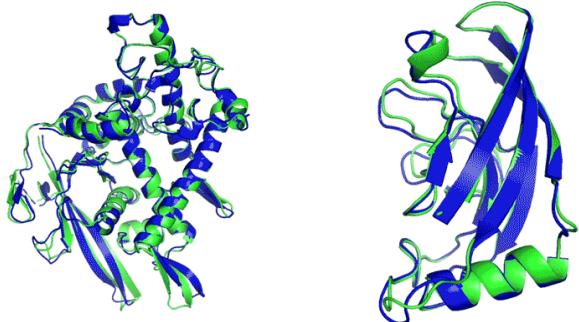


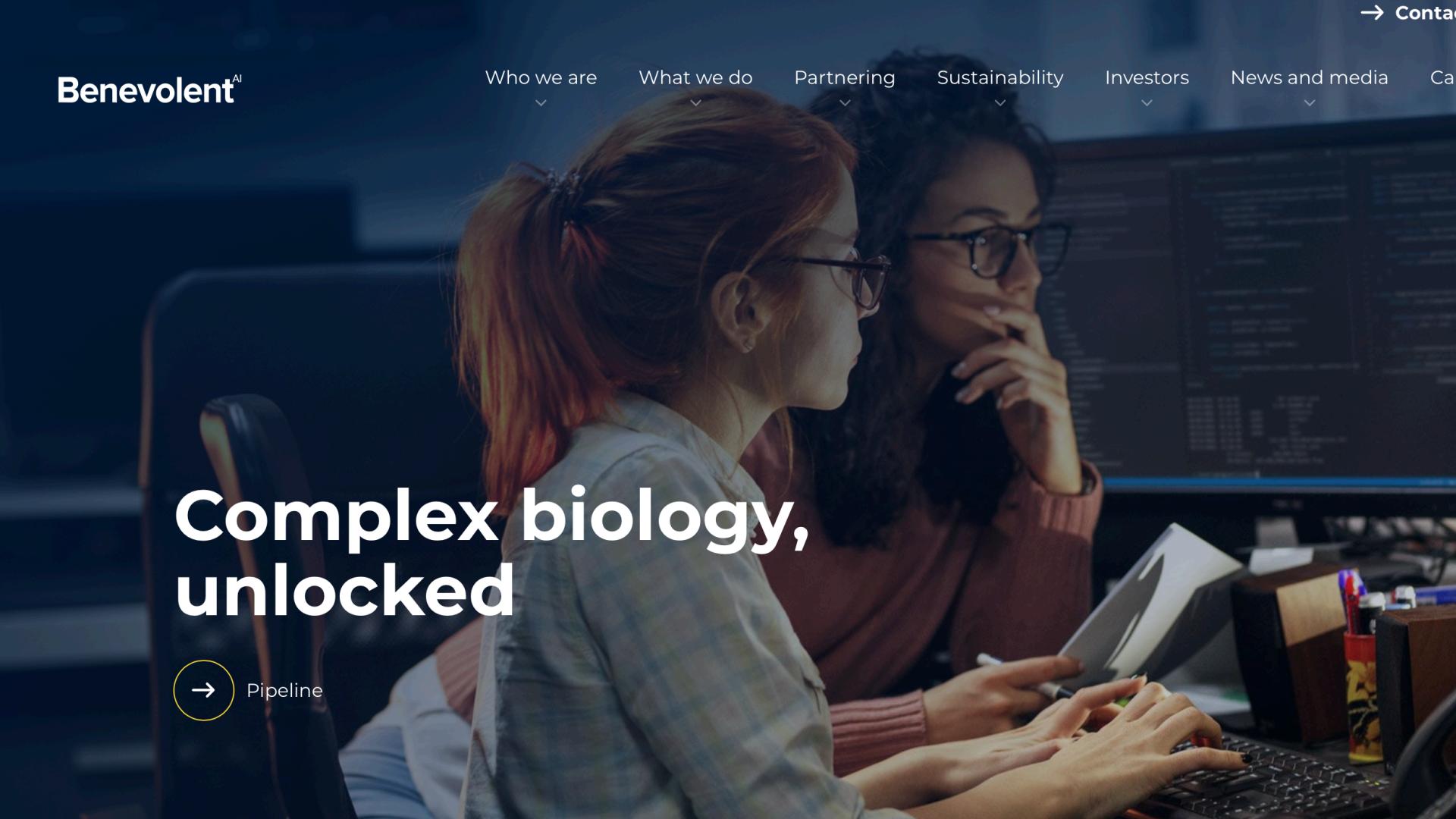
# Justification, transparency and computational reliabilism



T1037 / 6vr4  
90.7 GDT  
(RNA polymerase domain)

T1049 / 6y4f  
93.3 GDT  
(adhesin tip)

- Experimental result
- Computational prediction

A photograph of two female scientists in a lab setting. One scientist with red hair tied back is in the foreground, wearing glasses and a plaid shirt, looking down at a laptop. Another scientist with dark curly hair and glasses is behind her, also looking at the laptop screen. They are surrounded by computer monitors displaying code and data. The lighting is dramatic, with strong highlights on their faces.

# Complex biology, unlocked



Pipeline

# Wu & Zhang: Automated Inference on Criminality Using Face Images



(a) Three samples in criminal ID photo set  $S_c$ .



(b) Three samples in non-criminal ID photo set  $S_n$

Figure 1. Sample ID photos in our data set.

# Wu & Zhang: Automated Inference on Criminality Using Face Images

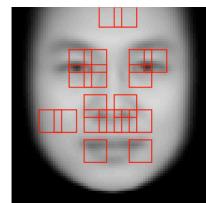
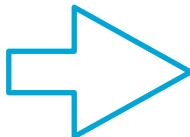


(a) Three samples in criminal ID photo set  $S_c$ .

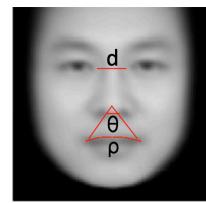


(b) Three samples in non-criminal ID photo set  $S_n$ .

Figure 1. Sample ID photos in our data set.



(a)



(b)

Figure 4. (a) FGM results; (b) Three discriminative features  $\rho$ ,  $d$  and  $\theta$ .

# Wu & Zhang: Automated Inference on Criminality Using Face Images

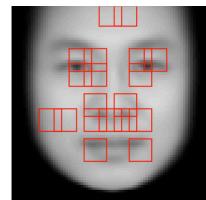
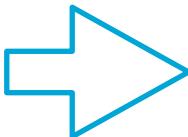


(a) Three samples in criminal ID photo set  $S_c$ .

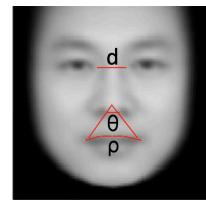


(b) Three samples in non-criminal ID photo set  $S_n$ .

Figure 1. Sample ID photos in our data set.



(a)



(b)

Figure 4. (a) FGM results; (b) Three discriminative features  $\rho$ ,  $d$  and  $\theta$ .



(a)



(b)

Figure 9. (a) The four subtypes of criminal faces; (b) The three subtypes of non-criminal faces.

