# Causality-Inspired Taxonomy for Explainable Artificial Intelligence

PEDRO C. NETO, INESC TEC, Portugal and Faculty of Engineering - University of Porto, Portugal

TIAGO GONÇALVES, INESC TEC, Portugal and Faculty of Engineering - University of Porto, Portugal

JOÃO RIBEIRO PINTO, Vision-Box, Portugal

WILSON SILVA, Department of Information and Computing Sciences, and Department of Biology - Utrecht University, The Netherlands and Department of Radiology - The Netherlands Cancer Institute, The Netherlands

ANA F. SEQUEIRA, INESC TEC, Portugal

ARUN ROSS, Michigan State University, United States

JAIME S. CARDOSO, Faculty of Engineering - University of Porto, Portugal and INESC TEC, Portugal

As two sides of the same coin, causality and explainable artificial intelligence (xAI) were initially proposed and developed with different goals. However, the latter can only be complete when seen through the lens of the causality framework. As such, we propose a novel causality-inspired framework for xAI that creates an environment for the development of xAI approaches. To show its applicability, biometrics was used as case study. For this, we have analysed 81 research papers on a myriad of biometric modalities and different tasks. We have categorised each of these methods according to our novel xAI Ladder and discussed the future directions of the field.

## 1 INTRODUCTION

Explainable Artificial Intelligence (xAI) and Causality have been two sides of the same coin for years. Over time, xAI has grown in interest and popularity due to the need to explain the ever complex artificial intelligence systems that have been developed and deployed by large companies, or research groups. Causality on the other end aimed to bring a common mathematical language to explain chains of relationships that cannot be simply described by statistics. In this document, we leverage this duality to propose a novel causality-inspired taxonomy for explainable artificial intelligence. To show its applicability, we have chosen biometrics as our case study. As such, the following document delves into biometrics and explainable artificial intelligence through the lens of causality.

Biometric systems keep growing and improving at an extremely fast pace. The SarS-COV-2 pandemic has further propelled the use of contactless biometric systems, which were already expected to be widely adopted in 2021 [119, 167]. The real-world applications of biometric systems are quite diverse and range from face detection for virtual reality filters [50] to offline smartphone authentication [215] and airport security [184, 237]. This wide adoption implies that these systems will be used by/on individuals of varied sexes, ethnicities, and demographics overall. Hence, it is important to ensure that the systems in production are capable of handling all its potential users in an equal and fair manner.

The word biometrics comes from Morris' definition in 1875 as the combination of the words Bio (life) and Metron (a measure). In other words, it consists of the usage of body traits and behavioural cues to perform measurement and analysis. This definition has endured over the years, and its applications evolved throughout the years. These systems are comprised of three major tasks: enrolment, authentication and identification (Fig. 1). Enrolment allows users to insert their biometric data into a gallery and link it to their identity. Authentication receives an identity claim and biometric information from the user, and performs a 1:1 verification with the provided information and the data in

the gallery linked to the claimed identity. Finally, for identification, the system receives only the biometric data and tries to compare it with the entire gallery (1:N comparisons); if it does not match any identity it returns no identity. Within these tasks, there are secondary tasks that aim to secure the system, such as presentation attack detection and morphing attack detection. As mentioned before, the widespread usage of biometrics implies that all these systems are consistently identifying and authenticating humans.
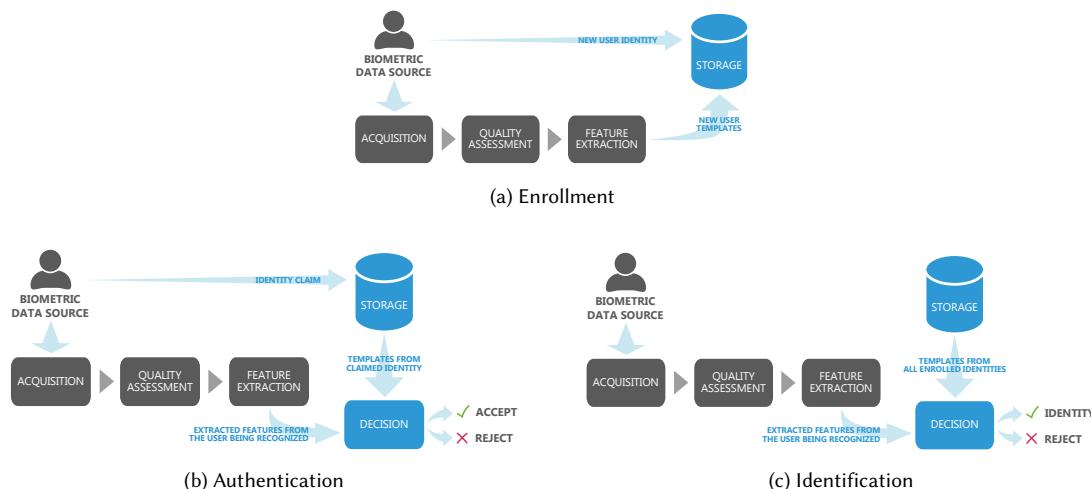


(a) Enrollment

(b) Authentication

(c) Identification

Fig. 1. Different tasks of biometric systems. In (a) users add their information to the database. In (b) users claim to be someone and and the system attempts an 1:1 comparison. In (c) the system tries to match the user with all the identities in an 1:N problem.

As seen in Fig. 1, the majority of the modules are common to all the three main tasks. This structure redistributes responsibilities and allows for smaller and incremental improvements on specific modules. Their responsibility within the system is the following:

- *User*: Without the user, there is no biometric system. It is the most central role in the entire system, it provides all the data, and systems are designed to be used by them. Due to the strong regulations imposed by the European General Data Protection Regulation (GDPR) [62], the user is, more than ever, the central part of these systems.
- *Acquisition*: The acquisition module is the bridge between the user data and the algorithm. Hence, as a bridge, it must comprise two different properties. First, it must be practical and benefit from a high usability score. Secondly, it must provide quality biometric data to the system and to the following modules. The acquisition process must be responsible for reducing the excessive noise and artefacts, which might require specialised hardware to read and digitise the biometric information.
- *Quality assessment*: Despite the efforts applied to the acquisition module, some noise might still be present in the data, which might be caused by an intentional attack on the system. This module aims to further remove the noise, and discard attacks or unusable samples.
- *Feature extraction*: Pattern recognition and signal processing based algorithms are capable of extracting meaningful features from input data. These extracted features represent the input signal mapped into a latent space that allows us to directly compare different samples. Researchers aim to create latent spaces that have a high

variance between samples from different identities, and a low variance between samples of the same identity. In other words, discrimination between different identities is facilitated by this representation.

- *Storage*: The storage of biometric templates is also an important element of biometric systems. On one hand, when the number of stored templates increases, the query response time and processing time are expected to increase too, and thus, this module is crucial for scalability. Moreover, there are security concerns and procedures that surround the storage of biometric templates, their link to the user identity and how they can be explored by harmful entities.
- *Decision*: This final model is responsible for the decision. This is one of the most application-dependent modules, since, different threshold decisions and different application goals are involved. It can either return the Boolean result of an identity claim (authentication) or an identity (identification). Some systems further include a confidence metric that supports the usage of a reject-option for uncertain samples.

Since the early 2010's these biometric systems have also benefited from the exponential growth of machine learning (ML) methods. ML-based systems excel in most of the artificial intelligence (AI) fields, and have outperformed other methods, as well as humans, in certain tasks. Undoubtedly, the success of AI systems is mainly due to three factors: i) improvements of deep learning (DL) methods; ii) availability of large databases; and iii) computational gains obtained with powerful graphics processing unit (GPU) cards [102, 113]. DL-based neural networks extract feature representations from data, which are characterised by high compactness and strong discriminability. In other words, the most affected module by the rise of deep learning systems has been the feature extraction module (Fig 1). It evolved from extracting handcrafted features to directly feed the inputs to deep neural networks. Feature selection is also optimised in a end-to-end fashion through the minimisation of the loss function. Moreover, the quality assessment module becomes less critical to the system since the features became significantly more robust to noise (for instance, models that can perform face recognition despite occlusions [151]). These performance gains were propelled by a significant increase in model complexity. In other words, the literature suggests that the opaqueness of a machine learning model may be positively correlated to its performance [178]. As seen in recent competitions dedicated to challenges related to face biometrics, the proposed solutions are mainly focused on deep learning approaches [22, 90, 147].

The existence of opaque models creates space for hidden biases, privacy issues, vulnerabilities to adversarial attacks and lack of transparency. And some of these vulnerabilities have already been shown to be present in some biometric systems [203, 209]. In fact, it is challenging to directly detect them due to the complexity of the models used. Hence, the focus is being redirected into a) using models that are inherently interpretable; b) developing methods to explain the predictions of black box systems; c) incorporating fairness and other similar metrics into the optimisation process of the model. The relevance of these topics and research directions creates the need for a study that gathers the findings, successes, and most relevant research questions for explainable artificial intelligence (xAI) within biometric systems.

The main contributions of this paper includes a chronological road-map of the shift of biometric systems into black box ones, the presentation of the current challenges of these systems, the description of the available methods for explainable artificial intelligence, and their limitations. Concluding with the anticipated future directions of xAI applied to biometric systems. It includes a review of why explainability is necessary, an analysis of works both within biometric and computer vision domains, and future directions. Hence, besides this introductory Section, the document is divided into four other major Sections. First, Section 2 discusses the current state of biometric systems as well as the need to integrate explainable artificial intelligence methods in their design. Afterwards, fundamental concepts and the background of xAI approaches are described in detail in Section 3, which finalises with methods that already started

to merge xAI and biometrics. This Section is followed by Section 4 that aims to illustrate the benefits of this merge and potential strategies to integrate both research areas. Finally, Section 5 concludes the main ideas of the paper and presents future research directions that aim to develop fair and transparent biometrics systems.

## 2 WHY BIOMETRICS IS NOT INTERPRETABLE (BUT SHOULD BE)

Before delving into explainable artificial intelligence concepts, it is important to understand how biometric systems evolved through time. Hence, this Section presents a brief introduction to the reasons that led biometric systems to become black box and the reasons that led researchers to focus on increasing the understanding of the inner workings of these systems.

### 2.1 How did biometric systems become opaque?

Biometric systems are everywhere, from smartphones and laptops to border control and other sensitive areas. Both experts and non-experts have constant contact with these systems, but do they really understand what is under the hood? Do they trust biometric systems? To better understand the deep implications of these questions it is necessary to understand the current state of biometrics systems and their evolution over time.

Biometric systems have been used since late 1800's, with the introduction of fingerprints in 1892 [60, 72]. After a few years their underlining methodologies have started to follow the automation path since the mid-1900s. From these initial applications to the current methods, there has been significant progress regarding the precision and accuracy of the overall system. More concretely, the first known biometric system to benefit from automation was fingerprint-based biometrics, which started with methodologies focusing on minutiae points [95], replicating the human analysis. Then the methodologies moved to textural approaches, and currently, the DL-based techniques are predominant for this biometric trait [131]. Face biometrics started as a technique based on specific regions of the face - fiducial points - and was the first biometric trait to truly leverage DL methods [200]. Iris recognition was tackled with texture-based methods due to the texture richness of these images. Afterwards, this trait has also moved towards DL approaches [132].

Over the previous decade, researchers, excited with the quick progress, focused on improving the performance of deep learning-based biometric systems. As reported in the Face Recognition Vendor Test conducted by NIST [79], the massive gains in face recognition accuracy have coincided with the usage of convolutional neural networks to extract an array of features from a face image. It is also worth noting how significant these improvements (from 2013 to 2018) are when compared the ones made during the previous period (from 2010 to 2013) [79]. For instance, in a four-year period the performance of the best performing face recognition algorithm improved ∼20x [78, 79]. From the initial False Negative Identification Rate (FNIR) of 0.041, achieved by EC0C in 2014, to a FNIR of 0.002, attained by NEC-3 in 2018. In practice, this translates into a ∼20x smaller error. The rapid improvements of these models correlated positively with their opaqueness. And, despite their remarkable performances, these systems are often deployed in use-cases with serious implications regarding wrong predictions. Moreover, discrimination and other harms might cause these errors to impact more a certain demographic group. As such, fairness and transparency must be tackled together with the development of better, more efficient and accurate systems.

### 2.2 Why Biometric systems need to reclaim transparency?

Following the idea of the previous section, explainability should become a central piece in the development of biometric systems. The predictions given by a biometric system should be sustained by an explanation or an interpretable behaviour. Yet, the transparency and the interpretability of these systems are, as of now, negatively correlated with

their performance [80]. DL tools can be a double-edged sword. On the one hand, it provided grounds for achieving human-level performance, on the other hand, it was a step away from white-box approaches. The complexity and capability to fit the data displayed by biometric systems is also a potential cause for their sensitivity to the underlying demographic characteristics of the training data, resulting in performances that vary with the demographic attributes of the user. For instance, this discriminatory behaviour can be harmful if a human is not aware of this behaviour, and it can be displayed through the usage of non-causal attributes or lack of discriminative features for that specific demographic group. Wang *et al.* [218] have shown that the usage of a shared optimisation process for face recognition can penalise underrepresented, or more difficult skin tones. Similarly, Hull [93] discussed the implications of spurious correlations of dirty data on the models' predictions. It is important to retain that, differently from other fields, on biometrics, sensitive information such as gender, age, ethnicity and health status can be implicitly encoded in the inputs [134, 136, 137], or in the extracted features [201] and are not easily removable. For such reason, these systems and their data require several layers of protection and privacy. Neto *et al.* [149] have shown that a simple classifier is able to learn how to separate different ethnicity groups on the latent space of face recognition deep neural networks.

Their black-box nature provides a grounding for targeted attacks. For instance, some adversarial attacks are capable of completely fooling a system by adding small noisy patterns to the image without affecting its realism to humans [53, 244] (Fig. 2). Other attacks consist of including patterns, for instance on the outfit of a person, that make that person undetectable [204, 228, 235]. Moreover, research can be found that focus on backdoor attacks, which encode behaviours hidden within the network and can be triggered by a certain input [207, 231]. These attacks can be difficult to detect, especially if we are not aware of the expected behaviour of the system. Besides these attacks, which are common within the deep learning literature, biometrics systems are targeted by several other attacks, such as: presentation [232, 233], morphing [89] and deepfake attacks [160]. Recently, beautification filters, similar to the ones found on popular social networks, have also shown a potential to impact these biometrics systems [135].

For better performances, researchers are also trying to create multi-modal systems based on more than one biometric trait [16]. The resulting model is expected to be more robust, but at the same time more complex. Considering the already difficult task of explaining and interpreting a model designed to work with a single modality, it might be significantly harder to understand a model with multiple traits. If one of the traits can be explained, then it might lead to improvements in the overall explainability of the remaining traits [164].

Opaque systems are not trustable [58, 216], sometimes not even if an attempted explanation is provided [112]. In other words, black-box systems reduce users' trust, especially after a wrong prediction, since it hinders the ability to follow the *thought process* of the system. Due to its prominence and effects on users' trust, the use of black-box systems becomes one of the most impactful factors for the adoption of new technologies. Educating users' in the recent advances of biometric technology is crucial but fails to compensate for the lack of explainability of these systems. As such, the adoption is conditioned on the ability to transpose the thought process of the model to a common language between the system and the user. Moreover, due to the heavy use of personal data, users must be protected through regulations and inspections. A strong protection, especially if perceived by the users, further increases the trust in the system [1, 130]. With this in mind, several countries and organisations worldwide have defined some degree of regulation of biometrics. For instance, the European Union (EU), with their General Data Protection Regulation (GDPR) [62], limits the processing of biometric information. Exceptions include law-related scenarios and if the user explicitly consents. In the United

States of America (U.S.A.), at least three different states (Illinois, Texas and Washington) have laws heavily regulating biometric systems [1].
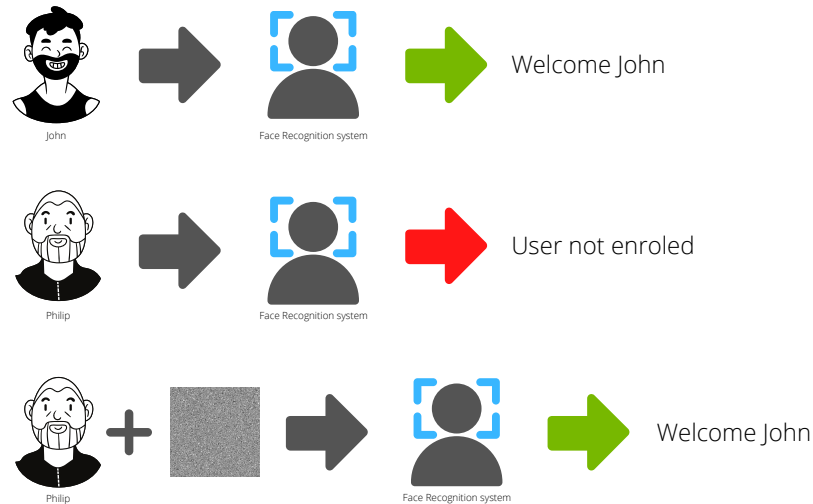


Fig. 2. Vulnerability of face recognition systems to natural-looking adversarial input images that drive the model to incorrect output predictions.

For years, researchers and engineers have joined forces to deploy several biometric systems at scale in real world scenarios. Despite their attempts to promote trust in those systems, and to prove that the systems have an unbiased and accurate behaviour, they have been shown to be wrong [46, 57, 106, 109, 198] and this technology has been described as not being production-ready. Srinivas *et al.* [198] have shown a performance difference between adults and children. Klare *et al.* [106] added both ethnicity and gender to the age as factors with an impact on the performance of face recognition systems. These consecutive and recurrent mistakes undermined the trustworthiness of these systems. In other words, the deployment of biased systems puts at risk the affected minorities, and further hurts the field as well as its growth potential. The first step for a better adoption of biometric systems is undoubtedly to provide users with an explanation of the algorithm decisions [49, 51].

The interpretation of black box models and their conversion to white systems has several purposes. Not only can they be used to provide explanations and answers that incite the user to believe in the answer given, but they can also be used to detect unexpected behaviour, attacks or biases. For this reason, developing interpretable models or explanation methods is crucial for the smooth implementation of biometric systems in high-risk situations. The relation between risk, trust and bias has been studied before by Lai *et al.* [111] and it highlights the need for transparent and explainable systems in order to achieve safe, trustworthy and unbiased models. With this in mind, researchers can aim

---

[1]https://www.huschblackwell.com/2023-state-biometric-privacy-law-tracker

to a careful deployment of biometric systems in large scale, further improving their capabilities and designing software that can work together with humans for more accurate solutions.

## 3 FUNDAMENTAL CONCEPTS AND BACKGROUND

As with any recent scientific field, the terminology of explainable artificial intelligence is surrounded by fog. There are several interpretations across the literature of the same terms. Hopefully, over time, the meaning of each term will converge to common ground. The first part of this Section is dedicated to defining these terms as they are used in the literature and to shed some light on how their definition can be improved. Once established the proposed interpretation of xAI terms, the second part of this Section discusses interpretability and explainability re-framed to a novel interpretation of these terms through the lens of causality. Finally, the third and last part of this Section presents existing biometric works that already include xAI concepts.

### 3.1 Towards a common terminology for xAI in Biometrics: Can we speak the same language?

The definitions of interpretability and explainability found in the literature are not consistent between different authors. Some authors use them interchangeably, while others make a distinction between them [139]. Within the group of researchers that make a distinction, the understanding of these terms also varies. Hence, we aim at establishing a working definition to be used in the following Sections of this document.

In our definition, explainability and interpretability are different but not independent concepts. Whose definition is linked to two main questions: "How?" and "Why?". The former is related to interpretability, which aims at providing a clear understanding of the inner workings of a model and the reasoning process that leads to an output. The latter is the question that, given a model's prediction, an explanation tries to answer [74]. However, it is possible to deduce the second question from the first, because if we fully understand the behaviour of a model, we also know the answer to the "Why?" question. The opposite is not true, because knowing the explanation of a certain prediction does not directly indicate how the model reached that prediction. We further introduce a third question, "Where?", which is frequently associated with explainability, while not providing concrete answers to the "Why?" question.

Due to the ability to answer to the "Why?" question, by answering the "How?" question, it can be assumed that knowing the answer to the first, it might be the case that we are moving towards the answer to the second, even if we do not have enough information to answer it completely. We believe that there is a strong and close relationship between xAI and causality [158].

