# Toward a Hermeneutics of Explainability: Unraveling the Inner Workings of AI Systems

*By Srihari Maruthi[1], Sarath Babu Dodda[2], Ramswaroop Reddy Yellu[3], Praveen Thuniki[4] & Surendranadha Reddy Byrapu Reddy[5]*

## 1. Introduction

For all the enthusiasm and sheer volume of research in explainability recently, there is curiously little consideration of the interpretive theory that forms the backdrop of many of the proposed methods. Virtually all accounts of explainability presume that there are right interpretations into which researchers should guide a particular audience when asked to provide a clear explanation of an outcome. In this paper, we argue that it is time for technical researchers to look to the humanities and social sciences traditions surrounding interpretation (with roots in the work of Gadamer in the 1960s and the hermeneutic circle) in order to ground our explainability efforts in a more informed, critical, and self-reflexive context. Indeed, in doing so, we will shed a more critical view of what is likely a commonplace task for human researchers that should not be taken lightly even when machine-based support strategies are deployed. Our take is that the core ideas of hermeneutics provide a template for understanding the relationality of interpretive acts. With these ideas in hand, AI researchers should be able to reason more coherently, and with greater humility and sensitivity, about what interpretative acts mean, and about how we might design systems and support strategies that help to realize specific ends in interpretive situations.

Over the past several decades, the emphasis within AI on building ever more powerful and capable systems has largely overshadowed concerns with understanding in a deeper sense what those systems are doing to generate particular outcomes. As the field has increasingly turned to the mysteries of so-called deep learning, in particular, the push toward intricate, heavily non-linear system designs has largely eclipsed work toward developing more sophisticated ways of understanding how these sophisticated mechanisms arrive at their conclusions. Yet the inherent unpredictability and uncertainty of AI systems, disclosed in a range of recent high-profile examples that deploy these complex methods, have stirred broader caution concerning the role that AI systems should hold within society. In this

---

[1] University of New Haven, West Haven, CT, United States
[2] Central Michigan University, MI, United States
[3] Independent Researcher & Computer System Analyst, Richmond, VA, United States
[4] Independent Researcher & Program Analyst, Georgia, United States
[5] Sr. Data Architect at Lincoln Financial Group, Greensboro, NC, United States

context, the pursuit of explainability, to disclose the inner workings of AI systems and to make the process from inputs to outputs interpretable, has taken on an almost existential level of importance.

### 1.1. Background and Rationale

When behavior, especially unintelligible behavior, impacts society, addressing AI unconventional studies becomes necessary. Societal understanding of AI is predominated by research in the empirical sciences. Consequently, it neglects every single AI subjective aspect, especially those related to built-in systemic properties, eventually leading explainability assessment to facing a multitude of unsolved problems. Since AI payments to society are unaligned with theoretical studies, the usefulness of theory is also questioned. Given its subjectivity, the study of AI behavior remains controversial and challenging. This work is a theoretical study of AI explainability. Since AI lacks conceptual representation and fails in achieving corpus explicative, this paper revisits AI and defines further research directions to develop AI-doc for helping society in building conceptions of AI. Due to the complexity of the AI universe, this topic is deeply interdisciplinary and the hermeneutical approach to AI without a hierarchical integration of sub-disciplines will be a futile put-on. The paper first revisits AI to define a conceptual landscape that is in conjunction with the processes that permeate AI. The occurrence of explainable behavior is accepted by AI artifacts. Indeed, this study investigates candidate AI explicative process to study to unravel the inner workings of AI as it faces society. Since AI has grown out of the natural philosophy of perception, a comprehensive approach to AI implicitly affects a plurality specific field belonging to the human knowledge either implicitly or explicitly related to explanation. Thus, the revision AI is studied according to a maintenance hermeneutic of its internal morphology, as the model is perceived in various scientific schools, in a comprehensive study that the contribution can help identify the conceptual tropism to insert AI inside a cultural tradition. Furthermore, AI explication needs both an etic and entic perspective to understand the model extrinsic and intrinsic regularities.

In the era of the fourth industrial revolution, innovation converges through technological advancements. The most characteristic of them concern the computerization of virtually all human activities and the associated advancement of artificial intelligence (AI). Understanding AI raises broader problems even when AI is developed based solely on human input and influenced by culture. Indeed, one can claim that it is derived from society and its own history, and is deeply linked with human society. Formalizing human cognitive abilities has been central to the history of AI from its inception. Consequently, implicitly or explicitly, AI is considered as an entity or even an imitation of human thought. Supporting society in uptake of AI requires addressing those problems in a tangible, responsible, and transparent way. Achieving an equal consensus definition of AI has been also shown

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

to be extremely difficult, let alone formalizing AI's behavior. Therefore, AI raises wide-reaching and society-felt subjective problems that must be addressed.