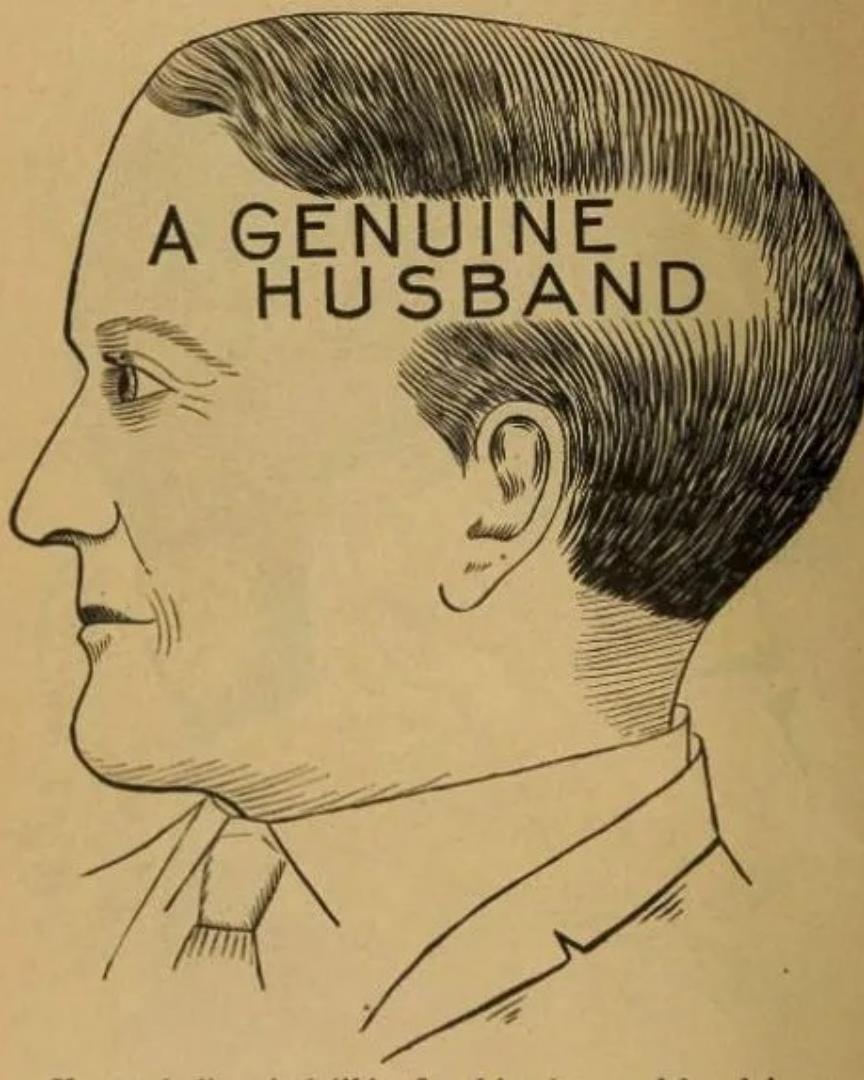
# Wu and Zhang

"The two manifolds formed respectively by the two data sets of criminal and non-criminal faces are concentric, with the manifold for the non-criminal face images lying in the kernel with a smaller span. This newly discovered knowledge suggests **a law of normality for faces of non-criminals: Given the race, gender and age, the faces of general law-abiding public have a greater degree of resemblance compared with the faces of criminals.** In other words, criminals have a significantly higher degree of dissimilarity in facial appearance than **normal population**" (2016, 2)

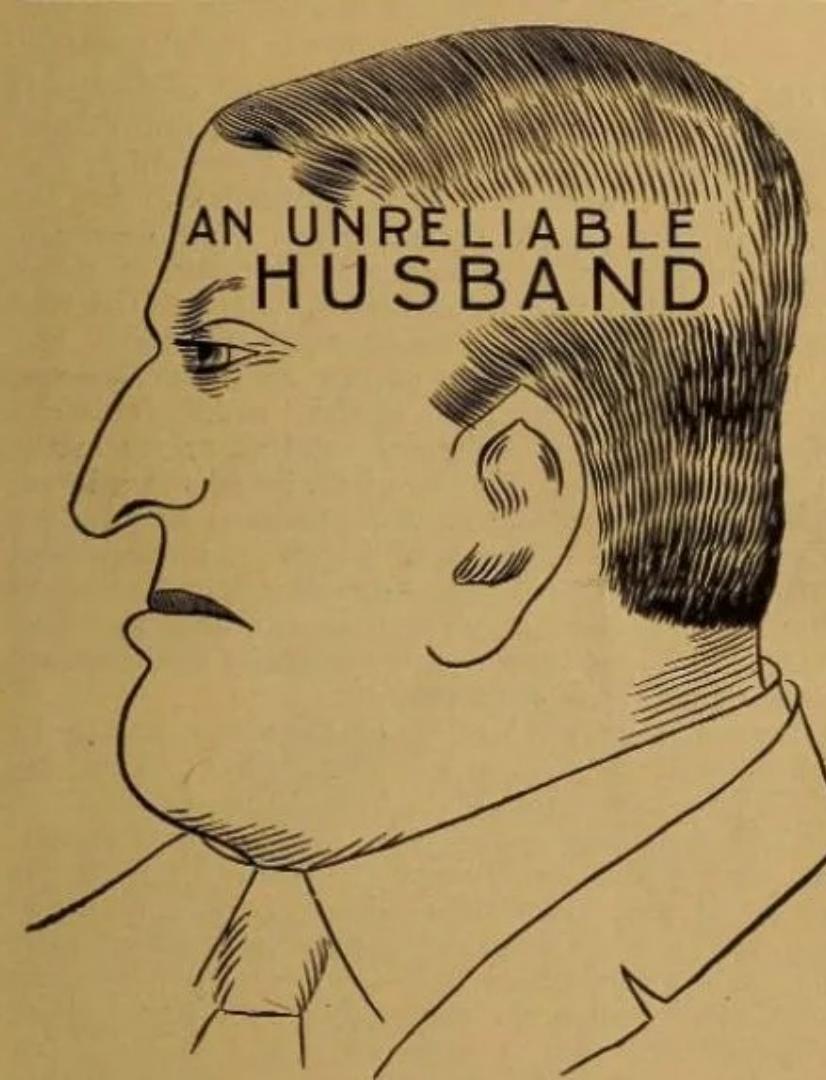
# Wu and Zhang

"The two manifolds formed respectively by the two data sets of criminal and non-criminal faces are concentric, with the manifold for the non-criminal face images lying in the kernel with a smaller span. This newly discovered knowledge suggests **a law of normality for faces of non-criminals: Given the race, gender and age, the faces of general law-abiding public have a greater degree of resemblance compared with the faces of criminals.** In other words, criminals have a significantly higher degree of dissimilarity in facial appearance than **normal population**" (2016, 2)

