," *arXiv preprint arXiv:2106.08971*, 2021.
- [59] S. Khorram and L. Fuxin, "Cycle-consistent counterfactuals by latent transformations," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2022, pp. 10 203–10 212.
- [60] M. Downs, J. L. Chu, Y. Yacoby, F. Doshi-Velez, and W. Pan, "Cruds: Counterfactual recourse using disentangled subspaces," *ICML WHI*, vol. 2020, pp. 1–23, 2020.
- [61] P. Rodríguez, M. Caccia, A. Lacoste, L. Zamparo, I. Laradji, L. Charlin, and D. Vazquez, "Beyond trivial counterfactual explanations with diverse valuable explanations," in *Proceedings of the IEEE/CVF International Conference on Computer Vision*, 2021, pp. 1056–1065.
- [62] J. Antorán, U. Bhatt, T. Adel, A. Weller, and J. M. Hernández-Lobato, "Getting a clue: A method for explaining uncertainty estimates," *arXiv preprint arXiv:2006.06848*, 2020.
- [63] Z. Chen, F. Silvestri, J. Wang, H. Zhu, H. Ahn, and G. Tolomei, "Relax: Reinforcement learning agent explainer for arbitrary predictive models," *arXiv preprint arXiv:2110.11960*, 2021.
- [64] S. Verma, K. Hines, and J. P. Dickerson, "Amortized generation of sequential algorithmic recourses for black-box models," in *Proceedings of the AAAI Conference on Artificial Intelligence*, vol. 36, no. 8, 2022, pp. 8512–8519.
- [65] A.-H. Karimi, G. Barthe, B. Balle, and I. Valera, "Model-agnostic counterfactual explanations for consequential decisions," in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2020, pp. 895–905.
- [66] S. Sharma, J. Henderson, and J. Ghosh, "Certifai: Counterfactual explanations for robustness, transparency, interpretability, and fairness of artificial intelligence models," *arXiv preprint arXiv:1905.07857*, 2019.
- [67] V. Kaffes, D. Sacharidis, and G. Giannopoulos, "Model-agnostic counterfactual explanations of recommendations," in *Proceedings of the 29th ACM Conference on User Modeling, Adaptation and Personalization*, 2021, pp. 280–285.
- [68] M. G. Kendall, "A new measure of rank correlation," *Biometrika*, vol. 30, no. 1/2, pp. 81–93, 1938.
- [69] S. Dandl, C. Molnar, M. Binder, and B. Bischl, "Multi-objective counterfactual explanations," in *International Conference on Parallel Problem Solving from Nature*. Springer, 2020, pp. 448–469.
- [70] A. Artelt and B. Hammer, "Efficient computation of counterfactual explanations and counterfactual metrics of prototype-based classifiers," *Neurocomputing*, vol. 470, pp. 304–317, 2022.
- [71] M. Tsang, H. Liu, S. Purushotham, P. Murali, and Y. Liu, "Neural interaction transparency (nit): Disentangling learned interactions for improved interpretability," *Advances in Neural Information Processing Systems*, vol. 31, 2018.
- [72] M. Tsang, D. Cheng, H. Liu, X. Feng, E. Zhou, and Y. Liu, "Feature interaction interpretability: A case for explaining ad-recommendation systems via neural interaction detection," in *International Conference on Learning Representations*, 2020. [Online]. Available: <https://openreview.net/forum?id=BkgnhTEtDS>
- [73] M. Grabisch and M. Roubens, "An axiomatic approach to the concept of interaction among players in cooperative games," *International Journal of game theory*, vol. 28, no. 4, pp. 547–565, 1999.
- [74] M. Hamilton, S. Lundberg, L. Zhang, S. Fu, and W. T. Freeman, "Axiomatic explanations for visual search, retrieval, and similarity learning," *arXiv preprint arXiv:2103.00370*, 2021.
- [75] R. Wang, X. Wang, and D. I. Inouye, "Shapley explanation networks," *arXiv preprint arXiv:2104.02297*, 2021.
- [76] I. Covert, C. Kim, and S.-I. Lee, "Learning to estimate shapley values with vision transformers," *arXiv preprint arXiv:2206.05282*, 2022.
- [77] P. Dabkowski and Y. Gal, "Real time image saliency for black box classifiers," *Advances in neural information processing systems*, vol. 30, 2017.
- [78] A. Kanehira and T. Harada, "Learning to explain with complementary examples," in *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, 2019, pp. 8603–8611.
- [79] W. Fu, M. Wang, M. Du, N. Liu, S. Hao, and X. Hu, "Differentiated explanation of deep neural networks with skewed distributions," *IEEE Transactions on Pattern Analysis and Machine Intelligence*, vol. 44, no. 6, pp. 2909–2922, 2021.
- [80] A. V. Konstantinov and L. V. Utkin, "Attention-like feature explanation for tabular data," *International Journal of Data Science and Analytics*, pp. 1–26, 2022.
- [81] R. Schwarzenberg, N. Feldhus, and S. Möller, "Efficient explanations from empirical explainers," *arXiv preprint arXiv:2103.15429*, 2021.