### 1.2. Scope and Objectives

To achieve the scope and objectives, we choose a communication-centered interpretative research method that adapts and integrates the outlines of ethnographic or "local interaction" method, "cultural" method, and the novel visualization of software architectures for AI systems operating in information-rich environments. We call this method AI hermeneutics and the various steps in the method constructionist hermeneutics, ethnographic hermeneutics, circulatory hermeneutics, and resistant hermeneutics. Importantly, the method focuses not only on understanding but also on the ability to interpret a domain of social action (specifically AI explainability) that is characterized by operational explication of AI-generated facts and their contextual relationship to the mission and objectives of intelligent systems.

The scope of this article is more modest. To paraphrase, the first objective is to begin to unravel the inner workings of key state-of-the-art AI systems and to assess their contributions to overall AI behavior. This includes a detailed examination of the inner workings of knowledge representations, inference and explanation modules, as well as their supporting architectures. The second objective is to begin to lay the predictive and evaluative research foundations vital to the development of stable classifications of various types of AI explanations in these and other AI systems. Our imperative is expressivist - to clarify and distinguish the different kinds of explanations that these AI systems can provide. In situational use, different principal-objective situated goals require different types of explanations. These support the 2-way communication loop between the AI system and the principal in the operational setting. Our intention is not regulative - to prescribe what kinds or what level of detail of explanation should be given in a given situation or for a principled AI more generally. These normative and regulative rules embody several types of ethical and liability considerations - potential harms; moral, legal and social responsibility; credibility and trust creation and maintenance; and legal, statutory, and regulatory disclosure requirements. These, too, require predictive research.

### 2. Foundations of Hermeneutics

### Hermeneutic Circle

According to Heidegger, every exegesis is tormented by a fundamental circularity, the one derived from the a priori of human existence, with regard to the worlds that stand before the projecting power of human action. This a priori is inherent in the distinction between ready-to-hand and present-at-hand

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

which makes it possible to interpret a thing or an entity in the world: humans can handle tools only if they will project a scenario of future situations in which the usefulness of the tools themselves will reveal itself. The Hermeneutic Circle is an inseparable part and parcel of the human interpretative horizon: the world exists for humans only if it is in some way interpreted, decoded, managed and planned. Only if in some way, the future makes a difference in the present, indeed, it can guarantee the stability and the cartography of our interpretative horizon.

**Heideggerian Hermeneutics**

The influence of Martin Heidegger on philosophical hermeneutics is perhaps unequaled. One of the core elements of Heidegger's philosophical activity is the critique of metaphysics within the context of a philosophical reflection carried out under the commitment to historicity. With his work, Heidegger reintroduces the fundamental importance of hermeneutics as an interpretative dimension of human behavior.

"If we do not explain what wisdom is, we have no right to pretend to use it."

**Key Concepts and Definitions**

Notable, and closing in with more emphasis on philosophy than technical performance, XAI differs from transparency, something we wish to emphasize because philosophical elements, and having clarity about model representations, are at the heart of XAI. More than describing system mechanisms, transparency formally includes details about model structures and weights. These can risk proprietary information, due to interest competition, and be used nefariously. Anything nefarious invites regulation and public oversight. Instead of transparency, XAI offers revelation: the knowing or showing of something that has previously been secret or that is unexpected.

Explainability (XAI) refers to the specification of machine learning system models in a transparent, understandable, human-interpretable manner. This differs from interpretability, which lies with the human. While interpretable systems imply that people can provide high-level explanations for why models display particular properties or responses, for the public, XAI emphasizes the ability to describe system models so that machine predictions in everyday applications are easier for humans to understand, and model mechanisms and behaviors are revealed with local, high-level insight. Explainable systems are at the heart of the data sciences, as they provide insights and understanding about communicated results, no matter whether the recipient possesses expertise. Prior mechanistic insight offers caveats and limitations to be kept in mind while managing contemporary predictive system use. Such cognitive reduction is quite timely now because society's ability to produce highly

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

interpretable models is being greatly outstripped by the volumes of observations and sophisticated algorithms available to understand mechanisms.

**Historical Development**

Recently, an academic AI gold rush to extend the state of the art of machine learning's object knowledge of pattern recognition of large unstructured data sets. These are deep learning, classifications, clustering, probability, data mining, reinforcement learning, and falsification, based on neural network architectures constructed from multiple layers, extending massively parallel backpropagation of errors during training. Deep learning can incorporate supervised, unsupervised, semi-supervised, and reinforcement learning, which have origins in statistics, neuroscience, genetics, cognitive sciences, and learning theory. The filtering, categorization, and prediction that ride on the shoulders of hidden structures are becoming an indispensable tool for knowledge discovery, knowing, and acknowledgment.