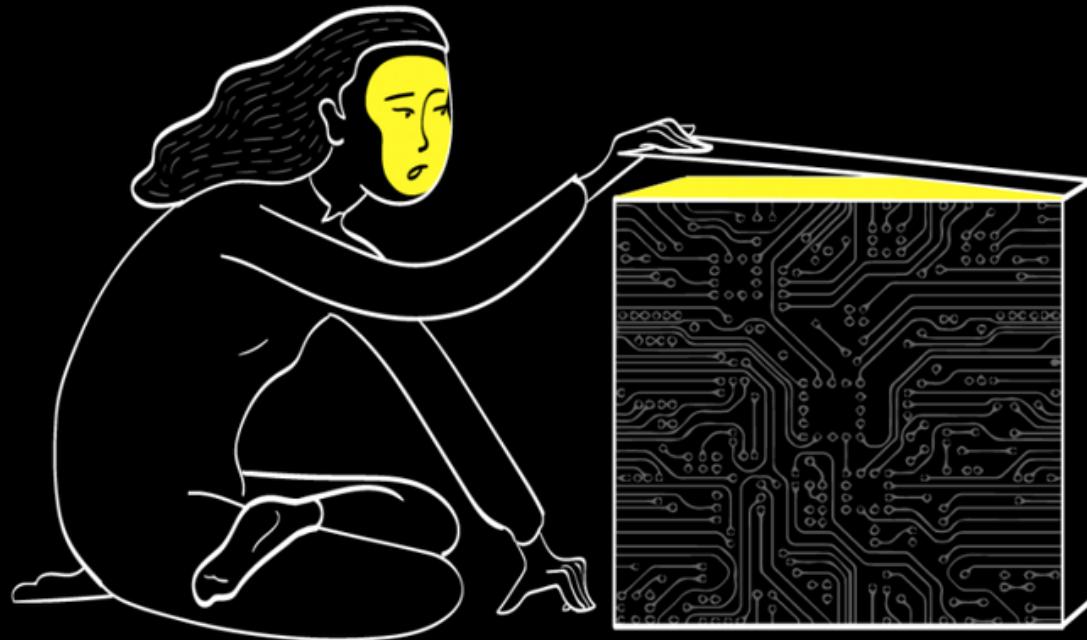
"With a Ph.D in computer science, we know all too well "garbage in and garbage out". However, some of our critics seemed to suggest that machine learning tools cannot be used in social computing simply because no one can prevent **the garbage of human biases from creeping in**. We do not share their pessimism. **Like most technologies, machine learning is neutral.**" (2017, 2)



A GENUINE  
HUSBAND



AN UNRELIABLE  
HUSBAND



# Two epistemologies for ML

# Two epistemologies for ML

- ➊ Transparency/Interpretability
- ➋ Computational Reliabilism

# Two epistemologies for ML

- ➊ Transparency/Interpretability

# Two epistemologies for ML

## 🔑 Transparency/Interpretability

- Requires to “open” the algorithm
- Justification is secured by a “third-party” algorithm (interpretable predictors, XAI)

# Two epistemologies for ML



## Transparency/Interpretability

Synthese (2021) 186:931–942  
https://doi.org/10.1007/s11229-020-02629-9

### The explanation game: a formal framework for interpretable machine learning

David S. Watson<sup>1</sup> · Luciano Floridi<sup>1,2</sup>

Received: 23 October 2019 / Accepted: 12 March 2020 / Published online: 3 April 2020  
© The Author(s) 2020

**Abstract**  
We propose a formal framework for interpretable machine learning. Combining elements from statistical learning, causal interventionism, and decision theory, we design an idealised *explanation game* in which players collaborate to find the best explanations(s) for a given algorithmic prediction. Through an iterative procedure of questions and answers, the players establish a three-dimensional Pareto frontier that describes the optimal trade-offs between explanatory accuracy, simplicity, and relevance. Multiple rounds are played at different levels of abstraction, allowing the players to explore overlapping causal patterns of variable granularity and scope. We characterise the conditions under which such a game is almost surely guaranteed to converge on a (conditionally) optimal explanation surface in polynomial time, and highlight obstacles that will tend to prevent the players from advancing beyond certain explanatory thresholds. The game serves as a descriptive and a normative function, establishing a conceptual space in which to analyse and compare existing proposals, as well as design new and improved solutions.

**Keywords** Algorithmic explainability · Explanation game · Interpretable machine learning · Pareto frontier · Relevance

### 1 Introduction

Machine learning (ML) algorithms have made enormous progress on a wide range of tasks in just the last few years. Some notable recent examples include mastering perfect-information games like chess and Go (Silver et al. 2018), diagnosing skin cancer (Esteva et al. 2017), and proposing new organic molecules (Seeger et al. 2018). These technical achievements have coincided with the increasing ubiquity of ML, which

### Transparency in Complex Computational Systems

Kathleen A. Creel<sup>\*†</sup>

Scientists depend on complex computational systems that are often ineliminably opaque, to the detriment of our ability to give scientific explanations and detect artifacts. Some philosophers have suggested treating opaque systems instrumentally, but computer scientists developing strategies for increasing transparency are correct in finding this unsatisfying. Instead, I propose an analysis of transparency as having three forms: transparency of the algorithm, the representation of the algorithm in code, and the way that code is run on particular hardware and data. This targets the transparency needed for a task, avoiding instrumentalism by providing partial transparency when full transparency is impossible.

**1. Introduction.** Scientists depend on complex computational systems to process their big data, but these systems are not always transparent. Physicists within the Large Hadron Collider's (LHC) Compact Muon Solenoid working group are considering using deep learning algorithms to sort particle collision events and discard the uninteresting ones (Duarte et al. 2018). The new algorithms for doing so, while faster than the old, are complex enough that their decisions cannot be reconstructed in terms of why some events were interesting and thus saved and why others were discarded

Received November 2018; revised October 2019.

\*To contact the author, please write to: University of Pittsburgh, Department of History and Philosophy of Science, 1101 Cathedral of Learning, 4200 Fifth Avenue, Pittsburgh, PA 15260; e-mail: kac28@pitt.edu

†I am grateful for helpful comments from and discussions with Holly Andersen, Robert Battiston, Paul Bunge, Brian Bright, Michael Cesa-Bianchi, Rocco Chiarini, Javier Leonelli, Jake Levine, Edoardo Madsen, Sanda Mitchell, Elmer Nichols, Kathleen Nichols, Anton Novikov, Olivia Ordóñez, William Penn, Rebecca Traister, Porter Williams, Eric Winsberg, and two anonymous reviewers. Thanks also to generous audiences at Philosophical Perspectives on Data-Intensive Science in Hanover; Models and Simulations 8 in Columbia, SC; the Machine Learning Workshop in Irvine, CA; and Science and Art of Science in Paris.

Philosophy of Science, 87 (October 2020) pp. 585–599, 0031-824X(2020)07074-00025.10.00

Copyright 2020 by the Philosophy of Science Association. All rights reserved.

568 i-arg.org/10.1086/709729 Published online by Cambridge University Press

<sup>\*</sup> David S. Watson  
david.watson@oii.ox.ac.uk

<sup>†</sup> Oxford Internet Institute, University of Oxford, 41 Saint Giles, Oxford OX1 3LW, UK

<sup>‡</sup> The Alan Turing Institute, British Library, 96 Euston Road, Kings Cross, London NW1 2DB, UK

Minds & Machines  
https://doi.org/10.1007/s11023-019-09502-w

ORIGINAL ARTICLE



### The Pragmatic Turn in Explainable Artificial Intelligence (XAI)

Andrés Páez<sup>1</sup>

Received: 11 March 2019 / Accepted: 27 May 2019  
© Springer Nature B.V. 2019

### Abstract

In this paper I argue that the search for explainable models and interpretable decisions in AI must be reformulated in terms of the broader project of offering a pragmatic and naturalistic account of understanding in AI. Intuitively, the purpose of providing an explanation of a model or a decision is to make it understandable to its stakeholders. But without a previous grasp of what it means to say that an agent understands a model or a decision, the explanatory strategies will lack a well-defined goal. Aside from providing a clearer objective for XAI focusing on understanding also allows us to relax the factivity condition on explanation, which is impossible to fulfill in many machine learning models, and to focus instead on the pragmatic conditions that determine the best fit between a model and the methods and devices deployed to understand it. After an examination of the different types of understanding discussed in the philosophical and psychological literature, I conclude that interpretive or approximation models not only provide the best way to achieve the objectual understanding of a machine learning model, but are also a necessary condition to achieve post hoc interpretability. This conclusion is partly based on the shortcomings of the purely functionalist approach to post hoc interpretability that seems to be predominant in most recent literature.