- [82] R. Hesse, S. Schaub-Meyer, and S. Roth, "Fast axiomatic attribution for neural networks," *Advances in Neural Information Processing Systems*, vol. 34, pp. 19 513–19 524, 2021.
- [83] N. Jethani, M. Sudarshan, Y. Aphinyanaphongs, and R. Ranganath, "Have we learned to explain?: How interpretability methods can learn to encode predictions in their interpretations." in *International Conference on Artificial Intelligence and Statistics*. PMLR, 2021, pp. 1459–1467.
- [84] C. W. Granger, "Investigating causal relations by econometric models and cross-spectral methods," *Econometrica: journal of the Econometric Society*, pp. 424–438, 1969.
- [85] S. Singla, B. Pollack, J. Chen, and K. Batmanghelich, "Explanation by progressive exaggeration," *arXiv preprint arXiv:1911.00483*, 2019.
- [86] S. Singla, B. Pollack, S. Wallace, and K. Batmanghelich, "Explaining the black-box smoothly—a counterfactual approach," *arXiv preprint arXiv:2101.04230*, 2021.
- [87] M. Pawelczyk, K. Broelemann, and G. Kasneci, "Learning model-agnostic counterfactual explanations for tabular data," in *Proceedings of The Web Conference 2020*, 2020, pp. 3126–3132.
- [88] S. Joshi, O. Koyejo, W. Vijitbenjaronk, B. Kim, and J. Ghosh, "Towards realistic individual recourse and actionable explanations in black-box decision making systems," *arXiv preprint arXiv:1907.09615*, 2019.
- [89] V. Guyomard, F. Fessant, T. Guyet, T. Bouadi, and A. Termier, "Vcnet: A self-explaining model for realistic counterfactual generation," in *Proceedings of the European Conference on Machine Learning and Principles and Practice of Knowledge Discovery in Databases (ECML PKDD)*, 2022.
- [90] D. Mahajan, C. Tan, and A. Sharma, "Preserving causal constraints in counterfactual explanations for machine learning classifiers," *arXiv preprint arXiv:1912.03277*, 2019.
- [91] S. Joshi, O. Koyejo, B. Kim, and J. Ghosh, "xgems: Generating exemplars to explain black-box models," *arXiv preprint arXiv:1806.08867*, 2018.
- [92] J. Ma, R. Guo, S. Mishra, A. Zhang, and J. Li, "Clear: Generative counterfactual explanations on graphs," *arXiv preprint arXiv:2210.08443*, 2022.
- [93] D. P. Kingma and M. Welling, "Auto-encoding variational bayes," *arXiv preprint arXiv:1312.6114*, 2013.
- [94] R. Mitchell, J. Cooper, E. Frank, and G. Holmes, "Sampling permutations for shapley value estimation," *arXiv preprint arXiv:2104.12199*, 2021.
- [95] H. B. McMahan, G. Holt, D. Sculley, M. Young, D. Ebner, J. Grady, L. Nie, T. Phillips, E. Davydov, D. Golovin *et al.*, "Ad click prediction: a view from the trenches," in *Proceedings of the 19th ACM SIGKDD international conference on Knowledge discovery and data mining*, 2013, pp. 1222–1230.
- [96] D. A. Teich and P. R. Teich, "Plaster: A framework for deep learning performance," Tech. rep. TIRIAS Research, Tech. Rep., 2018.
- [97] C. Ai and E. C. Norton, "Interaction terms in logit and probit models," *Economics letters*, vol. 80, no. 1, pp. 123–129, 2003.
- [98] J. L. C. Bárcena, M. Daole, P. Ducange, F. Marcelloni, A. Renda, F. Ruffini, and A. Schiavo, "Fed-xai: Federated learning of explainable artificial intelligence models," 2022.
- [99] M. Aledhari, R. Razzak, R. M. Parizi, and F. Saeed, "Federated learning: A survey on enabling technologies, protocols, and applications," *IEEE Access*, vol. 8, pp. 140 699–140 725, 2020.
- [100] Z. Jia, M. Zaharia, and A. Aiken, "Beyond data and model parallelism for deep neural networks." *Proceedings of Machine Learning and Systems*, vol. 1, pp. 1–13, 2019.
- [101] M. J. Kusner, J. Loftus, C. Russell, and R. Silva, "Counterfactual fairness," *Advances in neural information processing systems*, vol. 30, 2017.
- [102] J. Lin, L. Dang, M. Rahouti, and K. Xiong, "MI attack models: Adversarial attacks and data poisoning attacks," *arXiv preprint arXiv:2112.02797*, 2021.
- [103] M. Du, F. Yang, N. Zou, and X. Hu, "Fairness in deep learning: A computational perspective," *IEEE Intelligent Systems*, vol. 36, no. 4, pp. 25–34, 2020.
- [104] A. Chakraborty, M. Alam, V. Dey, A. Chattopadhyay, and D. Mukhopadhyay, "Adversarial attacks and defences: A survey," *arXiv preprint arXiv:1810.00069*, 2018.