Causality studies the relationship between causes and effects, and it aims to model the causal relationships that define an event or sequence of events [28]. Similarly, explainable artificial intelligence aims to understand the model inner mechanisms that led to a certain output when conditioned on a certain input. Pearl [158, 159] proposed the ladder of causality (Fig. 3). Within this framework, it is possible to characterise causality in three levels: (1) association, (2) intervention and (3) counterfactuals. Each of these levels aims to answer a certain type of questions, and they follow a hierarchical structure, in the sense that for answering the questions of a level $i$, we must have the information to answer the questions of all previous levels [158]. Pearl [158] goes beyond defining the degree of causality associated with each level, he further establishes the mathematical notation for each level. For the Association level (1) the usual representation is a conditional probability in the form of $P(y|x) = p$, which calculates the probability $p$ of the event $Y = y$ given that the event $X = x$ was observed. Bayesian networks are effective and efficient at calculating these conditional probabilities [157]. Intervention level (2) operations can be described, leveraging the *do* operator, as $P(y|do(x), z)$, which calculates the probability of the event $Y = y$, given that we intervene to set $X$ to $x$, and subsequently observed the event
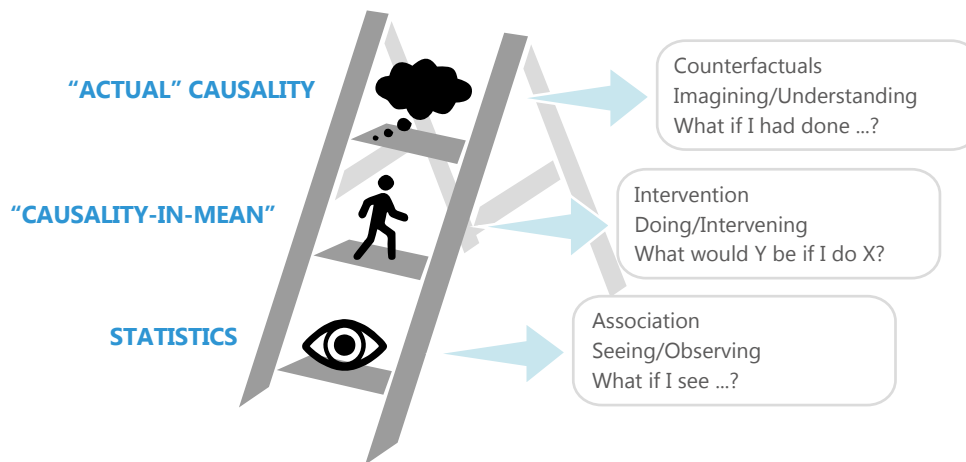
Fig. 3. The Ladder of Causality (adapted from [159]).

$Z = z$. The third and last level, the Counterfactuals, can be represented by $P(y_x|x', y')$, which states the probability $P(y|x)$ when we actually observed the events $X = x'$ and $Y = y'$. It is only possible to compute these if one has access to functional or Structural Equation models, or their properties.

The connection between these two scientific fields can be used to re-frame xAI, and we suggest this change in the following sections.

## 3.2 A brief introduction to xAI

Some of the terms used to define explainable artificial intelligence methods have overlapping definitions and concepts. Some authors aim at dividing them into pre-, in-, and post-model [105], others in methods specific and agnostic to the deep model, and even categorizing them as post hoc or not.

By dividing the interpretability into three smaller subtopics, researchers allow for the creation of different goals for each task. The pre-model analysis focuses on data knowledge and understanding. It significantly leverages visualization techniques that allow researchers to understand the distributions of the input samples [213]. This is an approach frequently used in Business Intelligence (BI) use cases. Traditionally, this technique led to the improved capability of creating handcrafted features for simpler models; nowadays, it is quite relevant to detect some biases in the data, which can be described through the analysis of the skewness in the dataset [225]. Understanding the data used for training allows for increased confidence in the posterior decisions and explanations.

On the other hand, there are also techniques that can be applied to the model and its outputs, during and after training. When the training process is conditioned by some constraint, these are known as in-model techniques and they focus on the direct integration of interpretability into the model through the usage of these conditioning priors. Some families of machine learning models already include inherently/intrinsically interpretable mechanisms in their design. Hence, their use is also considered to be an in-model approach. Some of these interpretable models are rule based [27, 175, 219], per-feature based [66, 84] and some simple versions of linear models such as linear regressors. They are, however, limited by the semantic meaning of the original features and the size/complexity/depth of the model. Nonetheless, it is possible to increase, to a certain degree, the interpretability of a complex model. Constraints to these
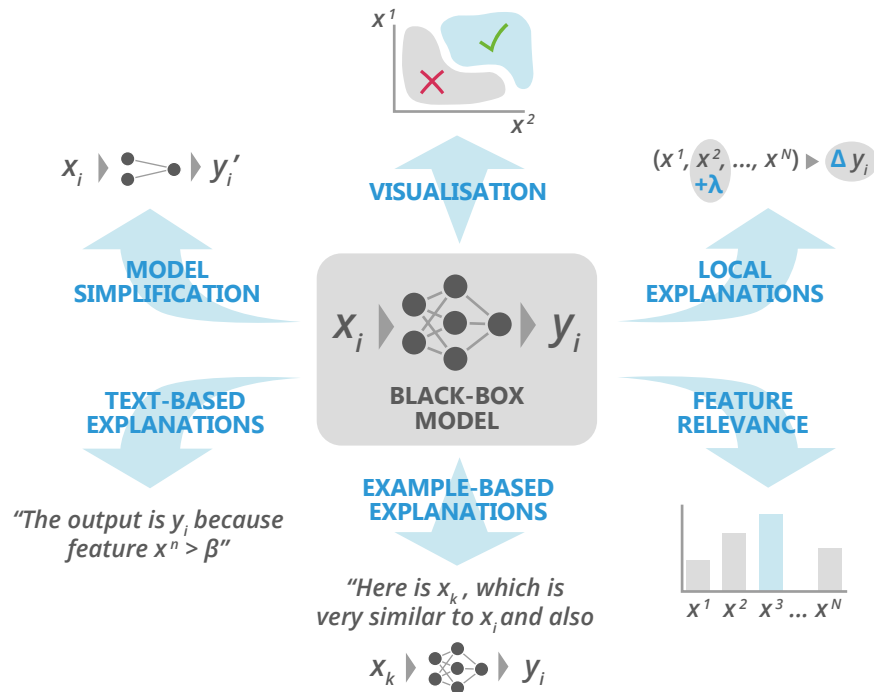
Fig. 4. How we can explain black box models: a summary of current *post hoc* explainability approaches (adapted from [17]).

models, frequently expressed in the form of regularisation terms (such as L1 regularisation [76]), are the most common approach to increase the interpretability of those models. Other approaches include constraining the model to have certain properties, such as monotonicity [190]. These constraints can, sometimes, be applied or enforced in specific modules. While, it is known that achieving global holistic interpretability is an extremely difficult task, analysing the model in modules for global modular interpretability is much more manageable [138]. The exploration of these in-model approaches within the context of biometric systems is rather limited due to the potential sacrifice of predictive power. Instead, the focus has been on *post hoc* methods, which is a different term to post-model approaches. It is, however, important to note that, if the constraints applied to the model represent meaningful priors, it can be the case that a better solution can be found, since the search space is reduced to exclude worse solutions [150, 227].

Aligned with the efforts towards continuing the development of more accurate models, the focus has been partially redirected towards post-model interpretability. This approach is frequently known as explainability since these methods generate explanations to a given prediction of a trained model. While it is more common to have model-specific in-model approaches and model-agnostic post-model approaches, this not always the case. For instance, gradient-based methods, which leverage gradient information to identify areas of the image that contribute to the final decision, are specific to deep neural networks [15, 193, 196, 197]. One of the common approaches to post hoc interpretability adds perturbations to the input in order to assess the impact on the output [65]. As discussed in the following Sections, these methods do not have a rich semantic meaning. Hence, to enrich the explanations produced, some methods map the latent features of a model to the input space [56, 240] or to a representation of semantic concepts [18].

Lately, there has been research on a post-model approach that closely resembles ideas from in-model methods. While training a simpler model is somewhat complicated if the task is too complex, it is possible to teach a simpler model to mimic the behaviour of the more complex (deeper) model [13]. Due to the lower parameter count, simpler models present a more interpretable view of the decisions made by the complex model. This is known as model distillation. This and other compression techniques, such as compression [23] and pruning [116] can help to reduce the complexity of an overly complex neural network architecture [33].

In biometrics, the prevalence of pre- and in-model approaches is rather limited when compared to post-model methods. Due to the high requirements for accuracy and other performance metrics, the usage of post hoc methods is much more convenient, requires less engineering effort, and provides an acceptable trade-off between accuracy and these efforts. Nonetheless, while it provides immediate convenience, it is not closer to a proper path to xAI than the other two approaches. Moreover, the current taxonomy is limited as it does not explore the real meaning of understanding a specific model. Rather, current approaches focus on explanation methods that focus on understanding the prediction, but not the model. Furthermore, this taxonomy fails to include several in-model approaches that improve the interpretability of the system through the introduction of domain knowledge within the training process. As such, the following Section proposes a revised taxonomy for Explainable AI and how it encapsulates the methods seen in Figure 4.

### 3.3 The causal relationships of a deep model

Causality's main goal is to understand and learn the causal model of the *world*. Once this model is found, it is possible to make predictions and informed guesses with a tremendous degree of certainty. It is possible to grasp the implication of the smallest details on the deviations of a certain event. As previously mentioned, the ladder of causality represents a simplification of the steps one must take towards getting a causal model of the world. In practice, the last step on the ladder represents a mental representation of the *world* that allows humans to replay, recreate and even imagine new events on their heads. One can formulate the question of "What would happen if I stop taking my antibiotics now?", and humans are capable of answering through imagination.
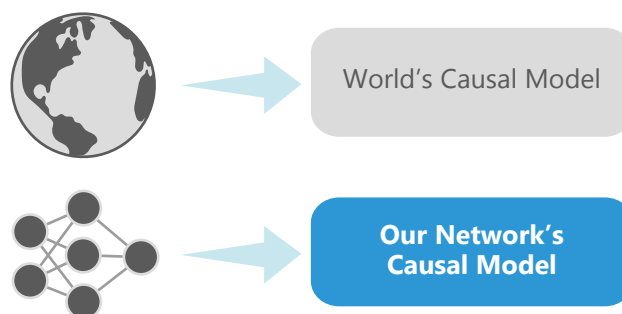


Fig. 5. A figure illustrating the World's Causal Model which represents all the causal relationships that we observe and rule the functioning of the world. Below we see a representation of the Causal Model of a Deep Neural Network. While the latter represents the causal relationships within the "rules" learnt by the DNN, the former represents the true causal relationships of the world. These two models are very likely to disagree, and the latter does not aim to explain the causal relationships of the world, just the specific deep learning model it represents.

Explainable artificial intelligence shares the same goal as causality approaches, in the sense that it aims to understand the model and provide explanations for the relation between two or more events (e.g. input and prediction). This means

that, in a sense, we can try to frame xAI within the context of causality. And, for that, it is necessary to answer the following question "Which causal relationships should our causal model learn?". Since we are only trying to understand a specific model, our causal system does not need to represent the *world*. Instead, we want to find a causal system that models the causal relations of our deep model, as show in Figure 5.

When compared to the previous approach, there are significantly more advantages to this approach. First, it is possible to measure the transition between the different steps of causality. It is also possible to quickly detect spurious connections learnt by our model, through a comparison between its causal model with reality. For instance, if a model uses the skin tone for a task that should not use that specific feature, the causal model of the deep model shall provide this information. Hence, hidden biases can be detected and explained in a straightforward manner. In other words, we can look for misalignments between the causal relationships of *reality* and the causal relationships of a deep model. So, in a biased model, we would find, for instance, a causal relationship between gender and lower wages, and through a careful comparison with the *world* causal relationships, we can deem this specific relationship as biased.
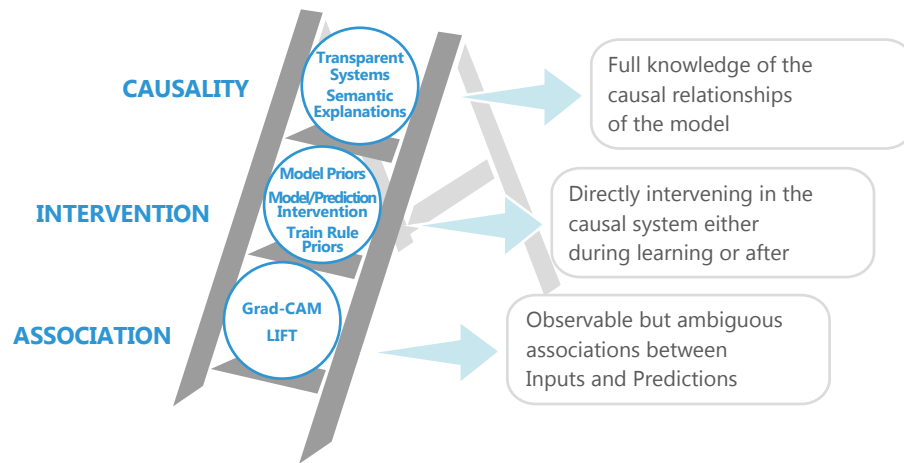


Fig. 6. xAI Ladder: causality ladder repurposed to categorise different types of xAI approaches according to the understanding of the model provided by them.

For this we revisited the ladder of causality (Fig. 3) and propose a similar structure for the ladder of explainable artificial intelligence as seen in Figure 6. This structure for categorising different xAI approaches shows how they contribute to answering the questions that we introduced in Section 3.1. Furthermore, it highlights the importance of each technique and the journey ahead. On the first step of the ladder, we can see observational explanations, such as visualisation maps, which are quite powerful at locating elements that were relevant for the decision, but that do not hold semantic value or any information regarding what was in fact important for the decision. This presents itself closely linked with the associations at the first step of the ladder of causality. Similarly, pre-model approaches, that rely on the understanding of the data, are unaware of the inner mechanisms of the model, hence, these approaches cannot, by definition, go beyond the first step on the ladder.

The second step on the xAI ladder encompasses some of the more advanced approaches currently used in the literature. Furthermore, it is the case that an approach might be between two different rungs. For instance, one might insert prior causal knowledge to condition the training process. Thus, it is possible to know the inner workings of the

model regarding that prior, which means that it is possible to reduce the set of possible explanations while partially knowing a subset of the causal relationships of our deep network. As such, the rung climbing is not a discrete process, but instead a continuous one. Some of the most common methods on the second rung are case-based explanations, semantic explanations, LIME (which closely resembles the $do(.)$ operation [158] by intervening in the input and analysing its effects. But in fact, the majority of applications is in between the first and second rung, where they can perfectly answer the "Where?" question regarding the importance of the variables, but cannot provide a strong and complete answer the "Why?" question. Climbing from the second to the third rung is the final goal of researchers working on xAI. This last rung requires a complete understanding of the inner causal relationships of the deep network. When this happens, one can fully answer the "How?" question, and can further answer to the "Why?" and "Where?" questions. Moreover, knowing these relationships means that producing counterfactuals is not a difficult task, for instance, for a skin-tone biased model we can quickly verify that a change in the skin tone would impact the final decision. And for this, we are not required to provide a new empirical sample to the model so we can empirically check the result (as happens with LIME for instance).

The implications, nuances and adaptations of this new framework to categorise are vast, and require a careful analysis and continual study. As such, we provide in the following sections a few notes on each rung.

*3.3.1   The First Rung: Observational methods as a first order explanation.* As with the association within the causal domain, in xAI, visualisation approaches are the most common. This tendency can be explained by the fact that most of the visualisation approaches require little computational effort, are less dependent on the model architecture and work well for several tasks. This new definition categorises visualisations as first order explanations, meaning that these are the most immediate explanations one can generate to quickly inspect the reasons behind a prediction. Furthermore, they provide less information than second and third order explanations.

One of the first strategies for the visualization of deep learning models was proposed by Zeiler *et al.* [240]. This was mainly motivated by the scarce insights that researchers had about the functioning and behavior of the models that led them to achieve high performances. In their approach, the authors used a multi-layered deconvolutional network [241] to project the feature activations into the input space. Besides, they also proposed several occlusion tests to evaluate the sensitivity of the classifier. They occluded portions of the input image and revealed the feature activations, thus showing the parts of the inputs that were relevant to the classification. Later, Simonyan *et al.* [193] proposed two techniques that focus on the respective class of the input to generate the feature activations. This is achieved through the computation of the gradient of the class score concerning the input image. In the first use case, the authors built their work on top of Erhan *et al.* [61]. Given a score, and a learned model, the idea is to generate an image that maximizes the score, outputted by the model. This synthetic image is thus representative of the class appearance learned by the model [193]. The second use case applies a similar strategy to generate a class *saliency map*. In this approach, the idea is to generate the saliency map that maximizes the class of a given input image. The result of this method is a mask that has pixels with high intensities in the locations that are related to that class.

Springenberg *et al.* [197] revisited the work behind deconvolution and performed experiments with a model composed solely of convolutional layers. Known as *guided backpropagation*, this approach proposed a modified way to handle backpropagation through rectified linear (ReLU) nonlinearity: instead of removing the negative values corresponding to negative entries of the top gradient (*i.e.*, deconvolution) or bottom data (*i.e.*, backpropagation), they remove the values for which at least one of these entries is negative. According to the authors, this technique prevents the backward flow of negative gradients, which often correspond to the neurons which decrease the activation of the higher layer unit one

aims to visualize. Sundararajan *et al.* [199] argued that most attribution methods did not respect two fundamental axioms (*sensitivity* and *implementation invariance*), thus representing an important flaw of such methods. Taking this premise into account, the authors proposed a method called *integrated gradients*, which does not require any modification to the original network architecture, while being related only to the gradient operator. Following a similar research line, Shrikumar *et al.* [189] created DeepLIFT. This method is based on the backpropagation of the contributions of all the neurons in the network to all the features of the input and works by comparing the activation of each neuron to a given *reference activation* and assigning contribution scores according to this difference. DeepLIFT also allows practitioners to give separate consideration to positive and negative contributions, hence allowing them to find other dependencies.

Bach *et al.* [14] studied the influence of pixel-level contributions through the introduction of the *layer-wise relevance propagation* (LRP) technique which considers that a deep model can be divided into several layers of computation which may be responsible for extracting features or performing classification. Following the work of Zhou *et al.* [245] on *class activation mapping* (CAM), Selvaraju *et al.* [183] developed the *gradient-weighted class activation mapping* (Grad-CAM), which allows the practitioner to generate saliency maps with respect to a specific class (*i.e.*, class-discriminative), a property that may help researchers unveiling some visual concepts related to a given class.

Frequently, these visualisations have not been seen as explanation methods, instead, they have been seen as a tool that can be used to provide additional information to the end-user [39, 180]. In a sense, it can be seen as a parallel of distillation methods frequently applied to complex models, but in this case, used on complex high-order explanations. Even more than traditional explanations, the interpretation of outputs from a visualisation method is highly subjective and user-dependent, despite some attempts to objectively evaluate them [179]. The former aligns with the perspective of assuming visualisations as approaches to achieve first order explainability, whereas the second cements it through another parallelism between visualisations and associations, which are also subject to very different interpretations. Hence, currently, saliency maps and other methods of feature visualisation are a bridge between the limited perception of humans on any number of dimensions above two or three, and the complete user-specific semantic explanations.

The practicality of the visualisations as a tool for first order explanations led to a wide-spread implementation of this. And while researchers agree that using visualisation techniques has a positive effect on our understanding of the deep model, the amount of information contained within these explanations is rather limited and vague. For instance, we remove crucial information from these explanations, such as the importance of colour, transmitting a more open-ended problem to be discussed among humans. One of the most consensual statements between experts in explainable artificial intelligence is that humans can better understand counterfactual explanations [31, 211]. In other words, a counter-example of the one given as input, which introduces a notion of contrastiveness [117]. This behaviour is representative of the relationship between counterfactual examples and causality [30, 224]. In other words, to move towards more complete explanations it is necessary to keep climbing the ladder.

Despite their limitations and their current use, it is expected that in the near future explanation methods will keep on leveraging methods from the first rung. Their practicality is difficult to defeat, and several tasks and problems are not concerned with higher-order explainability, since they only focus on small assessments of the model's behaviour.

It is worth noting that while this first rung considers both post-model approaches (e.g. Grad-CAM) and pre-model approaches (e.g. data visualisation), these can still be differentiated as some of them already have some information regarding the model's architecture and gradient propagation. Hence, it can be said that while not providing an answer to the "Why?" question, they provide stronger insights into the behaviour of the model. The following sub-section keeps on climbing the ladder.