The idea to interpret the interpretation or explain the explanation is an expression of a much larger epistemological grand record. Traditionally, philosophies of explanation focus on the subjective psychological act engaging in explaining complex phenomena to peers in an intelligible manner. The 'correspondence view of truth' norms information theory, expressing how explanations express facts about the world. Philosophies of explanation now also live in the variation cognitive sciences, critical theory, and hermeneutics, which scrutinize the rationalization of AI systems. General Agreement builds around the meaning of 'black box,' disclosure communication, its role and impact in the domains of science, and society of machine learning algorithms, featuring expert systems, neural networks, and decision trees.

**3. Explainability in AI Systems**

Process is essential in the context of interpretability. A careful review of the literature will lead to the conclusion that interpretability is mostly confused with explanation. Most communities related to AI believe that interpretability is required by the decision maker in order to make a choice, or by the legal framework on rights about the explanations of automated decisions. This leads the research community to consider that the interpretable system is the one that employs a method able to explain the outcome. We have already noted that the sensors of interpretative status. Given that the interpreter can be characterized at three critical levels of abstractions, that constitute together the final definition of interpreter, the primary research problem to solve may be phrased as follows: how to design a suitable interpretative mechanism that could serve the notion of explainability. We do not believe the present paradigm that considers explanation in learning models to answer specific queries related to machine

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

learning models or predict future possible behavior actually scales to a scenario beyond a learning model. We define this as the central question of AI Interpretability due to the fact that we are confronted with the duality of human-machine intelligence.

We carefully surveyed generic case studies in order to support our claim that every interaction that is commonly referred to as explainable, even in the literature of Explainable AI, achieves explainability only on the surface. In most of the case studies, explainability is guided action; the interpreters were tailored by the modellers to follow a pre-planned route that is superfluous to what the interpreter might think. The traditional notion of explainability considered the explanation to be a delivery — information is handed over by the explicable to the interpreter. And indeed, this interpretation defines cognitive effort of the explicable, that is, the cognitive effort of the software that is achieving, guiding, tailoring or delivering what will be in the hands, perceptive sensors, abstract representations, and interpretative mechanism of the human interpreter.

### 3.1. Importance and Applications

Some researchers and practitioners use the terms interpretability, transparency, and explainability interchangeably, while others prefer a technical distinction, noting individual philosophical and practical contributions that particular notions emphasize and operationalize. In this report, we adhere to a distinction suggested by guidelines developed for the Department of Defense. Specifically, interpretability is the degree to which a human can understand the cause of a decision. Transparency is the degree to which a human can understand the algorithmic and processing steps used to render the decision, and explainability is the degree to which a human can understand the course of an autonomous intelligent agent to reach that point. These definitions enable greater focus and precision when formulating research questions and integrating insights.

In the U.S. Department of Defense, interpretability, transparency, and AI safety are key objectives to ensure that humans can understand and appropriately trust and use AI applications to support the Department's missions. Progress toward these goals will enable operators to manage AI systems more effectively even as their capabilities grow and become more heterogeneous.

Over the past two decades, a constellation of technological advances has enabled the development of AI systems with a wide variety of desired capabilities. These systems automate tasks, perform with super-human accuracy, and exhibit complex, robust, multistage behaviors. As these AI systems continue to advance in their capabilities, the gap between desired user-level functionality and human interpretability of AI system-level processes is widening. This has a range of practical implications for safety, transparency, trust, and utility.

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

### 3.2. Challenges and Limitations

The transparency and hence, in part, also the ethical acceptability of processing target data is an important aspect of privacy legislation. Efforts to enhance privacy should lead to domains employing data processing in a less invasive and more controlled manner. It is clear that these points can equally be applied to data processing by AI systems. We believe this never-ending quest for the explainability of AI systems will also promote trust in AI systems. Furthermore, breakthroughs in the development of these systems could bring forth benefits to humanity in the form of facilitated, though less invasive and privacy-protecting, processing of personal data of individuals. The concept of Explainable AI is, despite its current limitations, therefore an ethical concept.

The phenomena of mentally capable machines in today's society are almost entirely hidden behind the black curtain of their inner workings. We glimpse only sporadically the actual behavior of artificial intelligence systems that accounts for their at least superficial plausibility. In most cases, the particular inner workings of AI systems are sufficiently opaque that they cannot be readily examined, probed, and understood by others. Will we still be able to perform tasks related to our traditional human concept of thinking, assisted by AI systems, even though we no longer know the actual mental states of these systems? We argue that even though this might be achievable, we should strive for the greatest possible transparency and explainability of AI systems.