**Keywords** Explainable artificial intelligence · Understanding · Explanation · Model transparency · Post-hoc interpretability · Machine learning · Black box models

### 1 Introduction

The main goal of Explainable Artificial Intelligence (XAI) has been variously described as a search for explainability, transparency and interpretability, for ways of validating the decision process of an opaque AI system and generating trust in the

<sup>1</sup> Andrés Páez  
andrespaez@gmail.com

<sup>2</sup> Department of Philosophy, Universidad de los Andes, Carrera 1 No. 18A-12 (G-533), Bogotá, DC 111171, Colombia

AI & SOCIETY (2021) 36:585–595  
https://doi.org/10.1007/s00146-020-00166-z

OPEN FORUM

### Artificial intelligence and the value of transparency

Joel Walmsley

Received: 6 January 2020 / Accepted: 25 August 2020 / Published online: 8 September 2020  
© Springer Nature London Ltd. part of Springer Nature 2020

### Abstract

Some recent developments in Artificial Intelligence—especially the use of machine learning systems, trained on big data sets and deployed in socially significant and ethically weighty contexts—have led to a number of calls for “transparency”. This paper explores the epistemological significance of this concept, as well as surveying and discussing the range of ways in which it has been invoked in recent discussions. Within “functional” forms of transparency (concerning the relationships between an AI system, its developers, users and the public) we straightforwardly achieved what I call “functional” transparency about the inner workings of a system is, in many cases, much harder to attain. In those situations, I argue that contestability may be a possible, acceptable, and alternative so that even if we cannot understand how a system came up with a particular output, we at least have the means to challenge it.

**Keywords** Transparency · Explainability · Contestability · Machine learning · Bias

### 1 Introduction

Alongside, and arguably because of, some of the most recent technological developments in Artificial Intelligence, the last few years have seen a growing number of calls for various forms of transparency within and about the field. For example, the European Commission’s High-Level Expert Group on Artificial Intelligence (HLEG) for Trustworthy AI—features the notion of transparency prominently, and the European Union’s General Data Protection Regulation (GDPR) includes the stipulation that, when a person is subject to automated decisions that significantly affect him/her, he/she has the right to obtain human intervention, to express his or her point of view, to obtain an explanation of the decision reached after such assessment and to challenge the decision.<sup>1</sup> In these calls respond to an epistemic limitation: machine learning techniques, together with their use of “Big Data” for training purposes, are often considered to be black boxes, obscuring for a complete understanding, and faster and more powerful than human cognition (at least, on the relatively narrow set of tasks for which AI is designed). Of course, in many cases,

<sup>1</sup> Sometimes also discussed under the heading of “explainability,” “inherent” or “design” stances further from reference, also, to “accountability,” “intelligibility” and “interpretability.” See Walmsley (2019).

<sup>2</sup> General Data Protection Regulation, Recital 71, available at <https://gdpr-info.eu/recital/no-71/>.

<sup>3</sup> See Donell (1991).

<sup>4</sup> Springer

# Two epistemologies for ML

- ➊ Transparency/Interpretability
- ➋ Computational Reliabilism

# Two epistemologies for ML

🔑 Transparency/Interpretability

⌚ Computational Reliabilism

- Accepts “black box” algorithms
- Justification comes from securing the reliability of the algorithm
  - Depends on self-regulating practices, methods, and processes with reliability-conferring property — external to the algorithm.

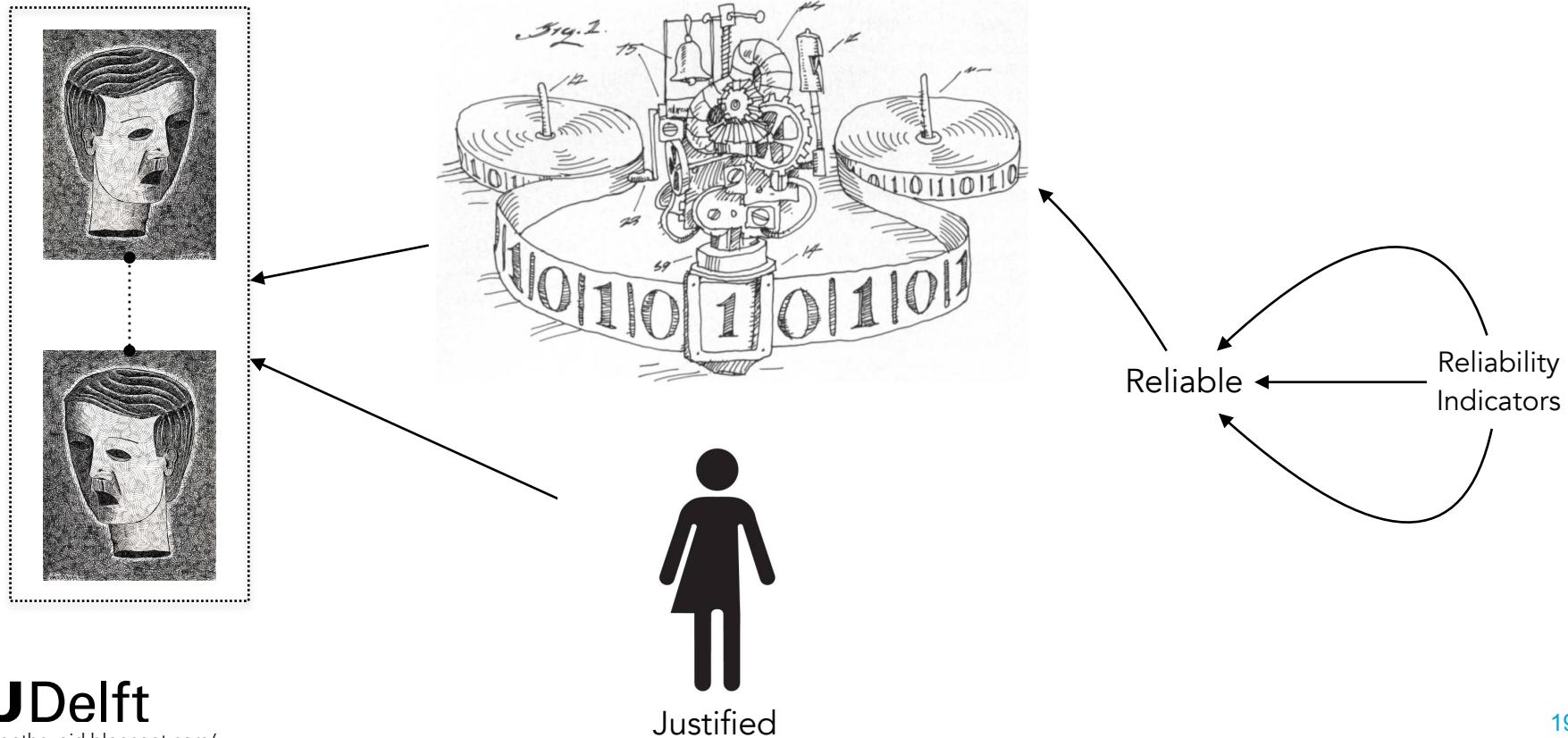
# What is process reliabilism? ( $\approx$ Goldman 1979)

**Process reliabilism:** An individual's belief is justified in case it is produced by a reliable belief-forming process (or sequence of processes)

A reliable belief-forming process has the tendency to produce beliefs that are true rather than false.