*3.3.2 The Second Rung: Higher Order Explanation with Priors and Semantic Information.* The climb from the first to the second rung presents itself as a task with a diverse set of pathways. Differently from the direct explanation presented in previous chapters regarding interventions, this operation can take many forms when applied to artificial intelligence. In a sense, we can describe any of the operations that will be discussed as an intervention to the deep neural network that we want to describe in causal terms. Nonetheless, there are fundamental differences between them. Moreover, some of these interventions require careful analysis, as confounding bias is a threat that remains present [195]. In contrast with the previous rung, the methods in this rung require their user to be, at least, partially aware of the causal relationships of the problem. Moreover, we further differentiate between the different elements in an AI system. Our overview of the system's causal structure can be seen in Figure 7. It is important to distinguish between the inferred label and the true label. The former is defined by an human annotator and is inferred from the input itself. Moreover, it has a direct impact on the model. The latter is the true label, it is not always available and is the cause of the input, as discussed in medicine, the disease causes the image alternation and not the opposite [35]. Hence, the latter has an indirect effect on the former too.
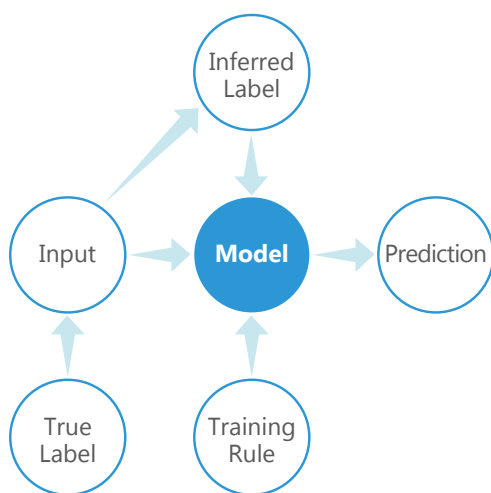


Fig. 7. Causal model representation of the training. The input, the label and the training rules directly affect the model, whilst the model affects the final prediction. The Inferred Label and the True Label are also separate, while the former is given to the model and inferred from the input, the latter is the cause of the input and is not always available.



Fig. 8. Causal model representation of the inference process. The Model node is a causal system by itself and it is the system that we are trying to understand.

In the second rung we can have both post-hoc and in-model methods. Within each of these two groups it is also possible to find distinct strategies to do interventions, or to partially find a causal representation of a deep learning model. For the post-hoc interventions, we highlight interventions in the input and interventions in the prediction.

- *Intervention in the Input:* As shown in Figure 7, the input given to the model has a direct arrow pointing towards the model. In other words, there is a causal effect between the input and the model, which means that changing the input will have consequences on the model's behaviour. It is possible that, by intervening on the

Inputs by modifying them to the desired values, we can observe and learn the causal relationships between the model and the prediction. This type of interventions allows for the interpretability assessment to occur after training. It can be particularly useful if one aims to detect the implications of skin tone on the final prediction. A transformation in the skin tone of the image might reveal some hidden causal relationships between the model and the prediction, that should not exist. The relevance of this type of interventions is growing with the current developments of generative artificial intelligence methods, such as Diffusion Models [161, 243].

- *Intervention in the Prediction:* Explaining the model is the main goal of xAI, nonetheless, in some scenarios explaining the decision might represent a good proxy. As such, one of the most common strategies to explain these predictions relies on finding similar samples in the input space that produced similar predictions. This behaviour might help the human to focus on the common areas of both inputs. The most relevant implementations of these interventions are usually related to medical images [140]. An extension of these case-based explanations can be used as a proxy to find counterfactual samples. Setting the prediction as $P$, it is possible to find the nearest input, $X$, with a different value for $P$. The difference between both inputs highlights some reasons that distinguish both classes.

On the other hand, for the in-model methods, we highlight the introduction of causal information in the model. This can be achieved through training or strong model constraints. The main advantage of these approaches is that they provide partial control of the model, which translates into an *expected behaviour* that can be explained through the causal "rules" directly given by researchers. In other words, as seen in Figure 8 acting on the Training Rule affects the model in a manner that can be predicted, and acting on the model also leads to predictable behaviour. This means that both approaches reveal a part of the causal model that will represent the final version of our system.

- *Model Priors:* Deep neural network architectures are, in the majority of the cases, shared across applications and tasks. This is possible due to their high expressiveness and capability to represent a myriad of complex functions. Nonetheless, this representativeness reinforces the opaqueness of the system. However, it is possible to insert certain decision choices in the model that reduce the function space, and lead to a predictable behaviour under certain circumstances. An example is the rotation equivariance on certain filters of the network [36, 77]. Another approach aligned with this strategy is the ProtoPNet [40, 54] which learns prototypes that can be used afterwards to explain model predictions.

- *Training Rule Priors:* In a similar fashion, these priors reveal certain parts of the causal system that represents the learnt deep neural network. However, these constraints are usually more relaxed than when applied to the model itself. Some examples are orthogonality constraints through the minimisation of the inner product of two vectors [150], occlusion robust losses that promote similar embeddings between occluded and non-occluded samples [146]. Due to their relaxed nature, these priors are seen as an approach to promote a desired behaviour, usually a behaviour data researchers already know to be true. For that reason, and since they do not completely enforce the constraint, the interpretability of the model is usually lower than with Model Priors.

The range of approaches that fit the second rung of the causal explainability ladder is wide and might benefit from one or several of the previously described strategies. For instance, it is possible to have strong model priors, and leverage them through interventions on the predictions to construct a naive counterfactual. A naive counterfactual can be achieved if a system is capable of providing to a user a close and plausible alternative to the original input that belongs to a different class. It is frequently achieved through a search on the latent space [104]. The term naive highlights the ignorance related to the inner relations of the model, and a focus on searching on the latent space. Some methods

promote a certain structure on the latent space, and thus, this search can be somewhat controlled for certain variables. This introduces the capability to have more human friendly explanations. However, this is not easy to devise in practice due to humans' limited range of perception and control over the majority of the variables [125]. Thus, what a human searches for in a counterfactual explanation might not be fully represented in the provided explanation, leading to a trial and error search.

*3.3.3 The Third Rung: Imagining/Extrapolating.* Constructing and understanding the causal model that represents a deep neural network is not trivial. At the time of writing this article, there was no approach capable of fully explaining a deep neural network. As such, the third rung is yet to be achieved. The majority of the efforts have been devoted to the first two rungs, as such explanations suffice for the vast majority of tasks. Nonetheless, in critical domains, such as biometrics and healthcare applications, a incomplete knowledge of a deep learning model might prove itself to be a challenge in terms of security and user acceptance.

The third rung of the explainability ladder is the last rung in the latter, which means that explanations provided by these methods are the highest possible degree of explanation. One very relevant example to compare this rung with the second rung is the construction of counterfactual samples. Differently from a naive counterfactual, a true counterfactual (or just counterfactual) originates from the complete understanding of the models' inner working mechanisms. Instead of searching on the latent space, one can wonder "What would be my prediction had my input image been given with a different head pose?". Moreover, this highlights the capabilities for detection and mitigation of biases, through simple queries, such as "Had the ethnicity been different, would it represent a worse prediction?". And we do not need to directly compute these values, as we can directly try and find the relationship between skin tone and certain predictions.

Locating problems within a deep learning model is a particularly difficult task. However, when the causal model of that neural network is known, it is possible to quickly debug it and look for unwanted/spurious relationships through the comparison of that causal model and a partial model of the "world". And while some relationships might be spurious in real life, they can represent causal knowledge to the learnt deep neural network. Hence, highlighting information wrongly acquired by that same network.

## 3.4 Work in xAI for Biometrics

Recently, the number of applications of methods that merge techniques of explainable artificial intelligence and biometrics has grown. As such, biometrics is a great case study for analyzing these techniques and how they fit within the proposed taxonomy. In this review, we analyzed 81 different papers published between 2018 and 2024. These studies have focused on distinct modalities and applications of biometrics with varying degrees of complexity. The curation of the paper list follows two main steps. First, a comprehensive keyword-guided search was conducted on papers indexed on Google Scholar. The keywords used are related to biometrics and explainable artificial intelligence. Subsequently, we manually scrapped the papers published at the main Computer Vision and Biometrics conferences, which were selected from the Google Scholar conference ranking[2]. From the selected research documents, 53.1% were published between 2022 and 2023, highlighting progress and interest in the field. In comparison, only 37.0% of these studies were published between 2020 and 2021. A small proportion (approximately 9 %) was published between 2018 and 2019. As shown in Table 1, the face is the most frequent modality (76.5% of the papers), with applications in recognition, morphing, presentation, and deepfake attack detection. Owing to its easier acquisition and familiarity, the face is an obvious target

---

[2]https://scholar.google.com/citations?view_op=top_venues&hl=en&vq=eng_computervisionpatternrecognition

for explainability research. Iris recognition and presentation attack detection have also had several studies focusing on bringing explainability to these methods.

Nearly half of the studies analyzed did not include any mechanism to promote interpretability in their designs. They rely exclusively on explanation-based methods that are used in the vast majority of situations to produce visualizations and observational information of a potential explanation for a prediction. Thus, 37 of the 81 papers were placed on the first rung of the xAI ladder. For instance, Trokielewicz *et al.* published studies on recognition [206] and presentation attack detection [205] on cadaver irises. Both approaches rely exclusively on Grad-CAM to produce visualizations, similar to Hu *et al.* [87] for regular iris recognition. Nonetheless, not all visualizations are equal, and the approach used is an important element. Fu *et al.* [70] introduced differential activation maps that were tweaked to explain demographic biases from the perspective of face recognition systems. The activation maps were grouped according to the ethnicity or gender of the individual, and the mean and variation information about them was extracted. These were used to construct a final map that illustrates the deviation from the Caucasian class.

Although extensively used for two-dimensional inputs, visualization methods are not restricted and can be applied to different inputs. For instance, Pinto and Cardoso [165] studied different components of an ECG signal to determine the most relevant component for a recognition task. Four distinct visualization methods were used to dissect ECG signals. Similarly, Lim *et al.* [115] proposed the use of visualization techniques, such as Deep Taylor, to study the differences in rhythm and pitch between DeepFake and human voices, and LRP to uncover characteristics not seen through human visual perception. Deepfake applications rely mostly on approaches that belong to the first rung. Mazaheri and Roy-Chowdhury [128] proposed the use of CAM-based methods to detect inconsistencies in a manipulated image. Their method achieved state-of-the-art results for both DeepFake and face-manipulation detections. In the attack detection domain, Seibold *et al.* [182] and Xu *et al.* [229] presented methods for detecting morphing and deep fake attacks, respectively. The former proposed Focused Layer-wise Relevance Propagation (FLRP), which is an extension of Layer-wise Relevance Propagation (LRP), to explain the predictions of the network to a human. This new method focuses on intermediate layers to generate a visualization map so that it can capture the artefacts that characterize morphing attacks. Its use has proven to be particularly beneficial in situations where the backbone neural network displays uncertainty. The novel contrastive solution proposed by Xu *et al.* [229] for deep fake detection can be generalized for distinct generators of attacks. More importantly, through the use of heat maps and Uniform Manifold Approximation and Projections (UMAP), it is possible to provide a simple visualization of the model decision. Genovese *et al.* [73] introduced a method to simulate facial ageing. While the proposed generative method does not include any interpretability element, the authors proposed a novel method to explain the behaviour of generative adversarial (GAN) models. The method known as the Cross-GAN Filter Similarity Index (CGFSI) uses extracted filters from two generators, classifies them as belonging to the same or different categories and finally gives a similarity score for both generators. Sharma and Ross [188] tackled the problem of iris presentation attack detection (PAD). They proposed the integration of an Optical Coherence Tomograph (OCT) within the systems that aim to detect presentation attacks. Besides the promising results, the performance of the model was studied from the explainable artificial intelligence perspective. The t-distributed stochastic neighbour embedding (t-SNE) and the Grad-CAM algorithm are used to create a visualization of the prediction in a smaller latent space and to generate a heatmap to spatially highlight the most relevant areas for the prediction.

In contrast to CAM-based methods, LIME-based approaches intervene in the input through slight variations. These interventions allow for an effective assessment of the importance of certain features and how local changes in these features affect the final performance of the system. Although not as effective in images as it is in explaining interpretable features, Rajpal *et al.* [171] proposed an LIME-based explanation for face recognition systems. This approach was tested

on four distinct neural network architectures and three different datasets. However, the proposed methodology is not comparable to the current literature because of the chosen architectures and experimental design. Both Bousnina *et al.* [21] and Lu and Ebrahimi [124] explored RISE-based [162] to explain face verification algorithms. The former generates randomly masked samples from the probe image and compares all samples to the reference image. Based on these comparisons, similar or dissimilar regions were highlighted for positive and negative matches, respectively. In the latter work, Lu and Ebrahimi [124] proposed a variation of Rise called S-RISE. Their approach consisted of comparing an unmasked probe with a masked reference image and a masked probe with an unmasked reference image. Based on these comparisons, it was possible to generate explanations that were evaluated quantitatively through insertion and deletion processes. Huber *et al.* [92] expanded on patch-replacement approaches as a form of intervention on the input. Their work focused on distinguishing between the contributions to the decision of the different features in the feature vector, and based on those individual contributions, it was possible to generate explanations that aligned with the final prediction. For positive matches, only similar features between the two images matter to find an explanation, whereas in negative cases, only dissimilar features matter. They further use this information and propose a new dataset called Patch-LFW, which allows them to assess explanations through the questions "What would be my explanation if I swapped this part of the face?". A similar patch-replacement strategy was studied by Knoche *et al.* [107], who also proposed an intervention on the prediction approach that was further expanded by Neto *et al.* [153]. The latter approach generates an interpretable and probabilistic score for the final decision of a verification system based on genuine and impostor distributions.

Interventions on the input have also been shown in fingerphoto verification and post-mortem iris recognition. Ramachandra and Li [172] tackled the complex task of fingerprint recognition through a mobile camera photo, and additionally, to the state-of-the-art performance, this study included several visualization techniques, such as Grad-CAM and gradient attribution, as well as interventions on the input such as LIME and Occlusion Sensitivity [240]. Boyd *et al.* [26] expanded the explainability approaches to post-mortem iris recognition proposed by Kuehlkamp *et al.* [110] that is aware of the areas that are altered by the lack of blood flow. . Their new proposal included human annotators that indicated which areas of the two distinct iris images matched. Based on these annotations and features, the final algorithm was able to determine which areas were a strong match between two images, in accordance with human beliefs. Leveraging this information allows this method to be perfectly aligned with the expectations of the final users, which are forensic examiners.

While the majority of the previous approaches do not necessarily reflect a change in the training process of the deep neural network, some other methods have focused on introducing regularization, priors, and assumption constraints to their learning rule. Specifically, in Face Recognition, Neto *et al.* [146, 148] tackled the problem of masked face verification through the incorporation of additional terms in the loss that represent some domain knowledge previously known by the researcher. In this case, acknowledging that the mask is not relevant to the final prediction, the losses constrain the model output feature vector to be invariant to a potential mask. Huber [88] presented a similar approach for designing a knowledge-distillation framework. Similarly, Franco *et al.* presented two studies [67, 68] that targeted the mitigation of biases within the face recognition system. Following the Demographic Parity constraint [34] the authors designed several regularization terms based on the Tikhonov/Rigde regularizer that were initially introduced by ONeto *et al.* [154]. Considering that a face has two distinct components, Han *et al.* [82] developed a training rule in which a face generates a kernel that is subsequently disentangled into two distinct kernels. The first commonality component represents the shared characteristics among the subjects in the reference set used for optimization. The second component represents the special features of a certain person and is composed of the residuals between the difference of the initial kernel and

the commonality kernel. This method showed a significant improvement over conventional CNNs on the three different datasets. Yin *et al.* [236] presented one of the first approaches for interpretable face recognition. The proposed algorithm was optimized using two novel losses: Feature Activation Diversity (FAD) loss and Spatial Activation Diversity (SAD) loss. The former promotes filter robustness against occlusions, whereas the latter encourages the inclusion of semantic information. Finally, predictions can be analyzed by studying the locations of the filter responses.

Williford *et al.* [222] presented one of the first face recognition approaches that included interventions on both the input and predictions as well as a trained rule prior. They investigated a new form of explanation, based on triplets and an inpainting game. The main goal was to identify regions that were similar between the probe and the positive and dissimilar regions between the probe and the negative. Furthermore, during training, two other approaches to promote interpretability were studied: subtree Excitation Backpropagation (EBP) and Density-based Input Sampling for Explanation (DISE). The former computes the gradients of a triplet for every node of the original network and ranks them into a subtree. These are subsequently used to construct a saliency map that supports these explanations. The latter leverages these outputs to perform prior-based sampling of perturbations in the input image. This is an extension of a previous method known as Randomized Input Sampling for Explanation (RISE) [162], which did not leverage any prior for the sampling. Later, Liu *et al.* [118] addressed the heterogeneous face recognition problem using adversarial training to disentangle different face modalities. Through this approach, their method, called heterogeneous face interpretable disentangled representation (HFIDR), can introduce semantic and interpretable information in the latent space. In this representation, information about identity and information related to modality are explicitly separated. Thus, conditioned on prediction, the author managed to reconstruct the original image in a different modality space. Recently, Terhörst *et al.* [202] used pixel quality indicators to explain face recognition systems. One of the most novel proposals of this approach is the adaptability of the method to different face recognition systems via the propagation of an input image 100 times through stochastic variation of the original network. The authors also proposed inpainting strategies to mitigate the impact of low-quality pixels in a face image. Although not completely interpretable pixel-quality can refer to obstacles in the recognition process, such as occlusions.

Early studies on explainable AI applied to presentation attack detection (PAD) aimed to understand what makes an attack different from a genuine image. While it might be easier to check, given an appropriate clue, if an image is indeed an attack, it is not trivial to define the differences between attacks and genuine images. Jourabloo *et al.* [99] attempted to model spoof information induced by each attack. For this, it was necessary to make certain assumptions regarding the properties of noise, such as its ubiquity. Through this method, it was possible to decompose the original image into noise and facial features. Liu *et al.* [122] approached the facial PAD problem from a similar perspective. Their disentanglement approach was flexible enough to not only detect spoof and live samples through the disentangling of spoof traces but also to reconstruct a live version of that same input or synthesize a new spoof image with the same style and a different reference live image. Wang *et al.* [220] method is another example of this type of approach. By decomposing the input image into a multi-scale frequency image with 12 channels, Fang *et al.* [64] achieved state-of-the-art results in cross-dataset scenarios. This ability to generalize aligns with the fact that the frequency representation is invariant to the dataset and that attacks display distinct frequency patterns. The architecture of the detection network was further designed using a hierarchical attention mechanism to support the detection of these different patterns [75]. Yang *et al.* [234] proposed a method that works with both temporal and spatial features. Specifically on the spatial features, the author splits a single image into several patches of the same image, which are processed by a shared-weight neural network. Subsequently, an attention module weighs the importance of each patch and provides a final classification. This training strategy favours an interpretable notion of the patches that influence the decision. Similarly, Bian *et*

*al.* [19] introduced features that can be interpreted by design. The proposed method was designed to work with both facial and environmental cues. However, the final decision uses a concatenation of interpretable and obscure features, making it difficult to trace the motives for the decision. Pan *et al.* [155] proposed a framework with an integrated attention mechanism and generator for Grad-CAM visualization. The authors proposed the use of visualizations to optimize the training of the attention mechanism to which they designed a set of experimental studies and comparisons. Additionally, their method can generate textual explanations from a set of predefined text options.

To solve Morphing Attack Detection (MAD), and knowing that morphing samples hold two identities, whereas genuine samples hold a single identity, Neto *et al.* [150] designed a deep neural network that outputs two different feature vectors. These vectors have the particularity that their optimization includes a term to minimize the inner product. Thus, orthogonality is promoted and there is no mutual information between them. Simultaneously, they were concatenated and used to make the final prediction using an MLP layer. After identifying some flaws in previous work, Caldeira *et al.* [32] proposed an improved version of the previous method. The orthogonality constraint was relaxed; instead, an autoencoder trained to reconstruct faces was used to create feature vectors. Hence, in the morphing scenarios, the training included two feature vectors of the original images produced by the autoencoder and two feature vectors produced by the proposed network from the morphed image. The new constraint consists of approximating the angle between the projections of each model. Additionally, owing to the intractability of the identity source from a concatenated vector, an interpretable score calculation was proposed. Each vector is used to predict the existence of an identity. When both indicated the presence of a face, the model flagged the input as a morphing sample.

Chen and Ross [42] presented an attention-based approach to tackle iris presentation attack detection. When compared to previous attention-based their approach achieved significant improvements, which are reinforced by Grad-CAM generated visualizations. On the other side, Fang *et al.* [63] developed a patch-based approach assisted by attention mechanisms to solve a similar task. The proposed approach has a loss that works at the pixel level as an auxiliary to the main task of detecting presentation attack inputs. As such, it is also possible to obtain a pixel score or, if intended, a patch score considering a group of pixels. It works especially well when combined with spatial attention blocks, as demonstrated in the Score-CAM visualizations. Joshi *et al.* [98] leveraged the stochastic nature of Monte Carlo dropout to generate uncertainty maps on the segmentation of the fingerprint regions of interest.