### 4. Philosophical and Ethical Considerations

This section explores the features and requirements of explainability and transparency, as these concepts could indirectly reflect certain ethical values, such as privacy or the right for human beings to share elements of rational thinking. As a critical issue regarding AI, this section addresses an issue of general interest not only for experts and stakeholders in AI, but also for society as a whole, as AI deployment transforms the norms and practices framing the right for an explanation. Its growing role in our life implies that its common benefits be maximized and its potential hazards minimized since AI can paradoxically lead to ethical louse failures.

At first glance, the questions related to the explainability and transparency of AI systems appear to only impact their effective use and safety. However, explaining AI tools raises more fundamental questions about applying them with wisdom. There is also a philosophical and ethical dimension to the way that explainability will be reached – if it is ever to be reached –, even if these requirements are usually unvoiced explicitly in the legal texts and the technical challenges overcome. As a matter of fact, use of the AI will materialize in concrete contexts and involve multiple implicit or explicitly shared societal values and benefits.

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

### 4.1. Interpretability vs. Transparency

In this sense, we revisit Burrell's typology of AI explainability: transparency, understandability, and persuasive. Goodman and Flaxman propose a fourth category of "manipulability," focusing mainly on the ability of the user to interact with the lessons drawn from the interpretable model. We argue that manipulability is actually part of the transparent AI system. Unlike Goodman and Flaxman, we directly address the underlying model complexity versus transparency debate using traditional metrics from multiple fields—statistics, computer systems, operation research, decision science, marketing, etc.

How does interpretability relate to transparency? Goodman et al. claim that "When AI is transparent, that is, understandable or predictable and capable of being interrogated, its functionality can be fairly judged. Interpretability is only valuable when transparency cannot be fully achieved." This implies that a transparent AI system is, in consequence, interpretable. However, to be explicit, while "transparent" is used to indicate that "models are understandable and capable of being interrogated," "interpretable" "engages with a particular audience and serves a particular purpose." We thus believe that interpretability is a unique, supratextual variable that is flexible, malleable, and designed according to the audience and the elicited purpose.

### 4.2. Bias and Fairness

If prediction quality becomes another force for human disempowerment in a period of rising work and social insecurity, then a backlash seems inevitable. Digital technologies, and big data in particular, have recently been fingered as culprits in discussions on the future of work. Turn the attention from automation to data and its hypes, AI and machine learning, our idols of the last years, do not look so benign. Much of the debate in economics has hitherto focused on the performance of AI and machine learning models according to their predictive competition on given data. Predictive model competition has dominated AI and machine learning conferences and contests since the mid-1990s. Now, QA-compatibility has come back into fashion as helping users to trust and understand models has become increasingly important. With acknowledged studies of AI systems, and not just curiosity about achieving machine intelligence or making profit out of predictive models. The increasing complexity of AI systems is entering the field by the front door and accounting for it is nothing short of essential.

Bias and fairness. One hot topic of late is the question of whether the prediction made by an algorithm introduces or replicates biases. Algorithms for hiring, for scoring products, to judge risk in lending, etc., are all subject to important discussions about potential biases. In general, the predictions of models can be unfair to specific individuals, and it can be debated whether algorithms that move decisions away from human agents (who might be biased) bring about more fairness or not. Substantial effort

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

has been dedicated to testing whether models replicate or fail to replicate patterns that would normally be associated with bias in a decision-making process. These fairness tests are laborious, and a more general understanding of the effects of the algorithm on the ways of decision in the mind of other agents could potentially identify a wider range of predictors of simple notions of unfairness. Such conclusions should resonate well with proponents of the original interpretable models. The objective is entered in a number of fairness situations using LIME in a dedicated fairness study. The LIME method was used in order to get an idea of where the fairness problems lie and look into the coefficient signs. These results were quite illustrative too.

## 5. Methodologies for Interpreting AI Systems

Both computing systems and minds, each in idiosyncratically different ways, produce what could be broadly construed as rational inferences. In the case of the human mind, rational analysis promulgates the "good reasons" hypothesis to explain away psychological actors attributed beliefs and intentions. To explain the behavior of a rational entity - in fact, the notion of explanation itself with its typical forms and styles of narrative and resort to counterfactual events and conditions - depends upon the fact that the entity is rational and possesses such a faculty, independently of its biological embodiment. Underway, while dwelling on the deep question of how we discover in the data of any source the principles that are being used by the AI. Finding that the AI learned some simple function of the data may not be profound, especially if that function was weakly validated and over-trained on spurious or irrelevant "patterns".