# Justification via CR (a sketch)



# Two epistemologies for ML

🔑 Transparency/Interpretability

⌚ Computational Reliabilism (type/token)

- RI<sub>1</sub> Technical robustness of algorithms
- RI<sub>2</sub> Computer-based scientific practice
- RI<sub>3</sub> Social construction of reliability

# (Type) Reliability indicators for ML

- RI<sub>1</sub> Technical robustness of algorithms
- RI<sub>2</sub> Computer-based scientific practice
- RI<sub>3</sub> Social construction of reliability

# Reliability indicators for ML

- RI<sub>1</sub> Technical robustness of algorithms

# (Token) Reliability indicators for ML

## RI<sub>1</sub> Technical robustness of algorithms

- Verification and validation methods ( $\rightarrow$  accuracy)
  - Robustness analysis
  - Data choices
  - Parametrization and hyper-parametrization choices
  - A history of (un)successful implementations
  - Error treatment (e.g., uncertainty quantification)

(Durán 2018; Durán & Formanek 2018)

Minds and Machines (2018) 28:645–666  
<https://doi.org/10.1007/s11023-018-9481-6>



# Grounds for Trust: Essential Epistemic Opacity and Computational Reliabilism

Juan M. Durán<sup>1</sup>  · Nico Formanek<sup>2</sup>

Received: 22 May 2018 / Accepted: 12 October 2018 / Published online: 29 October 2018  
© The Author(s) 2018

## Abstract

Several philosophical issues in connection with the assumption that results of simulations are reliable are discussed. The article starts with the debate on the experimental role of computer simulations (Barberousse and Vormund, 2009; Morrison et al., 2009). It then discusses the reliability of computer data (Barberousse and Vormund, 2009) and the changing face of scientific publishing (Humphreys, 2013; Humphreys and Vormund, 2013), and the changing face of scientific publishing (Barberousse and Vormund, 2013), and the explanatory power of computer simulations (Vormund, 2017). The aim of this article is to show that the reliability of computer simulations is based on reliable processes. After a short reconstruction of the debate on the reliability of computer simulations, the article elaborates extensively on process reliabilism with computer simulations, followed by a discussion of four sources for process reliabilism: theory, validation, robustness analysis and successful implementations, and the role of

**Keywords** Computer simulations · Re-validation · Robustness analysis · Hist

# Reliability indicators for ML

- RI<sub>1</sub> Technical robustness of algorithms
- RI<sub>2</sub> Computer-based scientific practice

# Reliability indicators for ML

● RI<sub>1</sub> Technical robustness of algorithms

● RI<sub>2</sub> Computer-based scientific practice

- Acting mechanisms (causality and otherwise)
- Model-Trainning practices / Epistemic and moral values
- Representation
- Some degree of theoretical, conceptual coherence (against JBDA)
  - Background knowledge (e.g., knowledge of the learned weights)
  - Some morphism between the implemented methodology and an established scientific methodology (e.g., use of natural kinds, accepted cut-off-values — e.g., when melanoma is carcinogenic)
  - Against JBDA
    - JBDA cannot contribute to the confirmation of theories over time
    - JBDA cannot enable inferences to similar phenomena

# (Type) Reliability indicators for ML

- RI<sub>1</sub> Technical robustness of algorithms
- RI<sub>2</sub> Computer-based scientific practice
- RI<sub>3</sub> Social construction of reliability

# (Type) Reliability indicators for ML

⌚ RI<sub>1</sub> Technical robustness of algorithms

⌚ RI<sub>2</sub> Computer-based scientific practice

⌚ RI<sub>3</sub> Social construction of reliability

- Peer and social acceptability of ML output
- Favorable, successful, scientifically valid use of ML/ML output (e.g., develop of new hypothesis or expand into new strategies)
- Coherence with established body of scientific beliefs
- Coherence with scientists' theoretical and other commitments
- Realization of the epistemic and ethical values

# Wu & Zhang: Automated Inference on Criminality Using Face Images



(a) Three samples in criminal ID photo set  $S_c$ .



(b) Three samples in non-criminal ID photo set  $S_n$

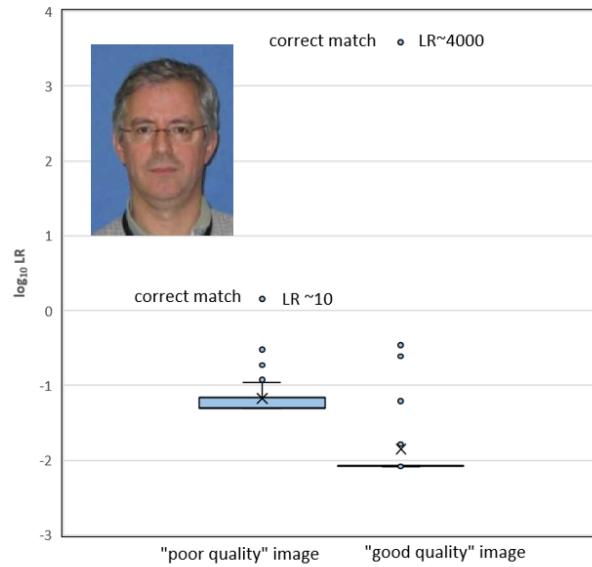
Figure 1. Sample ID photos in our data set.

- The system is forced to pick one out of two classes: “criminal” or “not-criminal”
- It fossilises concepts, such as “criminality”
- It fails to account for or include values (e.g., “innocent until proved guilty”)
- It is disconnected from a larger body of knowledge on:
  - Social construction of crime
  - The socio-economic basis of criminality
  - The psychology of criminals

# How does CR work? A case in Forensic AI

# Forensic face comparison

"Poor quality" image



"Good quality" image





# Nederlands Forensisch Instituut

## Ministerie van Justitie en Veiligheid

Tabel 2: Overzicht van alle gezichtbeeldvergelijkingen. Groen: geselecteerd voor gedetailleerde vergelijking. Rood: op basis van voorselectie uitgesloten.

	Zaak_01	Zaak_04	Zaak_05	Zaak_07	Zaak_13	Zaak_14	Zaak_15
	(2013)	(2012)	(2013)	(2013)	(2013)	(2013)	(2013)
Sofia  (2013, 2016)							
Diana  (2012, 2018)							
Djanija  (2012)							
Jasminka  (2016)							
Vahida  (datum onbekend, PV 2015)							
Vahida  (2001)							
Vera  (2013)							

Tabel 3: Samenvatting van de conclusies van de vergelijkingen (met cijfers volgens bijgaande legenda), waarbij a, b en c staat voor de conclusies van de afzonderlijke onderzoekers en 'Consensus' na besprekking van de resultaten (legenda: zie Tabel 4)

Zaak	Verdachte	a	b	c	Consensus	Zaak	Verdachte	a	b	c	Consensus
01	Sofia	2	2	2	2	13	Sofia	3	3	2	3
01	Diana	-3	-3	-3	-3	13	Diana	-3	-3	-3	-3
01	Djanija	-4	-3	-4	-4	13	Djanija	-3	-2	-4	-3
01	Vahida (jonger)	-2	-2	-4	-3	13	Vahida (ouder)	-3	-2	-3	-3
01	Vahida (ouder)	-3	1	-4	-2	14	Sofia	1	2	2	2
01	Vera	-4	-2	-4	-4	14	Diana	-2	-2	-3	-3
04	Sofia	4	3	3	4	14	Djanija	-4	-2	-4	-3
04	Diana	-4	-5	-4	-4	14	Jasminka	-1	-3	-4	-2
04	Djanija	-4	-3	-4	-4	14	Vahida (ouder)	-1	-2	-4	-2
05	Sofia	5	3	3	4	14	Vera	-3	-2	-4	-3
05	Diana	-4	-5	-4	-4	15	Sofia	3	4	3	3
05	Djanija	-4	-3	-5	-4	15	Diana	-3	-5	-4	-4
05	Vahida (ouder)	-4	-4	-4	-4	15	Djanija	-4	-5	-4	-4
07	Sofia	0	1	2	1	15	Vera	-4	-3	-4	-4
07	Diana	-2									
07	Djanija	-3									
07	Jasminka	-2									
07	Vera	-4									