In fact, despite the visible growth of xAI applied for biometrics, some important elements are often discarded or ignored. No work explored the effects of causality on the production of explanations, nor the fact that rarely there is only one interpretation for a given explanation/visualization. Moreover, there is no focus on abnormal or rare events that can be extremely useful to produce explanations. The usage of semantic explanations and the analysis of the human-friendliness of the explanations produced are nearly nonexistent. Thus, it is important to look into what is being done in other areas of computer vision regarding xAI, and how that work can be transposed or adapted into biometric applications. Furthermore, the implementation of these methods must be applied across the most diverse biometric traits and tasks. There is still a long path ahead for xAI methods applied to biometrics, and to introduce some consensus and uniformity between the needs of these methods. Despite a categorization of methods in the different rungs, not all the methods in a given rung are equally explainable, and it is important to pay attention to the methods that focus on going beyond simpler strategies, such as LIME.

Table 1. Methods that leverage explainability methods to inspect the predictions of biometric systems. Some of these methods, promote interpretability through the usage of certain techniques, whereas others focus mostly on explaining the predictions. Each method is placed in the respective rung in the xAI ladder. For the methods of the second rung they are further categorised into Model Prior (M. P.), Training Rule Prior (T. R. P.), Intervention on Input (I. I.) and Intervention on the Prediction (I. P.).

| Method | Year | Modality | xAI Ladder | Approach to xAI |
|---|---|---|---|---|
| **Face Recognition** | | | | |
| Huber *et al.* [92] | 2024 | Face Verification | Rung 2 - I. I. + I. P. | Decision-based Patch Replacement and Individual Dimension Contribution to Decision |
| Li [114] | 2023 | Face Recognition | Rung 2 - M. P. | Biologically Inspired Architecture + Attribute Weight for Decision |
| Correia *et al.* [47] | 2023 | Face Verification | Rung 2 - M. P. | Attention + Explainability Masks Construction |
| Bousnina *et al.* [21] | 2023 | Face Verification | Rung 2 - I. I. | RISE-Based Similarity/Dissimilarity Maps |
| Knoche *et al.* [107] | 2023 | Face Verification | Rung 2 - I. I. + I. P. | Patch Replacement |
| Lu and Ebrahimi [124] | 2023 | Face Recognition | Rung 2 - I. I. | S-Rise |
| Rajpal *et al.* [171] | 2023 | Face Recognition | Rung 2 - I. I. | LIME Maps |
| Neto *et al.* [153] | 2023 | Face Recognition | Rung 2 - I. P. | Interpretable Match Score |
| Rodriguez *et al.* [176] | 2023 | Quality Face Recognition | Rung 2 - T. R. P. | Enviromental Attributes |
| Terhörst *et al.* [202] | 2023 | Face quality for Face Recognition | Rung 2 - I. I. + I. P. | Pixel-level Quality Maps |
| Fu *et al.* [70] | 2022 | Face Recognition Bias | Rung 1 | Differential Activation Maps |
| Han *et al.* [82] | 2022 | Face Recognition | Rung 2 - T. R. P. | Personalized Kernels + Feature Maps |
| Roy *et al.* [177] | 2022 | Heterogeneous Face Recognition | Rung 2 - M. P. + I. I. | Invariant Features |
| Winter *et al.* [223] | 2022 | Face Recognition | Rung 2 - T. R. P. + I. P. | Local Patch-based Features with Eexplainable Boosting Machine (EMB) |
| Mery *et al.* [129] | 2022 | Face Verification | Rung 1 | Saliency maps AVG contours |
| Fu *et al.* [69] | 2022 | Face quality for Face recogniton | Rung 1 | Score-CAM[217] + Activation Maps |
| Franco *et al.* [67] | 2022 | Face Recognition | Rung 2 - T. R. P. | Tikhonov regularizers |
| Neto *et al.* [146, 148] | 2021 | Masked Face Recognition | Rung 2 - T. R. P. + I.I. | Mask-Invariant Loss + Grad-CAM Visualisation |
| Huber *et al.* [88] | 2021 | Masked Face Recognition | Rung 2 - T. R. P. + I.I. | Mask-Invariant Knowledge Distillation |
| Franco *et al.* [68] | 2021 | Face Recognition | Rung 2 - T. R. P. | Tikhonov regularizer |
| Liu *et al.* [118] | 2021 | Heterogeneous Face Recognition | Rung 2 - T. R. P. + I. P. | Modality/Identity Disentangled Representation + Cross-modality Synthesis |
| Anghelone *et al.* [9] | 2021 | Face Recognition | Rung 1 | Heatmap-based Visualisation |

Table 1. Methods that leverage explainability methods to inspect the predictions of biometric systems. Some of these methods, promote interpretability through the usage of certain techniques, whereas others focus mostly on explaining the predictions. Each method is placed in the respective rung in the xAI ladder. For the methods of the second rung they are further categorised into Model Prior (M. P.), Training Rule Prior (T. R. P.), Intervention on Input (I. I.) and Intervention on the Prediction (I. P.).

| Method | Year | Modality | xAI Ladder | Approach to xAI |
|---|---|---|---|---|
| Jiang and Zeng [97] | 2021 | Face Recognition | Rung 2 - T. R. P. | Face Components Representation and Self-Attention based Reconstruction + Components Similarity Matrix |
| Williford et al. [222] | 2020 | Face Recognition | Rung 2 - I.I. + T. R. P. + I. P. | DISE + subtree EBP + Inpainting game |
| Zee et al. [239] | 2019 | Face Recognition | Rung 1 | Guided Backpropagation & CAM |
| Yin et al. [236] | 2019 | Face Recognition | Rung 2 - T. R. P. + I. I. | Spatial and Feature Activation Diversity Losses |
| **Face Morphing** | | | | |
| Caldeira et al. [32] | 2023 | Face Morphing Detection | Rung 2 - M. P. + T. R. P. | Interpretable Score + Identity Disentanglement |
| Dargaud et al. [48] | 2023 | Face Morphing Detection | Rung 1 | PCA-based Visualisation |
| Dwivedi et al. [59] | 2023 | Face Morphing Detection | Rung 1 | Ensemble of Visualisations |
| Neto et al. [150] | 2022 | Face Morphing Detection | Rung 2 - T. R. P. | Orthogonal Identity Disentanglement |
| Myhrvold et al. [143] | 2022 | Face Morphing Detection | Rung 1 | LRP |
| Seibold et al. [182] | 2021 | Face Morphing Attack Detection | Rung 1 | Focused LRP |
| **Face PAD** | | | | |
| Prasad et al. [166] | 2023 | Face PAD | Rung 2 - T. R. P. | Intermediate Depth Estimation |
| Pan et al. [155] | 2022 | Face PAD | Rung 2 - T. R. P. | Grad-CAM + Textual Explanations |
| Fang et al. [64] | 2022 | Face PAD | Rung 2 - T. R. P. | Multi-Scale Frequency Decomposition + Hierarchical Attention + t-SNE |
| Bian et al. [19] | 2022 | Face PAD | Rung 2 - T. R. P. | Auxiliary Learning with Interpretable Cues |
| Wang et al. [220] | 2022 | Face PAD | Rung 2 - T. R. P. | Live Features Disentanglement + t-SNE [210] |
| Sequeira et al. [185] | 2021 | Face PAD | Rung 1 | Grad-CAM Visualisation |
| Chen et al. [43] | 2021 | 3D Mask Face PAD | Rung 1 | Grad-CAM Visualisation and t-SNE |
| Neto et al. [152] | 2021 | Face PAD | Rung 1 | Grad-CAM Visualisation |
| Aghdaie et al. [3] | 2021 | Face PAD | Rung 2 - T. R. P. | Attention + Attention Maps Visualization |
| Yu et al. [238] | 2020 | Face PAD | Rung 2 - M. P. + T. R. P. | Central Difference Convolution + Depth-Map Estimation |
| Pinto et al. [163] | 2020 | Face PAD | Rung 1 | Grad-CAM Visualisation |
| Wang et al. [221] | 2020 | Face PAD | Rung 2 - T. R. P. | Temporal Depth-Map Estimation + t-SNE |

Table 1. Methods that leverage explainability methods to inspect the predictions of biometric systems. Some of these methods, promote interpretability through the usage of certain techniques, whereas others focus mostly on explaining the predictions. Each method is placed in the respective rung in the xAI ladder. For the methods of the second rung they are further categorised into Model Prior (M. P.), Training Rule Prior (T. R. P.), Intervention on Input (I. I.) and Intervention on the Prediction (I. P.).

| Method | Year | Modality | xAI Ladder | Approach to xAI |
|---|---|---|---|---|
| Shao *et al.* [186] | 2020 | Face PAD | Rung 2 - T. R. P. | Meta-learning with Depth Regularization + Attention Map Visualization |
| Liu *et al.* [122] | 2020 | Face PAD | Rung 2 - T. R. P. | Spoof-traces disentanglement |
| Yang *et al.* [234] | 2019 | Face PAD | Rung 2 - T. R. P. | Patch-based Attention + Grad-CAM + t-SNE |
| Jourabloo *et al.* [99] | 2018 | Face PAD | Rung 2 - T. R. P. | Spoof Noise Modeling + t-SNE |
| Liu *et al.* [121] | 2018 | Face PAD | Rung 2 - T. R. P. | rPPG Signal Extraction |
| **Face DeepFakes** | | | | |
| Mathews *et al.* [126] | 2023 | Face DeepFake Detection | Rung 1 | Grad-CAM Visualisations |
| Silva *et al.* [191] | 2022 | Face DeepFake Detection | Rung 1 | Attention Maps |
| Mazaheri and Roy-Chowdhury [128] | 2022 | Face Manipulation Detection | Rung 1 | CAM |
| Korshunov *et al.* [108] | 2022 | Face DeepFake Detection | Rung 1 | t-SNE |
| Xu *et al.* [229] | 2022 | Face DeepFake Detection | Rung 1 | UMAP and Heatmap Visualisation |
| **Face Emotion/Expression Recognition** | | | | |
| Nam *et al.* [144] | 2023 | Face Expression Recognition | Rung 2 - T. R. P. | Spatio-Temporal Attention Maps + Facial cues |
| Rathod *et al.* [173] | 2022 | Children Emotion Recognition | Rung 1 | Grad-CAM + Grad-CAM++ + SoftGrad |
| Cesarelli *et al.* [37] | 2022 | Face Emotion Recognition | Rung 1 | Grad-CAM Visualisations |
| Zhu *et al.* [246] | 2022 | Face Emotion Recognition | Rung 1 | LRP |
| Araf *et al.* [12] | 2022 | Face Emotion Recognition | Rung 1 | Grad-CAM Visualisations |
| **Fingerprint PAD/Recognition** | | | | |
| Rai *et al.* [170] | 2023 | Fingerprint PAD | Rung 1 | Heatmap Visualisation |
| Ramachandra and Li [172] | 2023 | Fingerphoto Verification | Rung 2 - I. I. | Grad-CAM + Occlusion Sensitivity Maps + LIME + Gradient Attribution |
| Chowdhury *et al.* [45] | 2020 | Fingerprint Recognition | Rung 1 | Grad-CAM Visualisation |
| **Iris Recognition** | | | | |
| Boyd *et al.* [26] | 2023 | Post-mortem Iris Recognition | Rung 2 - I. I. | Visual Patch-matching |
| Kuehlkamp *et al.* [110] | 2022 | Post-mortem Iris Recognition | Rung 1 | Segmentation Masks + CAM Visualisations |
| Hu *et al.* [87] | 2020 | Iris Recognition | Rung 1 | Grad-CAM Visualisation |
| Trokielewicz *et al.* [206] | 2018 | Cadaver Iris Recognition | Rung 1 | Grad-CAM Visualisation |

Table 1. Methods that leverage explainability methods to inspect the predictions of biometric systems. Some of these methods, promote interpretability through the usage of certain techniques, whereas others focus mostly on explaining the predictions. Each method is placed in the respective rung in the xAI ladder. For the methods of the second rung they are further categorised into Model Prior (M. P.), Training Rule Prior (T. R. P.), Intervention on Input (I. I.) and Intervention on the Prediction (I. P.).

| Method | Year | Modality | xAI Ladder | Approach to xAI |
|---|---|---|---|---|
| **Iris PAD** | | | | |
| Sharma and Ross [188] | 2021 | Iris PAD | Rung 1 | Grad-CAM Visualisation and t-SNE |
| Fang *et al.* [63] | 2021 | Iris PAD | Rung 2 - T. R. P. | Patch-based Supervision + Attention + Score-Weighted CAM [242] |
| Chen and Ross [42] | 2021 | Iris PAD | Rung 2 - T. R. P. | Channel and Position Attention + Grad-CAM |
| Sharma *et al.* [187] | 2020 | Iris PAD | Rung 1 | Grad-CAM Visualisation and t-SNE |
| Trokielewicz *et al.* [205] | 2018 | Cadaver Iris PAD | Rung 1 | Grad-CAM Visualisation and Guided Backpropagation |
| **Others** | | | | |
| Diaz *et al.* [52] | 2023 | Signature Verification | Rung 2 - M. P. + T. P. | Universal Background Model(UBM) + Explainable Features |
| Aquino *et al.* [11] | 2022 | Accelerometer User Identification | Rung 1 | Grad-CAM |
| Lim *et al.* [115] | 2022 | Voice DeepFake Detection | Rung 1 | Deep Taylor + LRP |
| Algermissen and Hörnlein [5] | 2021 | Footstep Person Recognition | Rung 1 | Grad-CAM |
| Alshazly *et al.* [8] | 2021 | Ear Recognition | Rung 1 | Guided Grad-CAM |
| Chandaliya and Nain [38] | 2021 | Face Age Estimation | Rung 2 - T. R. P. | Attention |
| Joshi *et al.* [98] | 2021 | Fingerprint ROI Segmentation | Rung 2 - T. R. P. + I. P. | Monte Carlo Dropout + Uncertainty Map |
| Pinto and Cardoso [165] | 2020 | ECG Biometric Identification | Rung 1 | Gradient-SHAP + DeepLIFT + Saliency Maps |
| Ahmed *et al.* [4] | 2020 | Ethnicity Estimation | Rung 1 | Grad-CAM |
| Genovese *et al.* [73] | 2019 | Face Aging | Rung 1 | Cross- GAN Filter Similarity Index |

### 3.5 Work in xAI computer vision that can benefit biometrics

The scenario of interpretable methods for biometrics has grown and highlights an increasing research interest in this domain. However, other computer vision domains had the opportunity to further develop and deeply explore different approaches to xAI. These approaches, independently of the rung where we place them, are, in most cases, yet to be applied to biometrics. Given its impact on medical image and other critical domains, it is undeniable that they might present themselves as attractive research directions for anyone looking towards improving the explainability of biometric systems. Hence, the following paragraphs delve into some roads already crossed in the computer vision literature. As seen in Table 1, one of the most common training rule constraints on second-rung approaches is the usage of attention mechanisms. These mechanisms have the advantage of being more interpretable, as they provide the importance of each element for their prediction. One of the most commonly used attention-based systems is the Transformer [214]. This architecture led to advances in both computer vision and natural language processing domains. However, it is usually data-hungry, expensive to run on inference and not yet ready for embedded devices. Additionally, the recent usage of transformer-based foundation models further highlights the course of complexity and the opaqueness of these systems as they scale.

In 2019, Chen *et al.* proposed a novel deep neural network, *prototypical part network* (ProtoPNet), that contains an inner reasoning process that may be considered interpretable due to its transparency [41]. This model is trained in an end-to-end fashion and it is optimized to learn *prototypes* (*i.e.*, in this case, prototypes are latent representations) that are specific to a given class. To achieve this, the authors proposed to use generalized convolution by including a *prototype layer* that computes the squared $L_2$ distance instead of the inner product. Moreover, they also proposed to constrain each convolutional filter to be identical to some latent training patch. And thus, allowing us to interpret the convolutional filters as visualizable prototypical image parts. When the model is shown a new image, it dissects the image to find parts that may look like a prototypical part of some class and makes its prediction based on a weighted combination of the similarity scores between parts of the image and the learned prototypes. This approach is fitting to identification and authentication problems in the biometrics domain. Due to the capability to create patch-related explanations, it becomes easier to provide an understandable explanation to the user. Besides, this pipeline could help algorithm developers to assure that their models respect the "right to explanation" [101] in the sense that this reasoning process is transparent to the end-user. Additionally it could be possible to find common and identity-related face elements. On the strongest challenges of these methods is their classification-based design, which makes them inappropriate for open-set tasks, such as face verification. Other works have also explored the use of prototypical networks [83, 133, 145, 194, 247].

Inspired by the Recognition-By-Components cognitive psychology theory proposed by Biederman, Saralajew *et al.* [181] proposed an architecture called Classification-By-Components network (CBC). The proposed network is optimized to learn generic characterization components for object detection. Moreover, it follows a class-wise reasoning approach with three distinct types of reasoning (positive, negative and indefinite) that is learnt to solve the classification task. The combination of the reasoning types branch into a probabilistic classifier and the decomposition of objects into generic components provides a precise interpretation of the decision process. This approach is even shown to be able of explaining the success behind an adversarial attack, which can be particularly useful for critical biometric systems. It could be explored if other forms of attacks to biometric systems can benefit from this approach. Other relevant application could be related to an extension of patch-based iris recognition work [26].

Qi *et al.* [169] proposed a novel Explanation Neural Network (XNN), which maps the output embedding into an explanation space. Their XNN algorithm is called Sparse Reconstruction Autoencoder (SRAE), and the produced

explanation embedding is capable of retaining the prediction power of the original feature embedding. This explanation space can be understood and visualized by humans when employed a visualization approach proposed by the authors. Recently, Utkin *et al.* [208] proposed a similar approach based on autoencoders. Current works in biometrics, such as the one presented by Caldeira *et al.* [32] could leverage this approach to further boost the explainability of the autoencoder in their proposed approach.

In 2021, Booth *et al.* [20] presented a novel model inspection method called BAYES-TREX. Instead of relying on potential abnormal behaviour in the test set, this method finds in-distribution examples with specified prediction confidence. This probabilistic method can reveal misclassifications that have high confidence and display the boundaries between the classes through ambiguous samples. Moreover, it can also predict the behaviour of the network when exposed to a sample from a novel class. The importance of the prediction score in verification tasks is of higher relevance. Some methods present scores that are difficult to interpret, and this type of approaches aligns with the proposals of Neto *et al.* [153] and Knoche *et al.* [107] to create an interpretable matching score.

Aiming at including some semantic meaning into the inspectability methods of deep neural networks, Losch *et al.* [123] proposed Semantic Bottlenecks (SB). These are integrated into pretrained models and align their channel outputs with visual concepts. The proposed method is rather impressive since even with a reduction of the number of channels by two orders of magnitude, the results are still on par with the state-of-the-art for segmentation challenges. SB have two variations, the unsupervised SBs (USB) and the concept-supervised SBs (SSB). Despite the similar performance, the latter shows a higher degree of inspectability.

Another challenge related to the training of deep neural networks is the learning and proper visualization of concepts that are important for a given task. There are several methods that attempt to see inside the hidden layers of DL models, however, they are applied *a posteriori* (*i.e.*, post-hoc analysis), which may lead to misleading conclusions regarding the properties of the latent space [44]. Therefore, to avoid relying on post-hoc analysis of neural networks, Chen *et al.* [44] proposed a mechanism called *concept whitening* (CW). This mechanism applies a whitening operation to the latent space (*i.e.*, decorrelation and normalization) thus forcing the axes of the latent space to be aligned with known concepts of interest. The authors showed that this approach may be a reliable alternative to the use of batch normalization without hurting predictive performance. This methodology may be of interest in the field of PAD since it could open a different path for the comprehension of the concepts related to *bona fide* and attack signals. Moreover, it could be used to align with soft-biometrics attributes, which would help the system to provide a semantic and meaningful explanation based on terms that the end user is familiar with.

Due to being, by nature, more interpretable than more complex models, the interest in simple models is gaining traction. However, training these models frequently leads to worse performance. Hence, researchers have investigated the possibility to change how these models learn. One of the most used approaches, is the distillation of knowledge learned by a complex model (teacher) into a smaller one (student) [6, 120]. It has been already been leveraged in biometrics, not as a promoter of interpretability, but as an auxiliary process for a different task [24, 88]. Hence, even if the training of inherently interpretable models is less stable, it is possible to optimise it by gathering information from a more complex, opaque model.

Other methods that are being explored and can contribute to the increase of interpretability of biometric methods include counterfactual explanations and influence functions. And there is some work in that direction [81, 94, 100, 226]. However, they are still growing and maturing. In their current form, they are similar to the naive counterfactuals described in previous chapters as a second rung approach. It is clear that xAI was already properly introduced into computer vision systems. Nonetheless, when comparing the variability of approaches presented in this Section with

the ones shown in Table 1 that biometric systems are far from being interpretable. Hence, this work also aims at motivating the need to accelerate the research of these methods so that biometric applications can reach the same level of interpretability shown by other computer vision fields.