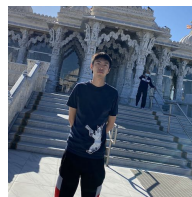
Yu-Neng Chuang is a second-year Ph.D. student in the Computer Science Department at Rice University. His research interest lies in trustworthy machine learning, including explainable artificial intelligence, machine learning fairness, and recommender systems. Currently, he is focusing on the efficiency of interpretable machine learning and fair modeling. Prior to Rice, his previous research also involves recommender systems via modeling user behaviors and sparse labeling of textual data. Yu-Neng received his B.S. and M.S. degrees in Mathematics and Computer Science from National Chengchi University, respectively in 2017 and 2020.



Guanchu Wang Guanchu Wang received the BS and MS degrees in electrical engineering and information science from Dalian University of Technology and University of Science and Technology of China, respectively. He is currently a Ph.D. student in the department of computer science at Rice university. His research focuses on efficient and interpretable machine learning.



Fan Yang Fan Yang is currently a final-year Ph.D. student in the Computer Science Department at Rice University. His research interests generally lie in the area of eXplainable Artificial Intelligence (XAI), with a major focus on model interpretation techniques, counterfactual explanation, and machine learning fairness. He is also interested in XAI-related downstream application, as well as its correlative intersections with Natural Language Processing and Human-Computer Interaction. Prior to Rice, Fan had research experiences on wireless communication and networking. Fan received his M.S. and B.S. degree in Xidian University, respectively in 2016 and 2013.



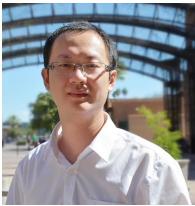
Zirui Liu Zirui Liu received BS and MS degrees in electrical engineering both from Harbin Institute of Technology. He is currently a Ph.D. student in the Department of Computer Science at Rice University. His research interests include large-scale machine learning and graph neural networks.



Xuanting Cai is a software engineer at Meta Platforms, Incorporated. Xuanting Cai earned his Ph.D. in Mathematics from Louisiana State University and B.S. in Mathematics and Economics from Peking University. Before joining Meta, he worked at Alibaba.com and Google as a software engineer.



Mengnan Du is an Assistant Professor in the Department of Data Science, New Jersey Institute of Technology (NJIT). Mengnan Du earned his Ph.D. in Computer Science from Texas A&M University. He has previously worked/interned with Microsoft Research (MSR), Adobe Research, Intel, Baidu Research, Baidu Search Science and JD Explore Academy. His research covers a wide range of trustworthy machine learning topics, such as model explainability, fairness, and robustness. He has had more than 40 papers published in prestigious venues such as NeurIPS, AAAI, KDD, WWW, ICLR, and ICML. He received over 2,300 citations with an H-index of 16.



Xia "Ben" Hu is an Associate Professor at Rice University in the Department of Computer Science. Dr. Hu has published over 100 papers in several major academic venues, including NeurIPS, ICLR, KDD, WWW, IJCAI, AAAI, etc. An open-source package developed by his group, namely AutoKeras, has become the most used automated deep learning system on Github (with over 8,000 stars and 1,000 forks). His papers have received several Best Paper (Candidate) awards from venues such as ICML, WWW, WSDM, ICDM, AMIA and INFORMS. He is the recipient of NSF CAREER Award and ACM SIGKDD Rising Star Award. His work has been cited more than 16,000 times.