Tabel 4: Legenda bij Tabel 3

### Legenda

Cijfer	De bevindingen van het onderzoek zijn:	Cijfer	De bevindingen van het onderzoek zijn:
5	extrem veel waarschijnlijker wanneer H1 waar is dan wanneer H2 waar is	-1	iets waarschijnlijker wanneer H2 waar is dan wanneer H1 waar is
4	zeer veel waarschijnlijker wanneer H1 waar is dan wanneer H2 waar is	-2	waarschijnlijker wanneer H2 waar is dan wanneer H1 waar is
3	veel waarschijnlijker wanneer H1 waar is dan wanneer H2 waar is	-3	veel waarschijnlijker wanneer H2 waar is dan wanneer H1 waar is
2	waarschijnlijker wanneer H1 waar is dan wanneer H2 waar is	-4	zeer veel waarschijnlijker wanneer H2 waar is dan wanneer H1 waar is
1	iets waarschijnlijker wanneer H1 waar is dan wanneer H2 waar is	-5	extrem veel waarschijnlijker wanneer H2 waar is dan wanneer H1 waar is
0	ongeveer even waarschijnlijk wanneer H1 waar is als wanneer H2 waar is		



# Nederlands Forensisch Instituut

## Ministerie van Justitie en Veiligheid

Tabel 2: Overzicht van alle gezichtbeeldvergelijkingen. Groen: geselecteerd voor gedetailleerde vergelijking. Rood: op basis van voorselectie uitgesloten.

	Zaak_01	Zaak_04	Zaak_05	Zaak_07	Zaak_13	Zaak_14	Zaak_15
	(2013)	(2012)	(2013)	(2013)	(2013)	(2013)	(2013)
Sofia  (2013, 2016)							
Diana  (2012, 2018)							
Djanija  (2012)							
Jasminka  (2016)							
Vahida  (datum onbekend, PV 2015)							
Vahida  (2001)							
Vera  (2013)							

Tabel 5: Vergelijking van conclusies vorige en huidige rapportage

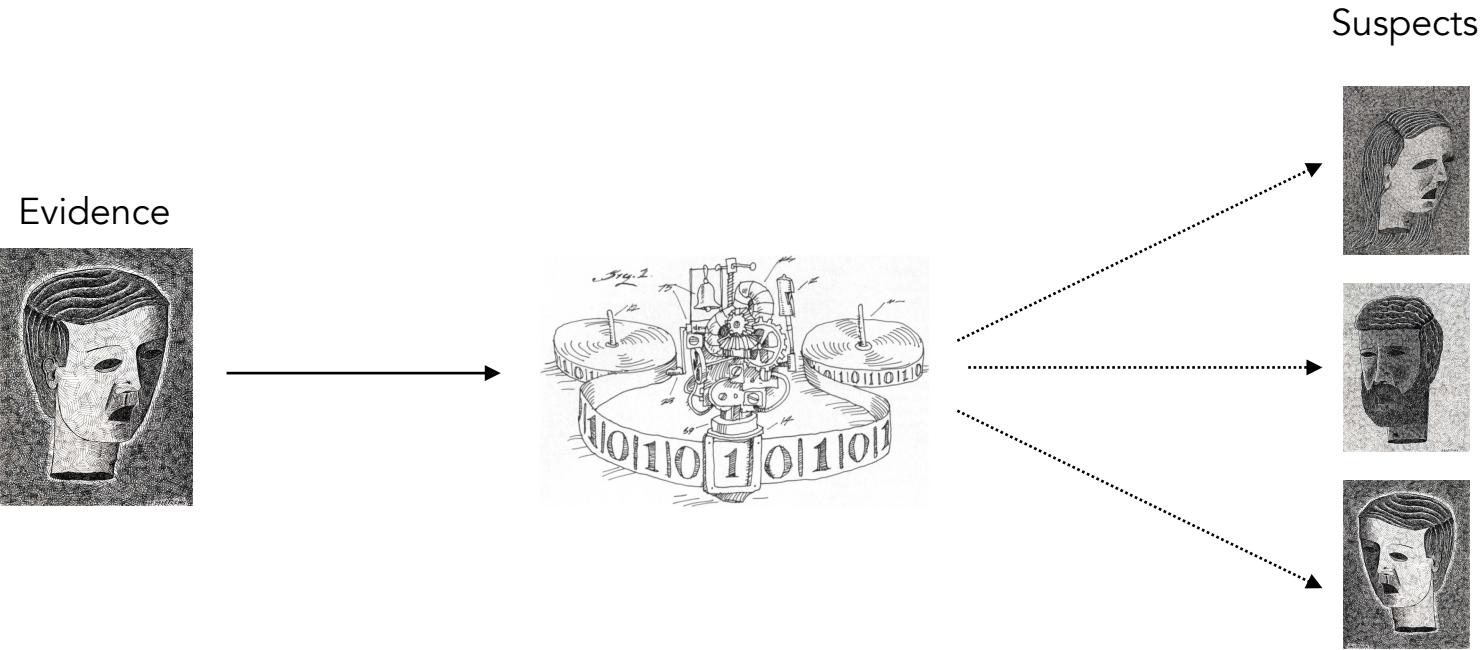
Zaaknr.	Rapportage aanvraag 002 (25 maart 2016)	Huidige aanvraag 003
Zaak 01	Waarschijnlijker	Waarschijnlijker
Zaak 04	Zeer veel waarschijnlijker	Zeer veel waarschijnlijker
Zaak 05	Zeer veel waarschijnlijker	Zeer veel waarschijnlijker
Zaak 07	--	Iets waarschijnlijker
Zaak 13	Veel waarschijnlijker	Veel waarschijnlijker
Zaak 14	Iets waarschijnlijker	Waarschijnlijker
Zaak 15	Veel waarschijnlijker	Veel waarschijnlijker

Tabel 3: Samenvatting van de conclusies van de vergelijkingen (met cijfers volgens bijgaande legenda), waarbij a, b en c staat voor de conclusies van de afzonderlijke onderzoekers en 'Consensus' na besprekking van de resultaten (legenda: zie Tabel 4)