There has been also attempts to build explainability pipelines capable of handling different types of machine learning models and distinct applications. IBM, presented "The AI 360 Toolkit"[3], which is one of the most advanced pipelines for model and application agnostic explanations. Despite the potential of these systems, they lack the domain information and design that is frequently needed for complete explanations. Hence, while it can be seen as a good starting point, it must not be seen as the definitive solution.

Conducting a careful analysis of the methods presented in this section, it is possible to categorise them accordingly to the xAI ladder proposed in previous sections. As such, it shows the robustness of this new taxonomy across multiple domains. Nonetheless, it is also visible that the field is still far from being saturated and further efforts are needed to achieve the third rung where researchers can understand the causal relations in their models.

## 4  ILLUSTRATING THE BENEFITS OF INTERPRETABLE BIOMETRICS

Interpretable biometrics are still being used exclusively for academic purposes. However, leveraging these techniques can impact lives outside the academic domain and, for instance, support vulnerable groups against discriminatory behaviour. Its integration with the systems that are currently deployed is crucial to the progression towards a more transparent artificial intelligence. Hence, this Section discusses the situations where interpretable biometric systems should be deployed, the benefits that shall arise from their usage, the limitations of the current approaches, the needs for innovation, and finally, it wraps up with some potential research directions.

### 4.1  Applications

The applications of biometric systems are vast, and thus, the applications of their interpretable counterpart are also vast. This raises the question of how it is possible to create a single interpretable system for all the tasks. Realistically, this is a question that shall not be answered for now. Before being solved holistically, interpretable biometrics must be tackled at each smaller subdomain. Progression into a more general approach occurs iteratively through the explanation of the different methods. It can, for instance, explore the validation of an interpretable approach on several tasks or traits. It is however essential to keep in mind that although it would be desirable to be an interpretable system that fits the most variate scenarios, in practice, it is known that there are no free lunches. Moreover, to achieve the third rung it is important to progress towards more complex methods within the second rung instead of focusing on the more simple approaches seen frequently. Besides the benefits previously mentioned, it is possible to discuss the impact of these systems in more areas. For this analysis, interpretable methods for biometrics are deconstructed into three areas where they can be integrated.

*4.1.1  Improving evaluation.* Metrics have been one of the most crucial elements for the raise of deep learning systems. Not only they are responsible for the understanding of how the systems are performing, but they are also used to establish comparisons between different methods. Moreover, measures such as the accuracy and model's confidence have been the golden standard to increase the public acceptance and deployment of these systems. Metrics are useful for and provide a significant portion of information, nonetheless, they also lack significant details.

---

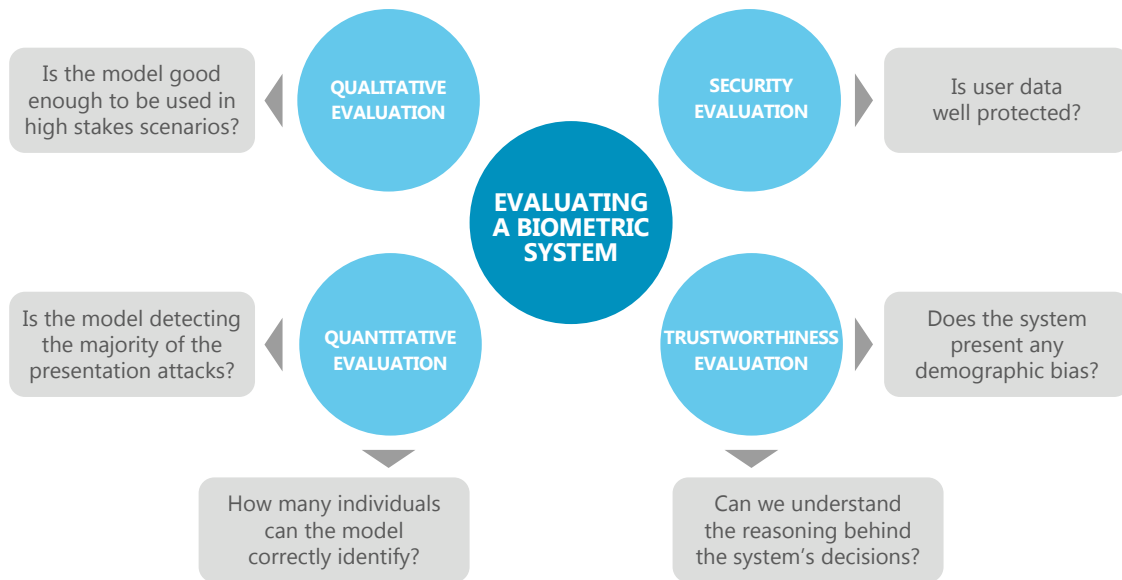[3]https://developer.ibm.com/articles/the-ai-360-toolkit-ai-models-explained/

Fig. 9. Proposed structure for a complete evaluation of the performance of biometric systems at four distinct levels. Quantitative, qualitative, security and trustworthiness evaluation. The latter is focused on the insights provided by explainable artificial intelligence methods.

Understanding how many predictions were correct and by which margin is needed for a quick analysis. However, a correct prediction might be correct for the wrong reasons, and a wrong prediction might be wrong for even worse motives. Hence, providing a single prediction with no explanation does not suffice to fully understand these hidden motives, it is necessary to show why. As such, the evaluation of deep neural network systems is far from complete, leading to flawed assessment of their performance. Leveraging our knowledge of the causal structure, attained through xAI methods, of the model that we have just built holds a higher potential for performance assessment. It allows for two types of evaluation. Retrospective evaluation where after getting a prediction, we try to find the "cause" of that prediction through our knowledge of the model. And also *apriori* evaluation, where researchers can inspect the learnt rules and infer if they are aligned with the expected rules.

We propose an evaluation of deep learning systems that can be divided into four major focal points as displayed in Fig. 9. First, it is essential to quantitatively evaluate our system, and for this we use our typical metrics such as accuracy, F1-Score and others. It is of no use to explain a "dumb" model. Secondly, one should assess in a qualitative manner the performance of its model. In that sense it is important to understand why the model takes the decisions that it takes, how he reaches the decision and which factors are important. This leads to a third evaluation process where we, now fuelled with the knowledge of the "why's" of our system, inspect for bias and trustworthiness problems. It is worth noting that Huber *et al.* [91] has shown that certain Rung 1 approaches suffer from biases themselves. Finally, it is essential to understand the implications of explainability on the privacy of the users and the privacy of the system itself. For instance, Matulionyte [127] explored the implications that having a transparent system has on trade-secrets and how it can be possible to navigate in the trade-off between these two concepts. Additionally, counterfactual explanations have

severe risks if not done properly. It could be possible that by proving such an explanation for a biometric recognition system we compromise the safety and privacy of another user.

*4.1.2 Providing feedback.* Whilst expressing different information, metrics and losses share a lot of similar details. The former, as previously described, supports the evaluation of deep learning models. While the latter provides information for the model optimization. Usually, the minimization of the latter improves the metrics. However, as already mentioned, a correct prediction might rely on incorrect cues. Hence, it is also important to have explainable artificial intelligence methods at the optimization level.
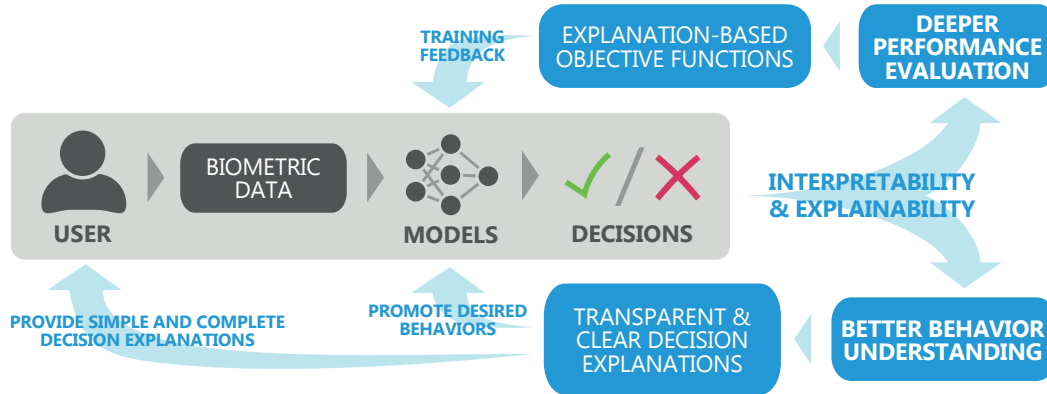


Fig. 10. Structure of a biometric system that leverages explainable artificial intelligence methods to improve its optimization function and to provide useful feedback to the user.

The integration of xAI and optimization processes can be done in two ways, as shown in Fig. 10. First, it is possible to promote certain behaviours or properties through the implementation of constraints to the model. These constraints can be applied to specific layers or to the output of any layer. The second approach is to design objective functions that incorporate feedback from the explanations produced from the predicted batch. This latter approach is significantly more challenging. Nonetheless, there are works aiming to tackle this fusion. For instance, the framework proposed by Pan *et al.* [155] tunes the weights of the attention mechanism through the visual explanations produced in the previous forward pass. Although it is a young field, these strategies have the benefit to directly affect the causal model learnt by the system. And, since we can control these strategies we can promote certain information to be learnt by our system leading to a lower knowledge gap between researcher and deep neural network.

*4.1.3 Protecting privacy.* Some approaches of explainable machine learning focus on image retrieval tasks. For instance, when explaining why two images do not match, it is possible to find examples of mated and non-matted images to create some sort of counterfactual explanation. However, this might expose users' private information. Silva *et al.* [192] leveraged interpretability methods to guide the search for the appropriate image on a similar task on the medical image domain. Further research on the medical domain also focused on the privacy-preserving component of image retrieval [141]. A potential extension for biometric applications was already briefly addressed [142]. We argue that privacy-preserving methods are of extreme importance due to the nature of biometric data. Furthermore, it ensures that users' trust and public opinion, which are crucial for the applicability of biometric systems, remain untouched. Hence, the exposure and potential leakage of important and private information is minimised. The integration of these

mechanisms in real-world situations depends on the desired degree of privacy and the efforts required to change the system.

## 4.2 The incompleteness of the current explanations

Rung 1 approaches such as saliency maps have been widely used to determine which parts of the image are being used by DL-models to perform a decision. However, it is important to acknowledge that knowing where the model is looking within the image does not tell the user what it is doing with that part of the image. A representative example of this is that we can obtain saliency maps for different classes, even if they have no relation to the true label. From all the "possible" explanations given to each class, only the predicted class maps are given to the user. This might create a false sensation of confidence in the decision, which might not be reflected in the overall explanations. For instance, the case where the explanations for multiple (or all) of the classes are identical. Moreover, these methods focus on presenting to the user the meaningful pixels, allowing the user to decide the reason behind the prediction based on those pixels. This raises a problem as humans are independent and subjective animals, and thus, from a simple heatmap $N$ different users can find $N$ different motives for the highlight of those pixels (Fig. 11). In a sense, Grad-CAM and similar methods give a small hint to the user regarding to where the model was looking at. The reason why those were the meaningful pixels remains unanswered and the prediction remains unexplained. Still from Figure 11, the model might have hidden biases against a demographic group, which are not carefully described by the observational explanation.
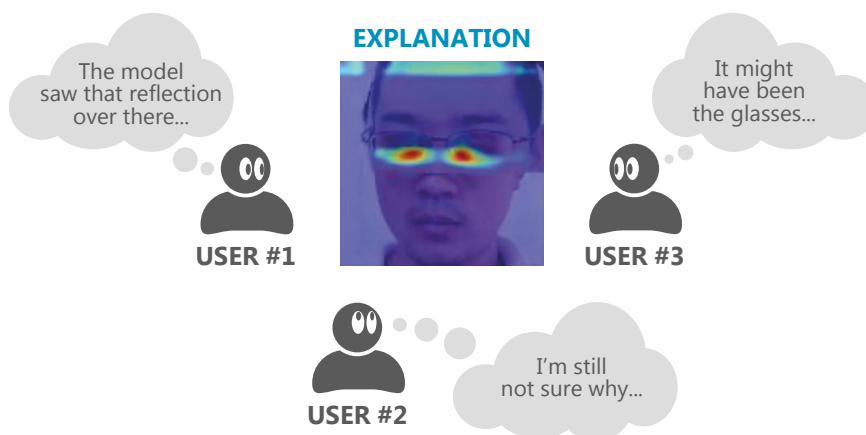


Fig. 11. Simple draft of a diagram exemplifying how Grad-CAM output as heatmap-based images can be interpreted in a different manner by humans.

Several models have been designed to generate textual explanations as part of their training process [74]. This approach has demonstrated interesting results in systems of visual question answering [10] and in algorithms for image classification [85]. To achieve human-understandable explanations, the modules that generate explanations are trained on large data sets of human-written explanations, thus generating models that explain their decisions using natural language [86]. This approach has been extended to multi-modal explanations, in [156]. In this use case, the model has an extra module (based on attention mechanisms) that generates visual explanations. Rio-Torto *et al.* [174] recently explored parameter efficient strategies to generate textual explanations for chest x-rays, leveraging the recently popular transformer architecture. There is, however, an important detail that we should take into account: the explanations

are generated based on the decision, after the decision of the network has already been made. Hence, given that the explanations are learned in a supervised manner and that the generated explanations are conditioned by the output of the network, should we consider these as real explanations? What differs these models from general classification models? How can we be sure that the textual explanation that was generated truly represents the reasoning process used by the model to output a prediction? Additionally, the current concern regarding hallucinations of large language models highlights some potential problems of explaining an obscure system through another obscure system [96, 230].

### 4.3 Unveiling biases

Biometric systems have at their disposal inputs that frequently encode demographic information or sensitive attributes. Not only that, but it has been shown that these attributes are also encoded in the latent space of these models. And while this information is useful for increased performance, one must be careful with the way it is used. On very important topic, especially in face recognition, is the measurement and mitigation of racial bias. This bias is characterised by an intra-ethnicity profile, in other words it relates to the capability of correctly matching the individuals within an ethnic group. However, if these embeddings are used for different tasks, for instance, as a secondary input to a credit approval system, the hidden demographic information can lead to an increased biased inter-ethnicity hindering the necessity to explain the information contained in the latent space.

While they can be measured, it is hard to directly infer the existence of biases simply by analysing the trained model. As such, several models are flagged as biased several days, weeks or even months after being released. This further reinforces the need for explainable AI, and through the proposed taxonomy we can see that the more causal relations of a model we understand, the easier it is to compare them to the causal relations of the *real world*. Through this comparison we can easily find relations that are not aligned and that could represent a source of bias, such as the usage of the skin tone to predict crime re-incidence or credit assessment. Kei *et al.* also argues that explanations can provide much-needed help in this matter [103]. It is important to detect and act quickly on those cases. Frequently edge cases are related to under-represented minorities in the dataset [29]. Moreover, some datasets with specific data from some minority include lower-quality samples [212]. Hence, they are the most exposed to potential harms by algorithmic mistakes. Explanations can serve as a tool to mitigate these nefarious effects and impel machine learning into fairness [55]. Terhörst *et al.* [202] and Boutros *et al.* [25] are examples of these quality assessments that can be used to construct balanced datasets.

With the rise of deep learning-based biometric systems developed with direct applications to police forces and law enforcement, unveiling bias becomes mandatory [2, 7, 71, 168]. Hence, these use cases have a high requirement for appropriate explanations and not incomplete ones. For instance, if a face presentation attack detection algorithm suffers from a demographic bias that is prominent on darker skin tones, it cannot be detected by visualization methods. Nor can it explain that the skin tone played a role in the decision. Moreover, if the explanation was given to a human, he would most likely focus his attention on the shape of the areas around the highlighted pixels, search for blurriness or artifacts. Hence, for the safe deployment of these systems one must attempt to learn the hidden rules that guide the model towards a decision. However, despite not being a complete explanation nor unveiling the hidden mechanisms of a deep neural network, rung 1 approaches are still useful for preliminary verification of the deep learning system.

## 5 CONCLUSION: IN WHICH DIRECTION SHOULD WE LOOK IN THE FUTURE?

Throughout this document we have covered several topics in explainable AI. More importantly, we have covered a novel taxonomy for explainable artificial intelligence that benefits from the experience and insight of causality. The proposed

taxonomy borrows concepts from the Ladder of Causality, leading to our proposed Ladder of xAI. In this, we have proposed a division of xAI method in three different rungs, where we consider that every deep learning model learnt has an inherent causal model that can represent its inner rules. The first rung, considers the simplest approach to xAI, visualisation methods, which bring observational information to the user. However, these methods are the equivalent of correlation for causality, as they answer a different question and are not guaranteed to be related to the model inner workings. On the second rung the number of different approaches grows significantly. First, we have approaches based on inserting prior to either the model or the training rule. Model priors, due to their strictness partially reveals the final causal representation of the learnt model, whereas training priors, due to its relaxed nature leads to a less effective revelation of that same model. Additionally, it is possible to make interventions at input and prediction levels. When there is an intervention at the input level, the researcher aims to know "How will this change affect my model predictions?". Future methods on the third rung will allow to answer this question immediately, and further answer to the "Why is this the behaviour of my model?". However, we are still far from reaching the third rung with deep learning-based methods.

The xAI ladder was applied to 81 different research papers published between 2018 and 2024. These papers proposed a myriad approaches from the first two rungs for different biometric modalities and tasks. From Face Recognition to Post-Mortem Iris Recognition and DeepFake Detection, the range of tasks that benefited from explainable approaches is outstanding. Moreover, there seems to be a growth trend. As such, it is important that researchers have a common language to address, categorise and assess explainable methods. Despite the current demonstration and application with biometrics works, the proposed taxonomy is sufficiently general to be applied to any domain of computer vision. Hence, other areas such as medical applications can benefit from this framework and converge to common ground. Additionally, this framework can be extended to fit the necessities of a specific domain as it is flexible and not bonded to a single strategy.

Additionally we have further discussed some points regarding explainable AI and how it could be used to improve different areas of biometrics. These discussed areas are frequently not seen as fields that can benefit from xAI, and for this we decided to illustrate some benefits that can be drawn from xAI approaches. In this sense, we have discussed the evaluation of models, the feedback given and the privacy protection of both user data and the trade-secrets behind an artificial intelligence model. Furthermore, we have discussed a few topics on unveiling biases with explainable AI. With the proposed taxonomy and strategy to approach xAI, the unveiling of biases begins during development and model assessment, reducing the amount of models deployed with hidden biases.

## 5.1 Future work

Despite the obvious progress, there are areas, specially in biometrics that were not fully covered in this survey. As such, we provide a glance on some relevant future worrk topics that might provide novel strategies in xAI, safer biometric systems or trustworthy models. And with the recent growing interest in artificial intelligence, mostly supported by the deployment of Large Language Models and Text-to-Image systems, it is of crucial importance that we remain aware of the dangers of these models, their misuse, and the strategies to mitigate them.

*5.1.1  Revisiting fundamental research in mathematics and computer science: Can we optimize our algorithms more efficiently?* Although a clear interpretation and analysis of some of the older machine learning algorithms was somewhat possible, the majority of that clearness has been lost with the deep learning uprising. It is not trivial to reformulate these solutions in order to decrease its opaqueness. However, as seen in the previous Sections, there are already

researchers working on variations of architectures and losses that promote a more controlled behavior. In an effort towards interpretable models, or models which can benefit from more complete explanations, it is crucial to redesign the current theory and algorithms with xAI in the scope.

### 5.1.2 Benefiting from multiple data sources: Can we build methodologies that take into account multi-modal data?

Working with multi-modal data might improve the performance of a unimodal method. However, this can be both a gift and a curse. When carefully implemented and designed, the integration of an additional data source can support explanations that are less ambiguous and more complete. Nonetheless, deep learning systems working on one modality are already increasingly complex, if an additional data source is incorporated, it might lead to more opaqueness. Hence, besides studying how to integrate multiple modalities, it is also extremely important to understand the compatibility of those modalities. The utilization of intrinsically interpretable models supports growth in the number of fused modalities without compromising the understanding of the model's inner workings.

### 5.1.3 What to believe: On the coherence of the explanations with human perception.

Humans are creatures of beliefs and prior knowledge, and despite their potential misalignment with reality, going against them creates a lot of entropy. And explanations are, for now, limited to not being capable of counterarguing someone that disagrees with it. Hence, an explanation that presents a different view than the one already acquired by the user might suffer from strong mistrust. Sometimes, this different perspective is on elements that are not the main focus of the explanation. And thus, to mitigate this, it is necessary to give explanations that also align with the users' view. To achieve this, it is necessary to craft the explanation method for the specific target user that is going to benefit from the machine learning model. Similarly, the subjectivity aspect of the explanations must be tolerated in certain use cases, whereas others require objective and precise explanations. Joining these concepts is not easy, and there is a significant amount of work to be done in order to direct the explanations to the target users.