The long-standing tradition of work in hermeneutics and interpretive social science can contribute to some cognitive clarity and precision about the way interpretability should be best understood, conceptualized, practiced, and critiqued as it specifically applies to AI systems. The philosophical and social science literature makes several distinctions and draws clear lines among such related concepts as explanation, understanding, reasoning (which is an aspect of computing at issue in AI research) and interpretability, or the enabling of explanations and meanings. A genuine practice of interpretability requires that there should be, or that we should aspire to one. However, the renowned "Duhem-Quine" thesis of underdetermination of theory by evidence also applies to interpretability, and therefore we must, as Garfinkel advised, proceed "documentarily" and specifically.

### 5.1. Model-specific Approaches

Concerning the processes of model structure and learning, two different branches of explainability research have evolved. Both of them suffer from an inherent trade-off between prediction quality and explainability. The main learners presented in this chapter strongly rely on entirely impenetrable black-

boxes, such as deep neural networks. They leverage the excellent performance of these models to solve challenging prediction tasks in domains such as computer vision, natural language processing, and speech recognition. According to the trade-off curve, they are computationally more costly than unexplainable models at a high-capacity model.

Model-agnostic approaches explicitly address the human-level explainability of a technique, thereby providing a model-independent explanation of a model's prediction. I deal with the model's structure and the process of learning the model separately. This separation is inherited from statistical learning where a process of hypothesis testing and model evaluation is used to select the model from a typically infinite hypothesis space.

An interesting line of research is XAI approaches based on salient maps. For computer vision tasks, such methods often outperform perturbation methods. Different from model-specific ones, perturbation-based explanation approaches search for an explanation by evaluating the predictiveness drop of a model when input features are perturbed.

An important category of explanation approaches directly addresses the specific challenges and opportunities associated with the inner workings of AI systems and offers model-specific methods to make predictions of a wide range of AI algorithms more understandable. For example, XAI methods modify the architecture of the model such that it becomes transparent to a human user. LIME (Locally Interpretable Model-agnostic Explanations) is an example of a model-agnostic CAM. Given any machine learning model, LIME provides an explanation for any individual prediction by learning an interpretable model locally around the prediction. This approach provides a human understandable explanation for any prediction at the cost of a considerable computational overhead.

**5.2. Post-hoc Techniques**

Supporting the AI system to provide explanations after it has already made decisions is also a continuing topic of research, involving post-hoc Explainable AI solutions. Conventional mechanisms are mainly model-agnostic and benefit from a broad variety of invisible glass boxes. This, however, is the low pathway to explainability as the genuine workings of a black-box system cannot be exposed. This paper provides a review of post-hoc techniques in the area of Explainable AI. It seeks to raise awareness of several key post-hoc problems, as well as rising points, pointing out potential avenues nested inside Explainable AI. Such a hermeneutical basis should sustain alleviate the Easter bunny effect and decrease reliance on invisible rabbits, directing the conversationist of the duck that needs to recognize and understand the intrinsic operations of the AI system.

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

Post-hoc techniques, also known as surrogate models, provide a new simpler model which approximates the underlying black-box model. This simple model is known to be interpretable, thus the integrity of the surrogate can be evaluated and knowledge can be gleaned from the results of the surrogate. The actual model can then be deemed interpretable based on the characteristics of the surrogate. Surrogate models work to reduce the complexity of the underlying black-box model into a simpler interpretable model, and expose the inner workings of the AI system. Given this model that would be questioned afterward, this post-hoc set of methods can be integrated into an Explainable AI system. Post-hoc Explainable AI can be applied to many branches of machine learning, as support for post-hoc techniques may be directly integrated within the machine learning algorithm, or as a wrapper considering local explainability is embedded within the process.

## 6. Visualization and User Interfaces

While visualization is a powerful tool for opening the black boxes of neural networks and providing context and immediate value for the users, it is not a panacea. The risk scenarios, introduced by Mediaec, are helpful when making a decision to trust the system or not. The overreliance on false-precision decision-making and confidence in AI and automation has been observed to have a risk of declining human proficiency. There is a possibility that visualization will confuse the users further by adding an extra layer of technical introspection, so beware of verificationism.

Apart from the neural network model debugger based on visualization, other visualization techniques have been developed for human inspection of the internal workings and decision making of neural networks. The TCAV algorithm generated aligned direction in input space to support interpretation of predictions in the network via visualization. The extension to this approach, Quantitative Direction Method (QDM), was introduced to provide the user with better control by centering its visualizations around a user-designed concept or a priori-aligned directions. Beyond directions, Heatmapping and Occlusion Maps were introduced for human-interpretable feature weights and attention within the CNN and LSTM for enhanced context-aware image and video captioning.