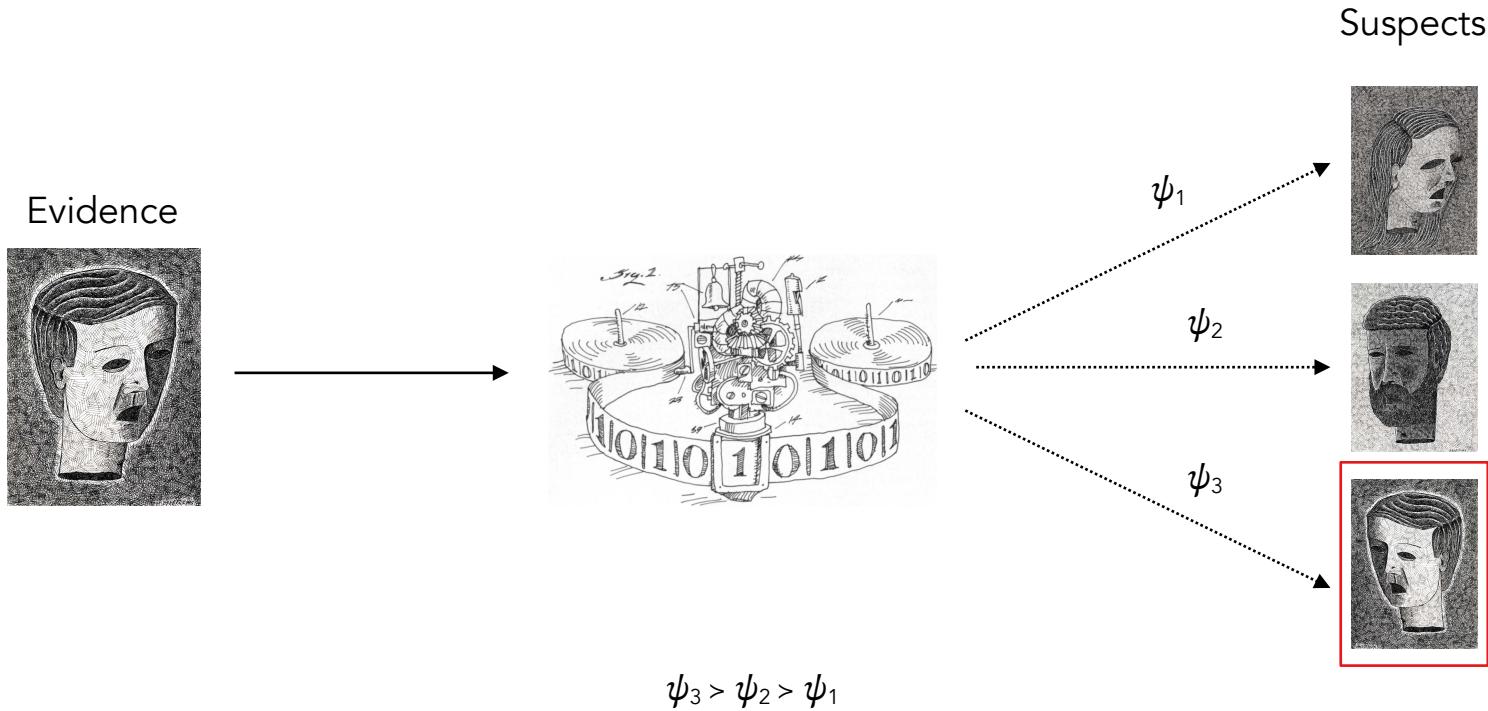
Zaak	Verdachte	a	b	c	Consensus
01	Sofia	2	2	2	2
01	Diana	-3	-3	-3	-3
01	Djanija	-4	-3	-4	-4
01	Vahida (jonger)	-2	-2	-4	-3
01	Vahida (ouder)	-3	1	-4	-2
01	Vera	-4	-2	-4	-4
04	Sofia	4	3	3	4
04	Diana	-4	-5	-4	-4
04	Djanija	-4	-3	-4	-4
05	Sofia	5	3	3	4
05	Diana	-4	-5	-4	-4
05	Djanija	-4	-3	-5	-4
05	Vahida (ouder)	-4	-4	-4	-4
07	Sofia	0	1	2	1
07	Diana	2	2	2	2

onderzoek zijn:	Cijfer	De bevindingen van het onderzoek zijn:
H2 waar is	-1	iets waarschijnlijker wanneer H2 waar is dan wanneer H1 waar is
H2 waar is	-2	waarschijnlijker wanneer H2 waar is dan wanneer H1 waar is
H2 waar is	-3	veel waarschijnlijker wanneer H2 waar is dan wanneer H1 waar is
H2 waar is	-4	zeer veel waarschijnlijker wanneer H2 waar is dan wanneer H1 waar is
H2 waar is	-5	extremee veel waarschijnlijker wanneer H2 waar is dan wanneer H1 waar is
H2 waar is		

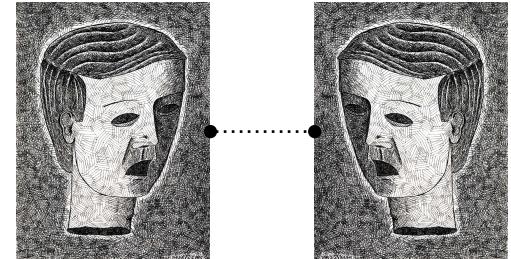
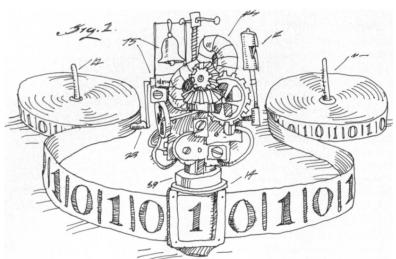
# Forensic AI



# Forensic AI



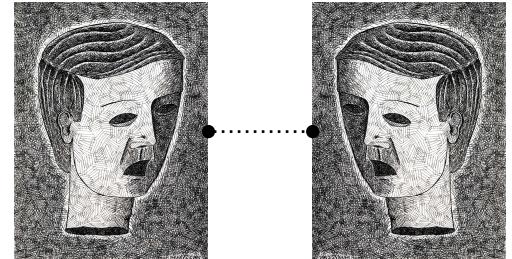
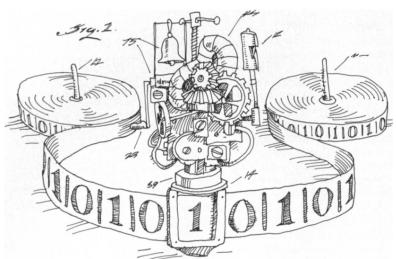
# Forensic AI



## RI<sub>1</sub> Technical robustness of algorithms

- Verification and Validation
- Robustness analysis
- Redundancy mechanisms
- Quality of data

# Forensic AI

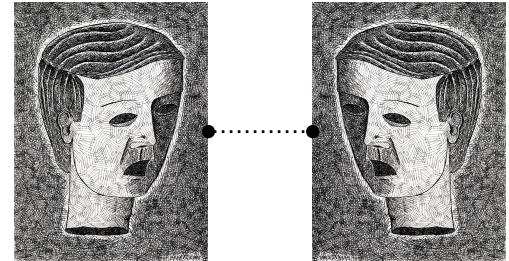
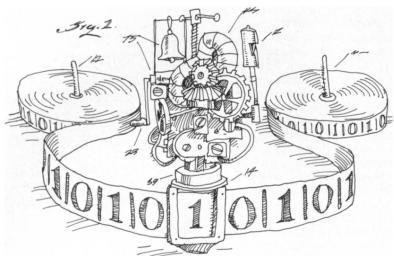


RI<sub>1</sub> Technical robustness of algorithms

RI<sub>2</sub> Computer-based scientific practice

- Verification and Validation
  - Robustness analysis
  - Redundancy mechanisms
  - Quality of data
- 
- Implementation of quality assurance methods
  - “We implement well-known and tested libraries”
  - Biometrical techniques

# Forensic AI



RI<sub>1</sub> Technical robustness of algorithms

- Verification and Validation
  - Robustness analysis
  - Redundancy mechanisms
  - Quality of data