### 5.1.4 Learning information: Correlation vs Causation.

Deep learning-based models are strong and complex function approximators also known as universal approximators. However, this complexity is not always beneficial. For instance, some background elements might happen more frequently on certain predictions, nonetheless, they are not the cause of that label. Learning these correlated relationships increases the chance that the network makes the inference based on them. The current framework, considers that each learnt model has an inherent causal model that explains its predictions, however, the learnt model is not guaranteed to be causal with the real world. In this sense, these correlations might represent causal knowledge when seen through the lenses of the deep model, potentially aggravating biases. An example of these is a face presentation attack detection system trained on images of people on a variety of environments and weathers. As shown by demographic information certain geographic areas are related to certain ethnic groups. A potential harmful correlation would be to classify images of people from certain ethnic groups as presentation attacks, just because the background is considerably different from the one used in the training data for the majority of that ethnic group. This relationship would present a causal rule within our learnt model that does not reflect the causal reality of the world.

Focusing on causes of a certain effect leads to more faithful and believable predictions. Therefore, there is a research direction on causal learning and causal inference. It is of utter importance to integrate these practices in biometric systems and biometric research. Recently, there were reported several cases of biased biometric systems. These biases can be unveiled by xAI and mitigated with causal inference. Hence, this research direction must continue to grow in future iterations of biometric systems.

*5.1.5 Intrinsically-interpretable models vs explanation methods.* While explanations methods represent the majority of xAI application in the biometrics domains, they do not represent a definitve solution for the black box problem. On the other hand, while intrinsically interpretable models, which are specifically designed to promote interpretability, are great tools to reduce the opaqueness of a model, they can still benefit of explanation generators. We believe that the future of explainable artificial intelligence relies in the joint work of in- and post-model approaches. It can be done through the development of differential post explanations that can provide feedback to the model, or through the design of explanation methods that rely on the intrinsic properties of a model.

## ACKNOWLEDGMENTS

## REFERENCES

[1] Ahmed Mahdi Abdulkareem and Anna Gordon. 2023. Evaluating the Usability and User Acceptance of Biometric Authentication in Different Applications. *Quarterly Journal of Emerging Technologies and Innovations* 8, 2 (2023), 1–10.

[2] Mohamed Abomhara, Sule Yildirim Yayilgan, Anne Hilde Nymoen, Marina Shalaginova, Zoltán Székely, and Ogerta Elezaj. 2020. How to Do It Right: A Framework for Biometrics Supported Border Control. In *E-Democracy – Safeguarding Democracy and Human Rights in the Digital Age*, Sokratis Katsikas and Vasilios Zorkadis (Eds.). Springer International Publishing, Cham, 94–109.

[3] Poorya Aghdaie, Baaria Chaudhary, Sobhan Soleymani, Jeremy Dawson, and Nasser M Nasrabadi. 2021. Attention Aware Wavelet-based Detection of Morphed Face Images. In *IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 1–8.

[4] Mazida A Ahmed, Ridip Dev Choudhury, and Kishore Kashyap. 2020. Race estimation with deep networks. *Journal of King Saud University-Computer and Information Sciences* (2020).

[5] Stephan Algermissen and Max Hörnlein. 2021. Person Identification by Footstep Sound Using Convolutional Neural Networks. *Applied Mechanics* 2, 2 (2021), 257–273.

[6] Raed Alharbi, Minh N Vu, and My T Thai. 2021. Learning Interpretation with Explainable Knowledge Distillation. In *IEEE International Conference on Big Data (Big Data)*. IEEE, 705–714.

[7] Miguel Almeida, Paulo Lobato Correia, and Peter Kastmand Larsen. 2016. BioFoV – an open platform for forensic video analysis and biometric data extraction. In *4th International Conference on Biometrics and Forensics (IWBF)*. 1–6. https://doi.org/10.1109/IWBF.2016.7449693

[8] Hammam Alshazly, Christoph Linse, Erhardt Barth, Sahar Ahmed Idris, and Thomas Martinetz. 2021. Towards Explainable Ear Recognition Systems Using Deep Residual Networks. *IEEE Access* 9 (2021), 122254–122273.

[9] David Anghelone, Cunjian Chen, Philippe Faure, Arun Ross, and Antitza Dantcheva. 2021. Explainable Thermal to Visible Face Recognition Using Latent-Guided Generative Adversarial Network. In *16th IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 1–8.

[10] Stanislaw Antol, Aishwarya Agrawal, Jiasen Lu, Margaret Mitchell, Dhruv Batra, C Lawrence Zitnick, and Devi Parikh. 2015. Vqa: Visual question answering. In *Proceedings of the IEEE international conference on computer vision*. 2425–2433.

[11] Gustavo Aquino, Marly GF Costa, and Cicero FF Costa Filho. 2022. Explaining One-Dimensional Convolutional Models in Human Activity Recognition and Biometric Identification Tasks. *Sensors* 22, 15 (2022), 5644.

[12] Tashreef Abdullah Araf, Ayesha Siddika, Sadullah Karimi, and Md Golam Rabiul Alam. 2022. Real-Time Face Emotion Recognition and Visualization using Grad-CAM. In *2022 Second International Conference on Advances in Electrical, Computing, Communication and Sustainable Technologies (ICAECT)*. IEEE, 1–5.

[13] Lei Jimmy Ba and Rich Caruana. 2014. Do Deep Nets Really Need to Be Deep?. In *Proceedings of the 27th International Conference on Neural Information Processing Systems - Volume 2* (Montreal, Canada) *(NIPS'14)*. MIT Press, Cambridge, MA, USA, 2654–2662.

[14] Sebastian Bach, Alexander Binder, Grégoire Montavon, Frederick Klauschen, Klaus-Robert Müller, and Wojciech Samek. 2015. On pixel-wise explanations for non-linear classifier decisions by layer-wise relevance propagation. *PloS one* 10, 7 (2015), e0130140.

[15] David Baehrens, Timon Schroeter, Stefan Harmeling, Motoaki Kawanabe, Katja Hansen, and Klaus-Robert Müller. 2010. How to explain individual classification decisions. *The Journal of Machine Learning Research* 11 (2010), 1803–1831.

[16] Neeru Bala, Rashmi Gupta, and Anil Kumar. 2022. Multimodal biometric system based on fusion techniques: a review. *Information Security Journal: A Global Perspective* 31, 3 (2022), 289–337.

[17] Alejandro Barredo Arrieta, Natalia Díaz-Rodríguez, Javier Del Ser, Adrien Bennetot, Siham Tabik, Alberto Barbado, Salvador Garcia, Sergio Gil-Lopez, Daniel Molina, Richard Benjamins, Raja Chatila, and Francisco Herrera. 2020. Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion* 58 (2020), 82–115. https://doi.org/10.1016/j.inffus.2019.12.012

[18] David Bau, Bolei Zhou, Aditya Khosla, Aude Oliva, and Antonio Torralba. 2017. Network Dissection: Quantifying Interpretability of Deep Visual Representations. In *CVPR*.

[19] Ying Bian, Peng Zhang, Jingjing Wang, Chunmao Wang, and Shiliang Pu. 2022. Learning multiple explainable and generalizable cues for face anti-spoofing. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 2310–2314.

[20] Serena Booth, Yilun Zhou, Ankit Shah, and Julie Shah. 2021. Bayes-TrEx: a Bayesian Sampling Approach to Model Transparency by Example. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 35. 11423–11432.

[21] Naima Bousnina, João Ascenso, Paulo Lobato Correia, and Fernando Pereira. 2023. A RISE-Based Explainability Method for Genuine and Impostor Face Verification. In *2023 International Conference of the Biometrics Special Interest Group (BIOSIG)*. IEEE, 1–6.

[22] Fadi Boutros, Naser Damer, Jan Niklas Kolf, Kiran Raja, Florian Kirchbuchner, Raghavendra Ramachandra, Arjan Kuijper, Pengcheng Fang, Chao Zhang, Fei Wang, et al. 2021. MFR 2021: Masked face recognition competition. In *IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 1–10.

[23] Fadi Boutros, Naser Damer, and Arjan Kuijper. 2022. Quantface: Towards lightweight face recognition by synthetic data low-bit quantization. In *2022 26th International Conference on Pattern Recognition (ICPR)*. IEEE, 855–862.

[24] Fadi Boutros, Naser Damer, Kiran Raja, Florian Kirchbuchner, and Arjan Kuijper. 2022. Template-Driven Knowledge Distillation for Compact and Accurate Periocular Biometrics Deep-Learning Models. *Sensors* 22, 5 (2022), 1921.

[25] Fadi Boutros, Meiling Fang, Marcel Klemt, Biying Fu, and Naser Damer. 2023. CR-FIQA: face image quality assessment by learning sample relative classifiability. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5836–5845.

[26] Aidan Boyd, Daniel Moreira, Andrey Kuehlkamp, Kevin Bowyer, and Adam Czajka. 2023. Human Saliency-Driven Patch-based Matching for Interpretable Post-mortem Iris Recognition. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 701–710.

[27] Leo Breiman, Jerome Friedman, Charles J Stone, and Richard A Olshen. 1984. *Classification and regression trees*. CRC press.

[28] Mario Bunge. 2017. *Causality and modern science*. Routledge.

[29] Joy Buolamwini and Timnit Gebru. 2018. Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification. In *Proceedings of the 1st Conference on Fairness, Accountability and Transparency (Proceedings of Machine Learning Research, Vol. 81)*, Sorelle A. Friedler and Christo Wilson (Eds.). PMLR, New York, NY, USA, 77–91. http://proceedings.mlr.press/v81/buolamwini18a.html

[30] Ruth M. J. Byrne. 2007. Précis of The Rational Imagination: How People Create Alternatives to Reality. *Behavioral and Brain Sciences* 30, 5-6 (2007), 439–453. https://doi.org/10.1017/S0140525X07002579

[31] Ruth M J Byrne. 2019. Counterfactuals in Explainable Artificial Intelligence (XAI): Evidence from Human Reasoning. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, [IJCAI-19]*. International Joint Conferences on Artificial Intelligence Organization, 6276–6282. https://doi.org/10.24963/ijcai.2019/876

[32] Eduarda Caldeira, Pedro C Neto, Tiago Gonçalves, Naser Damer, Ana F Sequeira, and Jaime S Cardoso. 2023. Unveiling the Two-Faced Truth: Disentangling Morphed Identities for Face Morphing Detection. *arXiv preprint arXiv:2306.03002* (2023).

[33] Eduarda Caldeira, Pedro C Neto, Marco Huber, Naser Damer, and Ana F Sequeira. 2024. Model Compression Techniques in Biometrics Applications: A Survey. *arXiv preprint arXiv:2401.10139* (2024).

[34] Toon Calders, Faisal Kamiran, and Mykola Pechenizkiy. 2009. Building classifiers with independency constraints. In *2009 IEEE international conference on data mining workshops*. IEEE, 13–18.

[35] Daniel C Castro, Ian Walker, and Ben Glocker. 2020. Causality matters in medical imaging. *Nature Communications* 11, 1 (2020), 3673.

[36] Eduardo Castro, Jose Costa Pereira, and Jaime S Cardoso. 2023. Symmetry-based regularization in deep breast cancer screening. *Medical Image Analysis* 83 (2023), 102690.

[37] Mario Cesarelli, Fabio Martinelli, Francesco Mercaldo, and Antonella Santone. 2022. Emotion Recognition from Facial Expression using Explainable Deep Learning. In *2022 IEEE Intl Conf on Dependable, Autonomic and Secure Computing, Intl Conf on Pervasive Intelligence and Computing, Intl Conf on Cloud and Big Data Computing, Intl Conf on Cyber Science and Technology Congress (DASC/PiCom/CBDCom/CyberSciTech)*. IEEE, 1–6.

[38] Praveen Kumar Chandaliya and Neeta Nain. 2021. Child Face Age Progression and Regression using Self-Attention Multi-Scale Patch GAN. In *IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 1–8.

[39] Angelos Chatzimparmpas, Rafael M Martins, Ilir Jusufi, and Andreas Kerren. 2020. A survey of surveys on the use of visualization for interpreting machine learning models. *Information Visualization* 19, 3 (2020), 207–233. https://doi.org/10.1177/1473871620904671

[40] Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett, Cynthia Rudin, and Jonathan K Su. 2019. This looks like that: deep learning for interpretable image recognition. *Advances in neural information processing systems* 32 (2019).

[41] Chaofan Chen, Oscar Li, Daniel Tao, Alina Barnett, Cynthia Rudin, and Jonathan K Su. 2019. This Looks Like That: Deep Learning for Interpretable Image Recognition. In *Advances in Neural Information Processing Systems*, H. Wallach, H. Larochelle, A. Beygelzimer, F. d'Alché-Buc, E. Fox, and R. Garnett (Eds.), Vol. 32. Curran Associates, Inc. https://proceedings.neurips.cc/paper/2019/file/adf7ee2dcf142b0e11888e72b43fcb75-Paper.pdf

[42] Cunjian Chen and Arun Ross. 2021. An Explainable Attention-Guided Iris Presentation Attack Detector. In *IEEE Winter Conference on Applications of Computer Vision Workshops (WACVW)*. 97–106. https://doi.org/10.1109/WACVW52041.2021.00015

[43] Shen Chen, Taiping Yao, Keyue Zhang, Yang Chen, Ke Sun, Shouhong Ding, Jilin Li, Feiyue Huang, and Rongrong Ji. 2021. A Dual-Stream Framework for 3D Mask Face Presentation Attack Detection. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 834–841.

[44] Zhi Chen, Yijie Bei, and Cynthia Rudin. 2020. Concept whitening for interpretable image recognition. *Nature Machine Intelligence* 2, 12 (2020), 772–782.

[45] Anurag Chowdhury, Simon Kirchgasser, Andreas Uhl, and Arun Ross. 2020. Can a CNN Automatically Learn the Significance of Minutiae Points for Fingerprint Matching?. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 351–359.

[46] Cynthia M. Cook, John J. Howard, Yevgeniy B. Sirotin, Jerry L. Tipton, and Arun R. Vemury. 2019. Demographic Effects in Facial Recognition and Their Dependence on Image Acquisition: An Evaluation of Eleven Commercial Systems. *IEEE Transactions on Biometrics, Behavior, and Identity Science* 1, 1 (2019), 32–41. https://doi.org/10.1109/TBIOM.2019.2897801

[47] Ricardo Correia, Paulo Correia, and Fernando Pereira. 2023. Face Verification Explainability Heatmap Generation. In *2023 International Conference of the Biometrics Special Interest Group (BIOSIG)*. IEEE, 1–5.

[48] Laurine Dargaud, Mathias Ibsen, Juan Tapia, and Christoph Busch. 2023. A principal component analysis-based approach for single morphing attack detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 683–692.

[49] Arun Das and Paul Rad. 2020. Opportunities and challenges in explainable artificial intelligence (xai): A survey. *arXiv preprint arXiv:2006.11371* (2020).

[50] Eline Dekyvere. 2016. Face/off: Is Snapchat stealing your face one cute dog filter at a time? - CITIP blog. https://www.law.kuleuven.be/citip/blog/faceoff-is-snapchat-stealing-your-face-one-cute-dog-filter-at-a-time/

[51] Tribikram Dhar, Nilanjan Dey, Surekha Borra, and R Simon Sherratt. 2023. Challenges of Deep Learning in Medical Image Analysis—Improving Explainability and Trust. *IEEE Transactions on Technology and Society* 4, 1 (2023), 68–75.

[52] Moises Diaz, Miguel A Ferrer, and Gennaro Vessio. 2023. Explainable offline automatic signature verifier to support forensic handwriting examiners. *Neural Computing and Applications* (2023), 1–17.

[53] Yinpeng Dong, Hang Su, Baoyuan Wu, Zhifeng Li, Wei Liu, Tong Zhang, and Jun Zhu. 2019. Efficient Decision-Based Black-Box Adversarial Attacks on Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition (CVPR)*.

[54] Jon Donnelly, Alina Jade Barnett, and Chaofan Chen. 2022. Deformable protonet: An interpretable image classifier using deformable prototypes. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10265–10275.

[55] Finale Doshi-Velez and Been Kim. 2017. Towards a rigorous science of interpretable machine learning. *arXiv preprint arXiv:1702.08608* (2017).

[56] Alexey Dosovitskiy and Thomas Brox. 2016. Inverting visual representations with convolutional networks. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 4829–4837.

[57] Pawel Drozdowski, Christian Rathgeb, Antitza Dantcheva, Naser Damer, and Christoph Busch. 2020. Demographic Bias in Biometrics: A Survey on an Emerging Challenge. *IEEE Transactions on Technology and Society* 1, 2 (2020), 89–103. https://doi.org/10.1109/TTS.2020.2992344

[58] Juan Manuel Durán and Karin Rolanda Jongsma. 2021. Who is afraid of black box algorithms? On the epistemological and ethical basis of trust in medical AI. *Journal of Medical Ethics* 47, 5 (2021), 329–335.

[59] Rudresh Dwivedi, Ritesh Kumar, Deepak Chopra, Pranay Kothari, and Manjot Singh. 2023. An Efficient Ensemble Explainable AI (XAI) Approach for Morphed Face Detection. *arXiv preprint arXiv:2304.14509* (2023).

[60] Josh Ellenbogen. 2012. *Reasoned and Unreasoned Images: The Photography of Bertillon, Galton, and Marey*. Penn State Press.

[61] Dumitru Erhan, Yoshua Bengio, Aaron Courville, and Pascal Vincent. 2009. Visualizing higher-layer features of a deep network. *University of Montreal* 1341, 3 (2009), 1.

[62] European Parliament, Council of the European Union. 2016. Regulation (EU) 2016/679 of the European Parliament and of the Council of 27 April 2016 on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation) (Text with EEA relevance). http://data.europa.eu/eli/reg/2016/679/oj. , 88 pages. http://data.europa.eu/eli/reg/2016/679/oj Online; accessed on 2021-08-07.

[63] Meiling Fang, Naser Damer, Fadi Boutros, Florian Kirchbuchner, and Arjan Kuijper. 2021. Iris Presentation Attack Detection by Attention-based and Deep Pixel-wise Binary Supervision Network. In *IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 1–8.

[64] Meiling Fang, Naser Damer, Florian Kirchbuchner, and Arjan Kuijper. 2022. Learnable Multi-level Frequency Decomposition and Hierarchical Attention Mechanism for Generalized Face Presentation Attack Detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 3722–3731.

[65] Kelwin Fernandes, Jaime S Cardoso, and Birgitte Schmidt Astrup. 2018. A deep learning approach for the forensic evaluation of sexual assault. *Pattern Analysis and Applications* 21, 3 (2018), 629–640.

[66] FICO. 2006. *Introduction to Scorecard for FICO Model Builder*. Technical Report. Fair Isaac Corporation.

[67] Danilo Franco, Nicolo Navarin, Michele Donini, Davide Anguita, and Luca Oneto. 2022. Deep fair models for complex data: Graphs labeling and explainable face recognition. *Neurocomputing* 470 (2022), 318–334.

[68] Danilo Franco, Luca Oneto, Nicolò Navarin, and Davide Anguita. 2021. Learn and Visually Explain Deep Fair Models: an Application to Face Recognition. In *2021 International Joint Conference on Neural Networks (IJCNN)*. IEEE, 1–10.

[69] Biying Fu and Naser Damer. 2022. Explainability of the Implications of Supervised and Unsupervised Face Image Quality Estimations Through Activation Map Variation Analyses in Face Recognition Models. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 349–358.

[70] Biying Fu and Naser Damer. 2022. Towards Explaining Demographic Bias through the Eyes of Face Recognition Models. In *2022 IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 1–10.

[71] Robert D Furberg, Travis Taniguchi, Brian Aagaard, Alexa M Ortiz, Meghan Hegarty-Craver, Kristin H Gilchrist, and Ty A Ridenour. 2017. Biometrics and Policing: A Protocol for Multichannel Sensor Data Collection and Exploratory Analysis of Contextualized Psychophysiological Response During Law Enforcement Operations. *JMIR Research Protocols* 6, 3 (mar 2017), e44. https://doi.org/10.2196/resprot.7499

[72] Francis Galton. 1892. *Finger prints*. Number 57490-57492. Cosimo Classics.

[73] Angelo Genovese, Vincenzo Piuri, and Fabio Scotti. 2019. Towards Explainable Face Aging with Generative Adversarial Networks. In *IEEE International Conference on Image Processing (ICIP)*. 3806–3810. https://doi.org/10.1109/ICIP.2019.8803616

[74] Leilani H Gilpin, David Bau, Ben Z Yuan, Ayesha Bajwa, Michael Specter, and Lalana Kagal. 2018. Explaining explanations: An overview of interpretability of machine learning. In *IEEE 5th International Conference on data science and advanced analytics (DSAA)*. IEEE, 80–89.

[75] Tiago Gonçalves, Isabel Rio-Torto, Luís F Teixeira, and Jaime S Cardoso. 2022. A survey on attention mechanisms for medical applications: are we moving towards better algorithms? (2022).