### 6.1. Types of Visualizations

It is also common to distinguish between internal and external visualizations, based on whether the data to be visualized comes from inside or outside the computational model, respectively. In empirical sciences, external visualizations, or scientific visualizations, have been used for decades. They are particularly emblematic in physics, climatology, biology, and neuroscience, and are often the type of statistics and data visualizations AI practitioners refer to when deploying visual explanations, since AI practitioners most often produce visual explanations from externally interpreted feature attribution

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

methods. The ubiquitous usage of visualizations in empirical science often goes unnoticed in AI conversations about AI explainability, perhaps because the type of visual explanation to be achieved is order of magnitude more complex than appealing to the philosophical argument for visual explanation validation.

The term 'visualization' is rather broad, and there have been many attempts to classify visualizations arising out of different tasks. Card et al. defines six distinct types of tasks which visualization techniques serve, each with a unique design myth: overview, zoom, filter, details-on-demand, relate, history.

### 6.2. Human Factors Considerations

Overloading users with too much information, however, can lead to confusion or reduce compliance. Knowing when help and guidance are important is an essential human factors lesson, and ATB provides a setting in which this becomes salient. For power users, overload may be a particular problem, especially if explanations span disparate user communities, where the same explanation might be necessary multiple times. Balancing the design and delivery of these explanations therefore necessitates a good understanding of user needs. Finally, if an individual is identified as unreliable, they may be left out of the provision of program support and are unlikely to receive explanations. Such individuals or groups may be marginalized in society, which is particularly harmful as people rely increasingly on automation.

When building explainable systems, it is important to keep people in the loop, forefront, and as the primary audience. End users are the ones, after all, who need to make sense of, understand, trust, and review automated tools. The literature suggests at least two primary factors that can be adjusted to increase satisfaction with AI system explanations and the perceived helpfulness of the underlying explanation techniques. This includes grounding explanations in individual or collective belief and increasing the overall perceived competence of the system, often by explicating the roles and responsibilities for each possible user. These factors are essential for improving both the models themselves and the system interfaces, however, and are far from a complete list of necessary human factors adjustments.

### 7. Case Studies and Practical Applications

As Haraway, Ester, and many other theorists of the artifact have suggested, our goals can indeed be achieved only if a dialogue is engaged between the artifact and the interpretation which describes it. This study is an attempt to respond to that call, and the proposed critique of explanation of artificial

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

and human artifacts is therefore a first step. We have analyzed the public statements by some of the pioneers. We have also tried to pinpoint their underlying assumptions with regard to society and technology, drawing inspiration from both contemporary and classical hermeneutic theories. Not surprising, we discovered that a certain hermeneutic blindness which was blinding them. Their study, either by omission or in the abstract, hides its social underpinnings, and product concealed as well its latent model of society. It therefore became necessary to unmask this model before examining the criticism in more depth. Admittedly, uncovering explanations and models is already a first step in this direction. First, the model of society being presented and a possible critique of society may already provide material for a deeper, more general critique of artificial artifacts. Second, it will enable the artifacts themselves to enter into the dialogue and take part in the process of deconstruction.

On the theoretical level, the hermeneutics of technical artifacts is a contemporary variant of the ancient hermeneutics of natural beings. It differs from it primarily through its abandonment of anthropology and through its proscription to the domain of specifically human understanding. The technique of interpretation has been replaced by a technique of reading, and in contemporary hermeneutics the artifact plays the role of an ally of the reader. For the most part, the hermeneutics of technical artifacts is concerned with interpretation. Its special task is to ascertain the inner workings of the artifact, or more specifically, its principle of interpretation. In this study, we have transformed this task into a critique of interpretation, because only thus can a critique of the models applied by the artifacts be targeted as well. As a result, the specific nature of the hermeneutics is minimized. It becomes a tool in the service of a more general hermeneutics which is concerned with interpretation itself.

### 7.1. Healthcare and Medicine

Explainability of medical AI is intrinsically high stakes but is not simply a matter of accuracy versus the imperatives of fairness, transparency, and accountability. Medical AI practitioners require a more nuanced understanding of explainability that reflects the more extreme power dynamic inherent in their sector. AI will very soon become wholly dependent on explainability for acceptance, adoption, and endorsement. The risks associated with failing to deliver this reflective benefit are not restricted to any one segment of society. On the contrary, they pose existential threats to many. Indeed, we are dealing with more than explanations required to allow the challenging coexistence of people and AI. We are dealing with broader societal norms, values, meanings, purposes, and reasons.