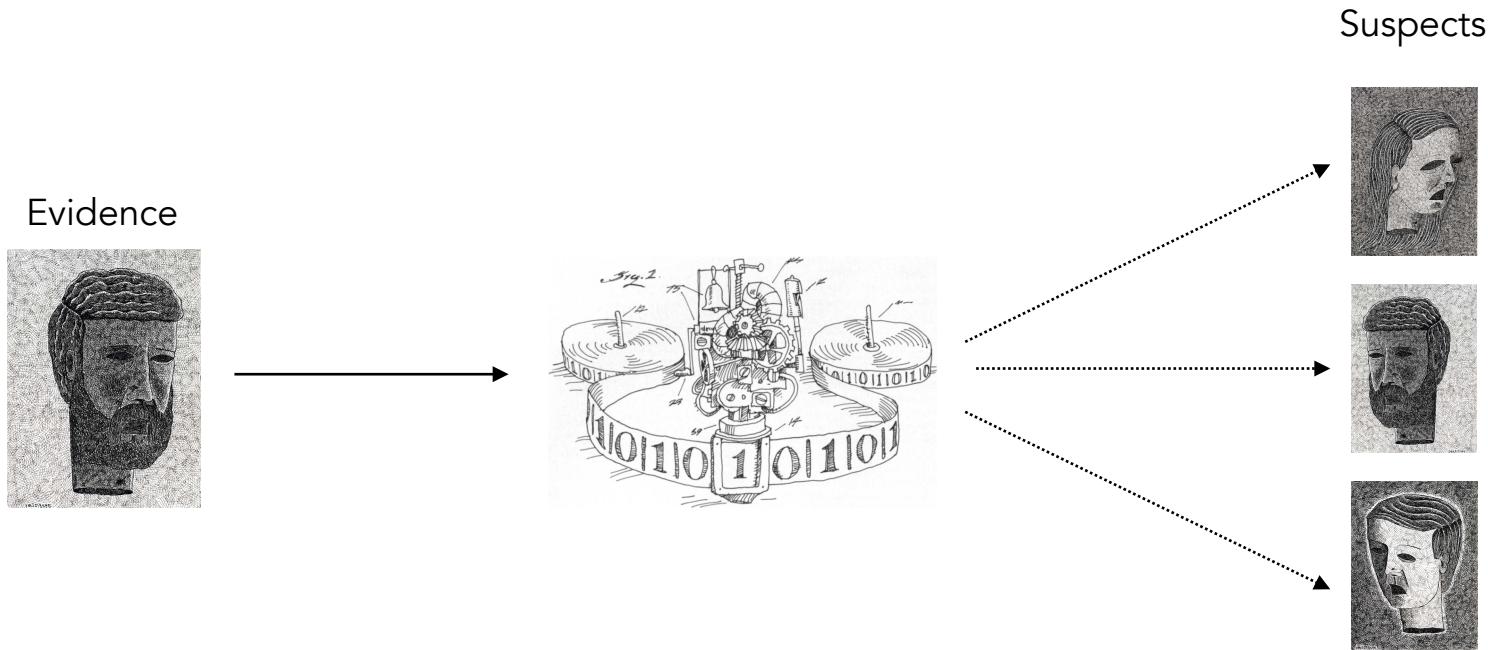
RI<sub>2</sub> Computer-based scientific practice

- Implementation of quality assurance methods
  - “We implement well-known and tested libraries”
  - Biometrical techniques

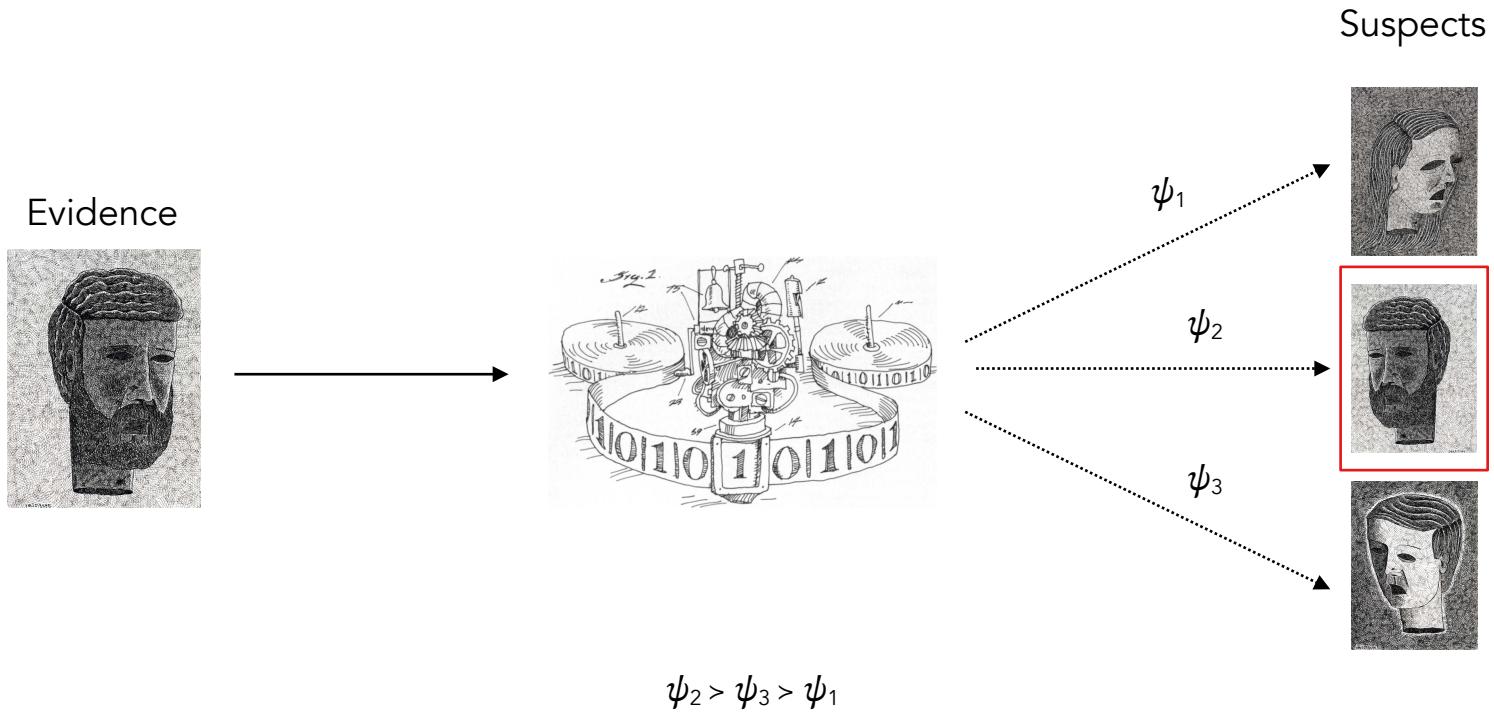
RI<sub>3</sub> Social construction of reliability

- Debates about the output  
(e.g., soundness, applicability, etc.)
  - Contrast with other non-computational methods  
(e.g., use of code review,  
expert assessment and report)

# Forensic AI



# Forensic AI



# Lights and shadows of CR

# Lights

- Externalist:
  - Credibility does not depend on third-party algorithms and are independent on the quality of our insight and our mental capabilities to understand it
- Each reliability indicator is independently justified
- There are various reliability indicators ("non-centralized")
- It detects markers for non-credible results
- It accommodates different needs: in some systems, V&V are more valuable, whereas in others is theoretical coherence

# Shadows

- It can't be automated (or the operationalization is very difficult)
  - It might hamper algorithmic utility
- It is unclear the precedence, order, and weight of each indicator
- *The tyranny of the few:* a few indicators might suffice for crediting reliability/non-reliability to an algorithm (when missing indicators or when conflict of indicators)



*View of Delft*, Johannes Vermeer, 1660-1661. Rijksmuseum, Amsterdam

Thank you!

j.m.duran@tudelft.nl

Draft paper: juanmduran.net & PhilPapers