[76] Ian Goodfellow, Yoshua Bengio, Aaron Courville, and Yoshua Bengio. 2016. *Deep learning*. Vol. 1. MIT press Cambridge.

[77] Margarida Gouveia, Eduardo Meca, Ana Rebelo, and Jaime S. Cardoso. 2023. Deep Minutiae Fingerprint Extraction Using Equivariance Priors. In *Proceedings of the International Conference on Bio-inspired Systems and Signal Processing (BIOSIGNALS)*.

[78] Patrick J Grother, Patrick J Grother, and Mei Ngan. 2014. *Face recognition vendor test (frvt)*. US Department of Commerce, National Institute of Standards and Technology.

[79] Patrick J Grother, Mei L Ngan, Kayee K Hanaoka, et al. 2018. Ongoing face recognition vendor test (frvt) part 2: Identification. (2018).

[80] Riccardo Guidotti, Anna Monreale, Salvatore Ruggieri, Franco Turini, Fosca Giannotti, and Dino Pedreschi. 2018. A survey of methods for explaining black box models. *ACM computing surveys (CSUR)* 51, 5 (2018), 1–42.

[81] Han Guo, Nazneen Fatema Rajani, Peter Hase, Mohit Bansal, and Caiming Xiong. 2020. Fastif: Scalable influence functions for efficient model interpretation and debugging. *arXiv preprint arXiv:2012.15781* (2020).

[82] Chunrui Han, Shiguang Shan, Meina Kan, Shuzhe Wu, and Xilin Chen. 2022. Personalized Convolution for Face Recognition. *International Journal of Computer Vision* (2022), 1–19.

[83] Peter Hase, Chaofan Chen, Oscar Li, and Cynthia Rudin. 2019. Interpretable image recognition with hierarchical prototypes. In *Proceedings of the AAAI Conference on Human Computation and Crowdsourcing*, Vol. 7. 32–40.

[84] Trevor Hastie and Robert Tibshirani. 1986. Generalized Additive Models. *Statist. Sci.* 1, 3 (1986), 297 – 310. https://doi.org/10.1214/ss/1177013604

[85] Lisa Anne Hendricks, Zeynep Akata, Marcus Rohrbach, Jeff Donahue, Bernt Schiele, and Trevor Darrell. 2016. Generating visual explanations. In *European Conference on Computer Vision*. Springer, 3–19.

[86] Lisa Anne Marie Hendricks. 2019. *Visual Understanding through Natural Language*. University of California, Berkeley.

[87] Qingqiao Hu, Siyang Yin, Huiyang Ni, and Yisiyuan Huang. 2020. An End to End Deep Neural Network for Iris Recognition. *Procedia Computer Science* 174 (2020), 505–517.

[88] Marco Huber, Fadi Boutros, Florian Kirchbuchner, and Naser Damer. 2021. Mask-invariant Face Recognition through Template-level Knowledge Distillation. In *16th IEEE International Conference on Automatic Face and Gesture Recognition*. IEEE, 1–8.

[89] Marco Huber, Fadi Boutros, Anh Thi Luu, Kiran Raja, Raghavendra Ramachandra, Naser Damer, Pedro C Neto, Tiago Gonçalves, Ana F Sequeira, Jaime S Cardoso, et al. 2022. SYN-MAD 2022: Competition on face morphing attack detection based on privacy-aware synthetic training data. In *2022 IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 1–10.

[90] Marco Huber, Fadi Boutros, Anh Thi Luu, Kiran Raja, Raghavendra Ramachandra, Naser Damer, Pedro C. Neto, Tiago Gonçalves, Ana F. Sequeira, Jaime S. Cardoso, João Tremoço, Miguel Lourenço, Sergio Serra, Eduardo Cermeño, Marija Ivanovska, Borut Batagelj, Andrej Kronovšek, Peter Peer, and Vitomir Štruc. 2022. SYN-MAD 2022: Competition on Face Morphing Attack Detection Based on Privacy-aware Synthetic Training Data. (2022). https://doi.org/10.48550/ARXIV.2208.07337

[91] Marco Huber, Meiling Fang, Fadi Boutros, and Naser Damer. 2023. Are Explainability Tools Gender Biased? A Case Study on Face Presentation Attack Detection. *arXiv preprint arXiv:2304.13419* (2023).

[92] Marco Huber, Anh Thi Luu, Philipp Terhörst, and Naser Damer. 2024. Efficient explainable face verification based on similarity score argument backpropagation. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 4736–4745.

[93] Gordon Hull. 2023. Dirty data labeled dirt cheap: epistemic injustice in machine learning systems. *Ethics and Information Technology* 25, 3 (2023), 1–14.

[94] Frederik Hvilshøj, Alexandros Iosifidis, and Ira Assent. 2021. ECINN: Efficient Counterfactuals from Invertible Neural Networks. *arXiv preprint arXiv:2103.13701* (2021).

[95] Anil Jain, Lin Hong, and Ruud Bolle. 1997. On-line fingerprint verification. *IEEE transactions on pattern analysis and machine intelligence* 19, 4 (1997), 302–314.

[96] Ziwei Ji, YU Tiezheng, Yan Xu, Nayeon Lee, Etsuko Ishii, and Pascale Fung. 2023. Towards mitigating LLM hallucination via self reflection. In *The 2023 Conference on Empirical Methods in Natural Language Processing*.

[97] Haoran Jiang and Dan Zeng. 2021. Explainable Face Recognition Based on Accurate Facial Compositions. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 1503–1512.

[98] Indu Joshi, Riya Kothari, Ayush Utkarsh, Vinod K Kurmi, Antitza Dantcheva, Sumantra Dutta Roy, and Prem Kumar Kalra. 2021. Explainable Fingerprint ROI Segmentation Using Monte Carlo Dropout. In *IEEE Winter Conference on Applications of Computer Vision Workshops (WACVW)*. 60–69. https://doi.org/10.1109/WACVW52041.2021.00011

[99] Amin Jourabloo, Yaojie Liu, and Xiaoming Liu. 2018. Face de-spoofing: Anti-spoofing via noise modeling. In *Proceedings of the European Conference on Computer Vision (ECCV)*. 290–306.

[100] Hong-Gyu Jung, Sin-Han Kang, Hee-Dong Kim, Dong-Ok Won, and Seong-Whan Lee. 2020. Counterfactual explanation based on gradual construction for deep networks. *arXiv preprint arXiv:2008.01897* (2020).

[101] Margot E Kaminski. 2019. The right to explanation, explained. *Berkeley Tech. LJ* 34 (2019), 189.

[102] Andrej Karpathy, George Toderici, Sanketh Shetty, Thomas Leung, Rahul Sukthankar, and Li Fei-Fei. 2014. Large-scale video classification with convolutional neural networks. In *Proc. of the IEEE CVPR*. 1725–1732.

[103] Frank Keil, Leonid Rozenblit, and Candice Mills. 2004. What lies beneath? Understanding the limits of understanding. *Thinking and seeing: Visual metacognition in adults and children* (2004), 227–49.

[104] Saeed Khorram and Li Fuxin. 2022. Cycle-consistent counterfactuals by latent transformations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10203–10212.

[105] Been Kim and Finale Doshi-Velez. 2017. Interpretable machine learning: the fuss, the concrete and the questions. *ICML Tutorial on interpretable machine learning* (2017).

[106] Brendan F. Klare, Mark J. Burge, Joshua C. Klontz, Richard W. Vorder Bruegge, and Anil K. Jain. 2012. Face Recognition Performance: Role of Demographic Information. *IEEE Transactions on Information Forensics and Security* 7, 6 (2012), 1789–1801. https://doi.org/10.1109/TIFS.2012.2214212

[107] Martin Knoche, Torben Teepe, Stefan Hörmann, and Gerhard Rigoll. 2023. Explainable model-agnostic similarity and confidence in face verification. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 711–718.

[108] Pavel Korshunov, Anubhav Jain, and Sébastien Marcel. 2022. Custom Attribution Loss for Improving Generalization and Interpretability of Deepfake Detection. In *ICASSP 2022-2022 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*. IEEE, 8972–8976.

[109] Adam Kortylewski, Bernhard Egger, Andreas Schneider, Thomas Gerig, Andreas Morel-Forster, and Thomas Vetter. 2018. Empirically Analyzing the Effect of Dataset Biases on Deep Face Recognition Systems. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR) Workshops*.

[110] Andrey Kuehlkamp, Aidan Boyd, Adam Czajka, Kevin Bowyer, Patrick Flynn, Dennis Chute, and Eric Benjamin. 2022. Interpretable Deep Learning-Based Forensic Iris Segmentation and Recognition. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 359–368.

[111] Kenneth Lai, Helder C. R. Oliveira, Ming Hou, Svetlana N. Yanushkevich, and Vlad P. Shmerko. 2020. Risk, Trust, and Bias: Causal Regulators of Biometric-Enabled Decision Support. *IEEE Access* 8 (2020), 148779–148792. https://doi.org/10.1109/ACCESS.2020.3015855

[112] Himabindu Lakkaraju and Osbert Bastani. 2020. "How do I fool you?" Manipulating User Trust via Misleading Black Box Explanations. In *Proceedings of the AAAI/ACM Conference on AI, Ethics, and Society*. 79–85.

[113] Yann LeCun, Yoshua Bengio, and Geoffrey Hinton. 2015. Deep learning. *Nature* 521, 7553 (2015), 436–444.

[114] Pengyu Li. 2023. BioNet: A Biologically-Inspired Network for Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 10344–10354.

[115] Suk-Young Lim, Dong-Kyu Chae, and Sang-Chul Lee. 2022. Detecting deepfake voice using explainable deep learning techniques. *Applied Sciences* 12, 8 (2022), 3926.

[116] Xiaofeng Lin, Seungbae Kim, and Jungseock Joo. 2022. Fairgrape: Fairness-aware gradient pruning method for face attribute classification. In *European Conference on Computer Vision*. Springer, 414–432.

[117] Peter Lipton. 1990. Contrastive Explanation. *Royal Institute of Philosophy Supplement* 27 (1990), 247–266. https://doi.org/10.1017/S1358246100005130

[118] Decheng Liu, Xinbo Gao, Chunlei Peng, Nannan Wang, and Jie Li. 2021. Heterogeneous Face Interpretable Disentangled Representation for Joint Face Recognition and Synthesis. *IEEE Transactions on Neural Networks and Learning Systems* (2021).

[119] Shanhong Liu. 2021. COVID-19 impact on the use of biometrics 2021 | Statista. https://www.statista.com/statistics/1264372/covid-impact-biometrics-use/

[120] Xuan Liu, Xiaoguang Wang, and Stan Matwin. 2018. Improving the interpretability of deep neural networks with knowledge distillation. In *IEEE International Conference on Data Mining Workshops (ICDMW)*. IEEE, 905–912.

[121] Yaojie Liu, Amin Jourabloo, and Xiaoming Liu. 2018. Learning deep models for face anti-spoofing: Binary or auxiliary supervision. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 389–398.

[122] Yaojie Liu, Joel Stehouwer, and Xiaoming Liu. 2020. On disentangling spoof trace for generic face anti-spoofing. In *European Conference on Computer Vision*. Springer, 406–422.

[123] Max Maria Losch, Mario Fritz, and Bernt Schiele. 2020. Semantic Bottlenecks: Quantifying and Improving Inspectability of Deep Representations. In *Pattern Recognition: 42nd DAGM German Conference, DAGM GCPR, Tübingen, Germany, Proceedings 42*. Springer, 15–29.

[124] Yuhang Lu and Touradj Ebrahimi. 2023. Explanation of Face Recognition via Saliency Maps. *arXiv preprint arXiv:2304.06118* (2023).

[125] Rogério Martins. 2016. Why Are We Not Able to See Beyond Three Dimensions? *The Mathematical Intelligencer* 38, 4 (2016), 46–51. https://doi.org/10.1007/s00283-016-9670-1

[126] Sherin Mathews, Shivangee Trivedi, Amanda House, Steve Povolny, and Celeste Fralick. 2023. An explainable deepfake detection framework on a novel unconstrained dataset. *Complex & Intelligent Systems* (2023), 1–13.

[127] Rita Matulionyte. 2022. Reconciling trade secrets and explainable AI: face recognition technology as a case study. *European Intellectual Property Review* 44, 1 (2022), 36–42.

[128] Ghazal Mazaheri and Amit K Roy-Chowdhury. 2022. Detection and Localization of Facial Expression Manipulations. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 1035–1045.

[129] Domingo Mery and Bernardita Morris. 2022. On Black-Box Explanation for Face Verification. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 3418–3427.

[130] Caroline Lancelot Miltgen, Aleš Popovič, and Tiago Oliveira. 2013. Determinants of end-user acceptance of biometrics: Integrating the "Big 3" of technology acceptance with privacy context. *Decision support systems* 56 (2013), 103–114.

[131] Shervin Minaee, Amirali Abdolrashidi, Hang Su, Mohammed Bennamoun, and David Zhang. 2023. Biometrics recognition using deep learning: A survey. *Artificial Intelligence Review* (2023), 1–49.

[132] Shervin Minaee, Amirali Abdolrashidiy, and Yao Wang. 2016. An experimental study of deep convolutional features for iris recognition. In *2016 IEEE signal processing in medicine and biology symposium (SPMB)*. IEEE, 1–6.

[133] Yao Ming, Panpan Xu, Huamin Qu, and Liu Ren. 2019. Interpretable and steerable sequence learning via prototypes. In *Proceedings of the 25th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining*. 903–913.

[134] Nelida Mirabet-Herranz and Jean-Luc Dugelay. 2023. LVT Face Database: A benchmark database for visible and hidden face biometrics. In *BIOSIG 2023, 22nd International Conference of the Biometrics Special Interest Group*.

[135] Nelida Mirabet-Herranz, Chiara Galdi, and Jean-Luc Dugelay. 2022. Impact of Digital Face Beautification in Biometrics. In *EUVIP 2022, 10th European Workshop on Visual Information Processing*.

[136] Nelida Mirabet-Herranz, Khawla Mallat, and Jean-Luc Dugelay. 2022. Deep Learning for Remote Heart Rate Estimation: A Reproducible and Optimal State-of-the-Art Framework. In *International Conference on Pattern Recognition*. Springer, 558–573.

[137] Nelida Mirabet-Herranz, Khawla Mallat, and Jean-Luc Dugelay. 2023. New Insights on Weight Estimation from Face Images. In *2023 IEEE 17th International Conference on Automatic Face and Gesture Recognition (FG)*. IEEE, 1–6.

[138] Christoph Molnar. 2019. *Interpretable Machine Learning*. https://christophm.github.io/interpretable-ml-book/.

[139] Grégoire Montavon, Wojciech Samek, and Klaus-Robert Müller. 2018. Methods for interpreting and understanding deep neural networks. *Digital Signal Processing* 73 (2018), 1 – 15. https://doi.org/10.1016/j.dsp.2017.10.011

[140] Helena Montenegro and Jaime S Cardoso. 2023. Anonymizing medical case-based explanations through disentanglement. *arXiv preprint arXiv:2311.04833* (2023).

[141] Helena Montenegro, Wilson Silva, and Jaime S Cardoso. 2021. Privacy-Preserving Generative Adversarial Network for Case-Based Explainability in Medical Image Analysis. *IEEE Access* 9 (2021), 148037–148047.

[142] Helena Montenegro, Wilson Silva, Alex Gaudio, Matt Fredrikson, Asim Smailagic, and Jaime S Cardoso. 2022. Privacy-preserving Case-based Explanations: Enabling visual interpretability by protecting privacy. *IEEE Access* (2022).

[143] Henning Myhrvold, Haoyu Zhang, Juan Tapia, Raghavendra Ramachandra, and Christoph Günther Busch. 2022. Explainable Visualization for Morphing Attack Detection. (2022).

[144] Borum Nam, Joo Young Kim, Yeongmyeong Kim, Jiyoon Kim, Soonwon So, Hyung Youn Choi, and In Young Kim. 2022. Facialcuenet: An Interpretable Deception Detection Model for Criminal Interrogation Using Facial Expression. *SSRN Electronic Journal* (2022). https://doi.org/10.2139/ssrn.4208671

[145] Meike Nauta, Ron van Bree, and Christin Seifert. 2021. Neural Prototype Trees for Interpretable Fine-grained Image Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 14933–14943.

[146] Pedro C Neto, Fadi Boutros, João Ribeiro Pinto, Naser Damer, Ana F Sequeira, and Jaime S Cardoso. 2021. FocusFace: Multi-task Contrastive Learning for Masked Face Recognition. *arXiv preprint arXiv:2110.14940* (2021).

[147] Pedro C Neto, Fadi Boutros, Joao Ribeiro Pinto, Naser Damer, Ana F Sequeira, Jaime S Cardoso, Messaoud Bengherabi, Abderaouf Bousnat, Sana Boucheta, Nesrine Hebbadj, et al. 2022. OCFR 2022: Competition on Occluded Face Recognition From Synthetically Generated Structure-Aware Occlusions. *arXiv preprint arXiv:2208.02760* (2022).

[148] Pedro C Neto, Fadi Boutros, João Ribeiro Pinto, Mohsen Saffari, Naser Damer, Ana F Sequeira, and Jaime S Cardoso. 2021. My Eyes Are Up Here: Promoting Focus on Uncovered Regions in Masked Face Recognition. In *International Conference of the Biometrics Special Interest Group (BIOSIG)*. IEEE, 1–5.

[149] Pedro C Neto, Eduarda Caldeira, Jaime S Cardoso, and Ana F Sequeira. 2023. Compressed Models Decompress Race Biases: What Quantized Models Forget for Fair Face Recognition. *BIOSIG 2023, 22nd International Conference of the Biometrics Special Interest Group* (2023).

[150] Pedro C. Neto, Tiago Gonçalves, Marco Huber, Naser Damer, Ana F. Sequeira, and Jaime S. Cardoso. 2022. OrthoMAD: Morphing Attack Detection Through Orthogonal Identity Disentanglement. https://doi.org/10.48550/ARXIV.2208.07841

[151] Pedro C. Neto, Joao Ribeiro Pinto, Fadi Boutros, Naser Damer, Ana F. Sequeira, and Jaime S. Cardoso. 2022. Beyond Masks: On the Generalization of Masked Face Recognition Models to Occluded Face Recognition. *IEEE Access* (2022), 1–1. https://doi.org/10.1109/ACCESS.2022.3199014

[152] P. C. Neto, A. F. Sequeira, and J. S. Cardoso. 2022. Myope Models - Are face presentation attack detection models short-sighted?. In *IEEE/CVF Winter Conference on Applications of Computer Vision Workshops (WACVW)*. IEEE Computer Society, Los Alamitos, CA, USA, 390–399. https://doi.org/10.1109/WACVW54805.2022.00045

[153] Pedro C Neto, Ana F Sequeira, Jaime S Cardoso, and Philipp Terhörst. 2023. PIC-Score: Probabilistic Interpretable Comparison Score for Optimal Matching Confidence in Single-and Multi-Biometric Face Recognition. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 1021–1029.

[154] Luca Oneto, Michele Donini, Giulia Luise, Carlo Ciliberto, Andreas Maurer, and Massimiliano Pontil. 2020. Exploiting mmd and sinkhorn divergences for fair and transferable representation learning. *Advances in Neural Information Processing Systems* 33 (2020), 15360–15370.

[155] Shi Pan, Sanaul Hoque, and Farzin Deravi. 2022. An Attention-Guided Framework for Explainable Biometric Presentation Attack Detection. *Sensors* 22, 9 (2022), 3365.

[156] Dong Huk Park, Lisa Anne Hendricks, Zeynep Akata, Anna Rohrbach, Bernt Schiele, Trevor Darrell, and Marcus Rohrbach. 2018. Multimodal explanations: Justifying decisions and pointing to the evidence. In *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition*. 8779–8788.

[157] Judea Pearl. 1985. Bayesian netwcrks: A model cf self-activated memory for evidential reasoning.

[158] Judea Pearl. 2009. *Causality*. Cambridge university press.

[159] Judea Pearl and Dana Mackenzie. 2018. *The book of why: the new science of cause and effect*. Basic books.

[160] Bo Peng, Hongxing Fan, Wei Wang, Jing Dong, Yuezun Li, Siwei Lyu, Qi Li, Zhenan Sun, Han Chen, Baoying Chen, et al. 2021. DFGC 2021: A deepfake game competition. In *2021 IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 1–8.

[161] Martin Pernuš, Mansi Bhatnagar, Badr Samad, Divyanshu Singh, Peter Peer, Vitomir štruc, and Simon Dobrišek. 2023. ChildNet: Structural Kinship Face Synthesis Model With Appearance Control Mechanisms. *IEEE Access* 11 (2023), 49971–49991. https://doi.org/10.1109/ACCESS.2023.3276877

[162] Vitali Petsiuk, Abir Das, and Kate Saenko. 2018. RISE: Randomized Input Sampling for Explanation of Black-box Models. In *British Machine Vision Conference 2018, BMVC 2018, Newcastle, UK, September 3-6, 2018*. BMVA Press, 151. http://bmvc2018.org/contents/papers/1064.pdf

[163] Allan Pinto, Siome Goldenstein, Alexandre Ferreira, Tiago Carvalho, Helio Pedrini, and Anderson Rocha. 2020. Leveraging shape, reflectance and albedo from shading for face presentation attack detection. *IEEE Transactions on Information Forensics and Security* 15 (2020), 3347–3358.