Society is right to be excited by the potential for the transformative power of AI in the area of healthcare and medicine. Here, AI systems are being used for applications that range from diagnosing disease to discovering new medicines, to the creation of patient-treatment plans. However, with this seismic shift, also come challenges. Top of the list, for healthcare perhaps, is that of social acceptability. How will

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

patients and their carers, clinicians and scientists respond and adapt to AI systems which they do not and cannot fully understand, despite the assurance of their capability and reliability from an engineering perspective? What happens when AI systems used for disease staging during radiological examinations, for instance, make suggestions during clinical decision-making without explanation or justification? When used to assess postsurgical progress, how will patients react to having AI doctors that are accepted as lawful holders of clinical knowledge but lack any communication or explainability skills?

## 7.2. Finance and Banking

A "sayable" AI system can consequently be based on some existing methods that follow different approaches characterized by how they generate explanations, that is, how they communicate over what the neural network has learned and what it is using to make its neural predictions. Such methods are either model-specific or are independent of the application while being concerned with model-specific information only. Following the latter strategy, one approach to make a pre-trained model explainable is to understand the internal mechanisms of the model and to produce explanations after interpreting these mechanisms. To do that, interpretable model architecture is needed that generates near-approximations for structural interpretation, like, for example, convoluted feature maps and pooling acts in CNNs, or hidden layers of LSTMs.

In finance and banking, AI systems like CRMs, fraud detection systems, or robo-advisors are built and put in place by comparing various sets of historical data to each other until the AI system finds a pattern or a significant fact. The goal of a fast explanation for the AI system's decision is, for example, to understand why and how it made that decision in order to make it auditable for regulatory purposes, to reassure the customer about privacy or fairness issues, to determine whether the system exploits biases in the dataset, or to cope with and manage human interactions. The association between empathy and explainability is most visible in a robo advisor, which is a digital investment advisor that utilizes algorithms to automatically allocate, construct, and manage clients' investment portfolios using passive indexing strategies. These advisors are designed to deliver stock market returns similar to indexes globally and to execute trades for registered discretionary portfolios.

## 8. Regulatory Frameworks and Standards

Regulatory frameworks and standards for AI-based systems are under development. They aim to empower individuals to understand, contest, and modify system-generated outputs. And for the automated systems to, in turn, be accountable. At a high level of understanding, the explainability of

**Journal of Artificial Intelligence Research and Applications**
**Volume 2 Issue 2**
**Semi Annual Edition | Jul - Dec, 2022**
This work is licensed under CC BY-NC-SA 4.0.

AI systems is necessary to ensure respect for democratic and human values since it should support the identification and fixing of model biases.

AI and machine learning in the context of the intelligible system. This includes mainly the public sector, where the use case for these technologies is such that they require high levels of transparency and understandability for specific individuals and groups, where also the social impact of the decision systems should be evaluated and acceptable. Nonetheless, the explainability of AI/ML models is a key issue also for the private sector. Actually, the requirement of providing explanations in turn pushes towards the development of models that are intrinsically better tailored to extract knowledge from data and to datasets that are good representations of reality.

### 8.1. Current Landscape

Since the mid-2010s, there has been an escalating, even feverish quest among both corporate and academic researchers to make AI models more transparent and comprehensible. The motivations driving this quest are substantial and varied: devising models that perform better in error-critical domains, fostering consumer trust, answering to legal requirements, and allowing different members of society to understand and to participate in decision processes in which they are involved or that otherwise affect them. Much of this quest has been expressed through the concept of "explanations," which are frequently characterized as belonging to one of two types: so-called "post-hoc" explanations derived from a trained model that tell how, or what parts of the input influenced, the model's decision (e.g., a heat map over a body scan to indicate why the model made a particular health-related inference); or simpler, inherently interpretable models working in tandem with, or as a proxy for, a black-box model.

### 8.2. Future Directions

The interpretive connections of an AI system can be regulated in various ways. One manner in which that process may be engineered is through formal methods that guarantee that another AI, or actually a group of AI, infers the proper interpretations. Bayesian networks, graph-based, and constraints-based classifications are the main categories of these formal methods. Notably, in the latter case, the interpretation will be based on a group Bayesian network. Another rival option, the one considered by authors like Salehi on decision trees, concerns patterns. To supply human-like interpretations, this sort of regulating method can encase the inside of AI in many shrouds on the outside. Furthermore, the similarity between the modification and the notion of reinforced conformity suggested by Foucault may be another possible avenue. The development of an ethics framework for explainability should augment the process or the analysis of robotic hermeneutics willing to endow some consequences to

the control dynamics induced by ethics. Finally, it will be important to underline where AI decision outputs need not to be transparent.

In this paper, we have provided an overview of what we term a hermeneutics of explainability. Whereas others focus on surface relationships, on the means of giving an interpretation of the inner workings of AI systems or on the ways in which humans come to understand outcomes of AI systems, we also underscore the verticality of our model in that interpretations may arise through interpretive interactions over time. We have also identified three new pathways informing contemporary AI research that are reframed by our model. But our analysis also highlights a series of complex questions that bother the peaceful waters of many works yet to come. For example, when and how to address the gap between human and AI interpretability; metrics; cognitive biases; new legal considerations vis-à-vis explainability; the role of ethics; or hackability.

## 9. Conclusion and Future Directions

Where next? An initial step is to consider how the work we have described in some kind of evaluation of explanation can generalize beyond the deployment context in which our identified connections focus. Secondly, we separate our connections as individual links to be unraveled. This work has been successful in building a pluralist and interdisciplinary perspective on explainability, but to move forward we need to consider how these are actually connected to show how and when they matter. We also need to consider and distinguish the different kinds of consequences of failing to meet these connections, which might have particular implications for the different societal norms and institutional practices of the deployers that we have identified. Yet in trying to paint a nuanced picture of the importance of understanding and transparent inner workings, we raise the question of where to place the threshold.

What we argue for in this paper is a step toward a "hermeneutics of explainability," driven by the position that understanding the inner workings of a model is not simply a form of explanation, but it should be a primary goal in AI systems. In particular, we identify what we see as critical connections between understanding and the potential for trust, credibility, human AI interaction, and interpretability. It is through the surfacing of these connections that AI researchers can gain better understanding of the needs of explainability that face them. This is particularly the case when it comes to AI systems being deployed in contexts where the stakes are raised, where decision-making is consequential, and where the implications for society have significant impact. There is of course plenty more to be said about these connections, such as how they manifest within different stakeholders and their differing needs, the recognition of which is necessary to develop AI systems which are truly trusted.

**9.1. Summary of Key Findings**

Given the ethical concerns that we previously detailed, we believe that hermeneutics could further investigate explainability methods and results. Our analysis focused on the study of the predictions of deep neural network models, trained on several datasets of retinal images for the binary classification of pathologic retinal diseases. We studied the attained prediction probabilities as model results, paying specific attention to images for which model predictions were poor or completely wrong. We explained the outcomes using Layer-wise Relevance Propagation, whose principles gave insights about the semiological relevance of input features in the obtained retinal image population. We also planned and conducted a clinical study which went beyond AI-driven technology effectiveness on small-scale clinical application, dealing with further aspects such as congruency. The final goal was to devise grounded strategies to avoid potential AI bias and misunderstandings in a wide-scale clinical deployment of AI explanation and prediction model results.

Given the increasing development and diffusion of AI systems into applications for various domains, including AI in medicine, their inner workings are a topic of debate. Several authors have dealt with the relationship between results stemming from explainability methods and their inference in terms of the development of trust, which is crucial in real-world applications. Our contribution to the debate stresses the necessity to deepen the study of the results coming from AI models, once they have been obtained and fully assessed in terms of their accuracy, generalization, and error. This is due to the potential impact they can have on our understanding of the examined phenomenon, the theoretical setting of the bio-mathematical model, the relevance of input features in the context of AI-driven real-world applications, as well as related ethical issues.

**9.2. Research Agenda**

In calling for both a broader notion of scientific explanation and a more restrictive notion of inference, elucidation, and all the manner of neo-positivist philosophemes – that we should be content supplying mostly a-priori constraint and guiding principle for our scientific algorithms – we ask what it is about our explainable models that may have a greater or unique use, and a configuration-distillation function that is yet to be discovered or leaves some or most thought-lab experimentally un-parse-able. This holds even when the explanation can take the form of an approximate reductionist solution that is transformable into even our simplest empirical models with a minimum of effort in the face of AI demagogue-ism.

In moving toward an ambitious and productive research agenda for the way in which we characterize, evaluate, and logically relate models of explainability, we suggest that it is productive to clarify the

"σὸλος" of explainability. In separating the idea of explainability as a characteristic of scientific models from the models themselves, we can come up with both logical relations and necessary conditions of each. In proposing that explainability can be fractured from the idea of AI systems as scientific models--hence removed from notional first-to-third (notional) or to-first (algorithmic) person status-–we can suggest a necessary re-analysis of non-empirical aspects of scientific models. As such, the how and what funny distinctions should have some epistemological structure or content. In asking what it is about our models as models other than tractable inferences and discoveries that can be used or re-purposed to explain them?