[164] João Ribeiro Pinto. 2023. Seamless Multimodal Biometrics for Continuous Personalised Wellbeing Monitoring. *arXiv preprint arXiv:2301.03045* (2023).

[165] J. R. Pinto and J. S. Cardoso. 2020. Explaining ECG Biometrics: Is It All In The QRS?. In *Proceedings of the 19th International Conference of the Biometrics Special Interest Group (BIOSIG)*. Darmstadt, Germany.

[166] Shyam Sunder Prasad, Naval Kishore Mehta, Ankit Shukla, Pranav Mahajan, Arshdeep Singh, Sumeet Saurav, and Sanjay Singh. 2022. A Self-explainable Face Anti-spoofing Solution Based on Depth Estimation. In *International Conference on Frontiers of Intelligent Computing: Theory and Applications*. Springer, 103–113.

[167] Mary K. Pratt. 2021. Biometric security technology could see growth in 2021. https://searchsecurity.techtarget.com/feature/Biometric-security-technology-could-see-growth

[168] S. Prema, Mohamed Riyas V. S. Deen, Murali V. P. Krishna, and S. Praveen. 2019. Vehicle And License Authentication Using Finger Print. In *5th International Conference on Advanced Computing Communication Systems (ICACCS)*. 737–740. https://doi.org/10.1109/ICACCS.2019.8728402

[169] Zhongang Qi, Saeed Khorram, and Li Fuxin. 2021. Embedding deep networks into visual explanations. *Artificial Intelligence* 292 (2021), 103435.

[170] Anuj Rai, Somnath Dey, Pradeep Patidar, and Prakhar Rai. 2023. EXPRESSNET: An Explainable Residual Slim Network for Fingerprint Presentation Attack Detection. *arXiv preprint arXiv:2305.09397* (2023).

[171] Ankit Rajpal, Khushwant Sehra, Rashika Bagri, and Pooja Sikka. 2023. XAI-FR: Explainable AI-Based Face Recognition Using Deep Neural Networks. *Wireless Personal Communications* 129, 1 (2023), 663–680.

[172] Raghavendra Ramachandra and Hailin Li. 2023. Finger-NestNet: Interpretable Fingerphoto Verification on Smartphone using Deep Nested Residual Network. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 693–700.

[173] Manish Rathod, Chirag Dalvi, Kulveen Kaur, Shruti Patil, Shilpa Gite, Pooja Kamat, Ketan Kotecha, Ajith Abraham, and Lubna Abdelkareim Gabralla. 2022. Kids' Emotion Recognition Using Various Deep-Learning Models with Explainable AI. *Sensors* 22, 20 (2022), 8066.

[174] Isabel Rio-Torto, Jaime S Cardoso, and Luis Filipe Teixeira. 2024. Parameter-Efficient Generation of Natural Language Explanations for Chest X-ray Classification. In *Submitted to Medical Imaging with Deep Learning*. https://openreview.net/forum?id=EWbMmmQnVy under review.

[175] Ronald L Rivest. 1987. Learning decision lists. *Machine learning* 2, 3 (1987), 229–246.

[176] Andrea Macarulla Rodriguez, Luis Unzueta, Zeno Geradts, Marcel Worring, and Unai Elordi. 2023. Multi-Task Explainable Quality Networks for Large-Scale Forensic Facial Recognition. *IEEE Journal of Selected Topics in Signal Processing* (2023).

[177] Hiranmoy Roy, Debotosh Bhattacharjee, and Ondrej Krejcar. 2022. Interpretable Local Frequency Binary Pattern (LFrBP) Based Joint Continual Learning Network for Heterogeneous Face Recognition. *IEEE Transactions on Information Forensics and Security* 17 (2022), 2125–2136.

[178] Cynthia Rudin. 2019. Stop explaining black box machine learning models for high stakes decisions and use interpretable models instead. *Nature Machine Intelligence* 1, 5 (2019), 206–215.

[179] Wojciech Samek, Alexander Binder, Grégoire Montavon, Sebastian Lapuschkin, and Klaus-Robert Müller. 2016. Evaluating the visualization of what a deep neural network has learned. *IEEE transactions on neural networks and learning systems* 28, 11 (2016), 2660–2673.

[180] Wojciech Samek, Thomas Wiegand, and Klaus-Robert Müller. 2017. Explainable artificial intelligence: Understanding, visualizing and interpreting deep learning models. *arXiv preprint arXiv:1708.08296* (2017).

[181] Sascha Saralajew, Lars Holdijk, Maike Rees, Ebubekir Asan, and Thomas Villmann. 2019. Classification-by-components: probabilistic modeling of reasoning over a set of components. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*. 2792–2803.

[182] Clemens Seibold, Anna Hilsmann, and Peter Eisert. 2021. Focused LRP: Explainable AI for Face Morphing Attack Detection.. In *WACV (Workshops)*. 88–96.

[183] Ramprasaath R Selvaraju, Michael Cogswell, Abhishek Das, Ramakrishna Vedantam, Devi Parikh, and Dhruv Batra. 2017. Grad-cam: Visual explanations from deep networks via gradient-based localization. In *Proceedings of the IEEE international conference on computer vision*. 618–626.

[184] David Semedo, David Carmo, Ruslan Padnevych, and Joao Magalhaes. 2021. Contact-free Airport Borders with Biometrics-on-the-Move. In *IEEE International Workshop on Biometrics and Forensics (IWBF)*. IEEE, 1–2.

[185] Ana F Sequeira, Tiago Gonçalves, Wilson Silva, João Ribeiro Pinto, and Jaime S Cardoso. 2021. An exploratory study of interpretability for face presentation attack detection. *IET Biometrics* (2021).

[186] Rui Shao, Xiangyuan Lan, and Pong C Yuen. 2020. Regularized fine-grained meta face anti-spoofing. In *Proceedings of the AAAI Conference on Artificial Intelligence*, Vol. 34. 11974–11981.

[187] Renu Sharma and Arun Ross. 2020. D-NetPAD: An explainable and interpretable iris presentation attack detector. In *IEEE International Joint Conference on Biometrics (IJCB)*. IEEE, 1–10.

[188] Renu Sharma and Arun Ross. 2021. Viability of Optical Coherence Tomography for Iris Presentation Attack Detection. In *25th International Conference on Pattern Recognition (ICPR)*. IEEE, 6165–6172.

[189] Avanti Shrikumar, Peyton Greenside, and Anshul Kundaje. 2017. Learning important features through propagating activation differences. In *International Conference on Machine Learning*. PMLR, 3145–3153.

[190] Joseph Sill. 1997. Monotonic networks. In *Advances in Neural Information Processing Systems (NeurIPS)*.

[191] Samuel Henrique Silva, Mazal Bethany, Alexis Megan Votto, Ian Henry Scarff, Nicole Beebe, and Peyman Najafirad. 2022. Deepfake forensics analysis: An explainable hierarchical ensemble of weakly supervised models. *Forensic Science International: Synergy* 4 (2022), 100217.

[192] Wilson Silva, Alexander Poellinger, Jaime S Cardoso, and Mauricio Reyes. 2020. Interpretability-guided content-based medical image retrieval. In *International Conference on Medical Image Computing and Computer-Assisted Intervention*. Springer, 305–314.

[193] Karen Simonyan, Andrea Vedaldi, and Andrew Zisserman. 2014. Deep inside convolutional networks: Visualising image classification models and saliency maps. In *In Workshop at International Conference on Learning Representations*. Citeseer.

[194] Gurmail Singh and Kin-Choong Yow. 2021. These do not Look Like Those: An Interpretable Deep Learning Model for Image Recognition. *IEEE Access* 9 (2021), 41482–41493.

[195] Andrea C Skelly, Joseph R Dettori, and Erika D Brodt. 2012. Assessing bias: the importance of considering confounding. *Evidence-based spine-care journal* 3, 01 (2012), 9–12.

[196] Daniel Smilkov, Nikhil Thorat, Been Kim, Fernanda Viégas, and Martin Wattenberg. 2017. Smoothgrad: removing noise by adding noise. *arXiv preprint arXiv:1706.03825* (2017).

[197] Jost Tobias Springenberg, Alexey Dosovitskiy, Thomas Brox, and Martin Riedmiller. 2014. Striving for simplicity: The all convolutional net. *arXiv preprint arXiv:1412.6806* (2014).

[198] Nisha Srinivas, Karl Ricanek, Dana Michalski, David S. Bolme, and Michael King. 2019. Face Recognition Algorithm Bias: Performance Differences on Images of Children and Adults. In *IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops (CVPRW)*. 2269–2277. https://doi.org/10.1109/CVPRW.2019.00280

[199] Mukund Sundararajan, Ankur Taly, and Qiqi Yan. 2017. Axiomatic attribution for deep networks. In *International Conference on Machine Learning*. PMLR, 3319–3328.

[200] Yaniv Taigman, Ming Yang, Marc'Aurelio Ranzato, and Lior Wolf. 2014. Deepface: Closing the gap to human-level performance in face verification. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 1701–1708.

[201] Philipp Terhörst, Daniel Fährmann, Naser Damer, Florian Kirchbuchner, and Arjan Kuijper. 2020. Beyond identity: What information is stored in biometric face templates?. In *2020 IEEE international joint conference on biometrics (IJCB)*. IEEE, 1–10.

[202] Philipp Terhörst, Marco Huber, Naser Damer, Florian Kirchbuchner, Kiran Raja, and Arjan Kuijper. 2023. Pixel-level face image quality assessment for explainable face recognition. *IEEE Transactions on Biometrics, Behavior, and Identity Science* (2023).

[203] Philipp Terhörst, Jan Niklas Kolf, Marco Huber, Florian Kirchbuchner, Naser Damer, Aythami Morales, Julian Fierrez, and Arjan Kuijper. 2021. A Comprehensive Study on Face Recognition Biases Beyond Demographics. *IEEE Transactions on Technology and Society* (2021), 1–1. https://doi.org/10.1109/TTS.2021.3111823

[204] Simen Thys, Wiebe Van Ranst, and Toon Goedemé. 2019. Fooling automated surveillance cameras: adversarial patches to attack person detection. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition workshops*. 0–0.

[205] Mateusz Trokielewicz, Adam Czajka, and Piotr Maciejewicz. 2018. Presentation attack detection for cadaver iris. In *IEEE 9th International Conference on Biometrics Theory, Applications and Systems (BTAS)*. IEEE, 1–10.

[206] Mateusz Trokielewicz, Adam Czajka, and Piotr Maciejewicz. 2019. Perception of image features in post-mortem iris recognition: Humans vs machines. In *IEEE 10th International Conference on Biometrics Theory, Applications and Systems (BTAS)*. IEEE, 1–8.

[207] Alexander Unnervik and Sébastien Marcel. 2022. An anomaly detection approach for backdoored neural networks: face recognition as a case study. In *2022 International Conference of the Biometrics Special Interest Group (BIOSIG)*. IEEE, 1–5.

[208] Lev Utkin, Pavel Drobintsev, Maxim Kovalev, and Andrei Konstantinov. 2021. Combining an Autoencoder and a Variational Autoencoder for Explaining the Machine Learning Model Predictions. In *28th Conference of Open Innovations Association (FRUCT)*. IEEE, 489–494.

[209] Fatemeh Vakhshiteh, Ahmad Nickabadi, and Raghavendra Ramachandra. 2021. Adversarial attacks against face recognition: A comprehensive study. *IEEE Access* 9 (2021), 92735–92756.

[210] Laurens van der Maaten and Geoffrey Hinton. 2008. Visualizing Data using t-SNE. *Journal of Machine Learning Research* 9, 86 (2008), 2579–2605. http://jmlr.org/papers/v9/vandermaaten08a.html

[211] Nicole Van Hoeck, Patrick D Watson, and Aron K Barbey. 2015. Cognitive neuroscience of human counterfactual reasoning. *Frontiers in Human Neuroscience* 9 (2015), 420. https://doi.org/10.3389/fnhum.2015.00420

[212] Kushal Vangara, Michael C King, Vitor Albiero, Kevin Bowyer, et al. 2019. Characterizing the variability in face recognition accuracy relative to race. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 0–0.

[213] Kush R Varshney, Jamie C Rasmussen, Aleksandra Mojsilović, Moninder Singh, and Joan M DiMicco. 2012. Interactive visual salesforce analytics. In *33rd International Conference on Information Systems (ICIS)*.

[214] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. 2017. Attention is all you need. *Advances in neural information processing systems* 30 (2017).

[215] Brandon Vigliarolo. 2020. Apple's Face ID: Cheat sheet. https://www.techrepublic.com/article/apples-face-id-everything-iphone-x-users-need-to-know/

[216] Warren J von Eschenbach. 2021. Transparency and the black box problem: Why we do not trust AI. *Philosophy & Technology* 34, 4 (2021), 1607–1622.

[217] Haofan Wang, Zifan Wang, Mengnan Du, Fan Yang, Zijian Zhang, Sirui Ding, Piotr Mardziel, and Xia Hu. 2020. Score-CAM: Score-weighted visual explanations for convolutional neural networks. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*. 24–25.

[218] Mei Wang, Yaobin Zhang, and Weihong Deng. 2021. Meta Balanced Network for Fair Face Recognition. *IEEE Transactions on Pattern Analysis and Machine Intelligence* (2021).

[219] Tong Wang, Cynthia Rudin, Finale Doshi-Velez, Yimin Liu, Erica Klampfl, and Perry MacNeille. 2017. A bayesian framework for learning rule sets for interpretable classification. *The Journal of Machine Learning Research* 18, 1 (2017), 2357–2393.

[220] Yu-Chun Wang, Chien-Yi Wang, and Shang-Hong Lai. 2022. Disentangled Representation with Dual-stage Feature Learning for Face Anti-spoofing. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 1955–1964.

[221] Zezheng Wang, Zitong Yu, Chenxu Zhao, Xiangyu Zhu, Yunxiao Qin, Qiusheng Zhou, Feng Zhou, and Zhen Lei. 2020. Deep spatial gradient and temporal depth learning for face anti-spoofing. In *Proc. IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5042–5051.

[222] Jonathan R Williford, Brandon B May, and Jeffrey Byrne. 2020. Explainable Face Recognition. In *European Conference on Computer Vision*. Springer, 248–263.

[223] Martin Winter, Werner Bailer, and Georg Thallinger. 2022. Demystifying Face-Recognition with Locally Interpretable Boosted Features (LIBF). In *2022 10th European Workshop on Visual Information Processing (EUVIP)*. IEEE, 1–6.

[224] J B Woodward. 2003. Making Things Happen: A Theory of Causal Explanation. *Making Things Happen: A Theory of Causal Explanation* (2003). https://doi.org/10.1093/0195155270.001.0001

[225] Haiyu Wu and Kevin W Bowyer. 2023. A Real Balanced Dataset For Understanding Bias? Factors That Impact Accuracy, Not Numbers of Identities and Images. *arXiv preprint arXiv:2304.09818* (2023).

[226] Tongshuang Wu, Marco Tulio Ribeiro, Jeffrey Heer, and Daniel S Weld. 2021. Polyjuice: Automated, general-purpose counterfactual generation. *arXiv preprint arXiv:2101.00288* (2021).

[227] Rui Xin, Chudi Zhong, Zhi Chen, Takuya Takagi, Margo I. Seltzer, and Cynthia Rudin. 2022. Exploring the Whole Rashomon Set of Sparse Decision Trees. In *NeurIPS*.

[228] Kaidi Xu, Gaoyuan Zhang, Sijia Liu, Quanfu Fan, Mengshu Sun, Hongge Chen, Pin-Yu Chen, Yanzhi Wang, and Xue Lin. 2020. Adversarial t-shirt! evading person detectors in a physical world. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part V 16*. Springer, 665–681.

[229] Ying Xu, Kiran Raja, and Marius Pedersen. 2022. Supervised Contrastive Learning for Generalizable and Explainable DeepFakes Detection. In *Proceedings of the IEEE/CVF Winter Conference on Applications of Computer Vision*. 379–389.

[230] Ziwei Xu, Sanjay Jain, and Mohan Kankanhalli. 2024. Hallucination is inevitable: An innate limitation of large language models. *arXiv preprint arXiv:2401.11817* (2024).

[231] Mingfu Xue, Can He, Jian Wang, and Weiqiang Liu. 2021. Backdoors hidden in facial features: a novel invisible backdoor attack against face recognition systems. *Peer-to-Peer Networking and Applications* 14 (2021), 1458–1474.

[232] David Yambay, Priyanka Das, Aidan Boyd, Joseph McGrath, Zhaoyuan Fang, Adam Czajka, Stephanie Schuckers, Kevin Bowyer, Mayank Vatsa, Richa Singh, et al. 2023. Review of iris presentation attack detection competitions. In *Handbook of Biometric Anti-Spoofing: Presentation Attack Detection and Vulnerability Assessment*. Springer, 149–169.

[233] David Yambay, Luca Ghiani, Gian Luca Marcialis, Fabio Roli, and Stephanie Schuckers. 2019. Review of fingerprint presentation attack detection competitions. *Handbook of Biometric Anti-Spoofing: Presentation Attack Detection* (2019), 109–131.

[234] Xiao Yang, Wenhan Luo, Linchao Bao, Yuan Gao, Dihong Gong, Shibao Zheng, Zhifeng Li, and Wei Liu. 2019. Face anti-spoofing: Model matters, so does data. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 3507–3516.

[235] Xiao Yang, Fangyun Wei, Hongyang Zhang, and Jun Zhu. 2020. Design and interpretation of universal adversarial patches in face detection. In *Computer Vision–ECCV 2020: 16th European Conference, Glasgow, UK, August 23–28, 2020, Proceedings, Part XVII 16*. Springer, 174–191.

[236] Bangjie Yin, Luan Tran, Haoxiang Li, Xiaohui Shen, and Xiaoming Liu. 2019. Towards interpretable face recognition. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*. 9348–9357.

[237] Frankie Youd. 2021. Touchless travel: the introduction of airport biometrics. https://www.airport-technology.com/features/touchless-travel-the-introduction-of-airport-biometrics/

[238] Zitong Yu, Chenxu Zhao, Zezheng Wang, Yunxiao Qin, Zhuo Su, Xiaobai Li, Feng Zhou, and Guoying Zhao. 2020. Searching central difference convolutional networks for face anti-spoofing. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 5295–5305.

[239] Timothy Zee, Geeta Gali, and Ifeoma Nwogu. 2019. Enhancing Human Face Recognition with an Interpretable Neural Network. In *Proceedings of the IEEE/CVF International Conference on Computer Vision Workshops*.

[240] Matthew D Zeiler and Rob Fergus. 2014. Visualizing and understanding convolutional networks. In *European conference on computer vision*. Springer, 818–833.

[241] Matthew D Zeiler, Graham W Taylor, and Rob Fergus. 2011. Adaptive deconvolutional networks for mid and high level feature learning. In *International Conference on Computer Vision*. IEEE, 2018–2025.

[242] Qinglong Zhang, Lu Rao, and Yubin Yang. 2021. Group-CAM: Group Score-Weighted Visual Explanations for Deep Convolutional Networks. *arXiv preprint arXiv:2103.13859* (2021).

[243] Wenliang Zhao, Yongming Rao, Weikang Shi, Zuyan Liu, Jie Zhou, and Jiwen Lu. 2023. DiffSwap: High-Fidelity and Controllable Face Swapping via 3D-Aware Masked Diffusion. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*. 8568–8577.

[244] Yaoyao Zhong and Weihong Deng. 2020. Towards transferable adversarial attack against deep face recognition. *IEEE Transactions on Information Forensics and Security* 16 (2020), 1452–1466.

[245] Bolei Zhou, Aditya Khosla, Agata Lapedriza, Aude Oliva, and Antonio Torralba. 2016. Learning deep features for discriminative localization. In *Proceedings of the IEEE conference on computer vision and pattern recognition*. 2921–2929.

[246] Hongbo Zhu, Chuang Yu, and Angelo Cangelosi. 2022. Explainable emotion recognition for trustworthy human–robot interaction. In *Proc. Workshop Context-Awareness Hum.-Robot Interact. Approaches Challenges ACM/IEEE HRI, Sapporo, Japan*.

[247] Pengkai Zhu, Ruizhao Zhu, Samarth Mishra, and Venkatesh Saligrama. 2021. Low Dimensional Visual Attributes: An Interpretable Image Encoding. In *International Conference on Pattern Recognition*. Springer, 90–102.