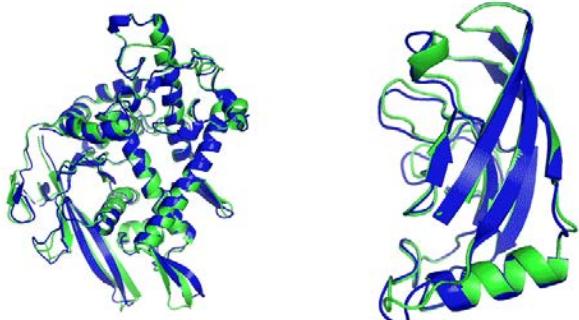


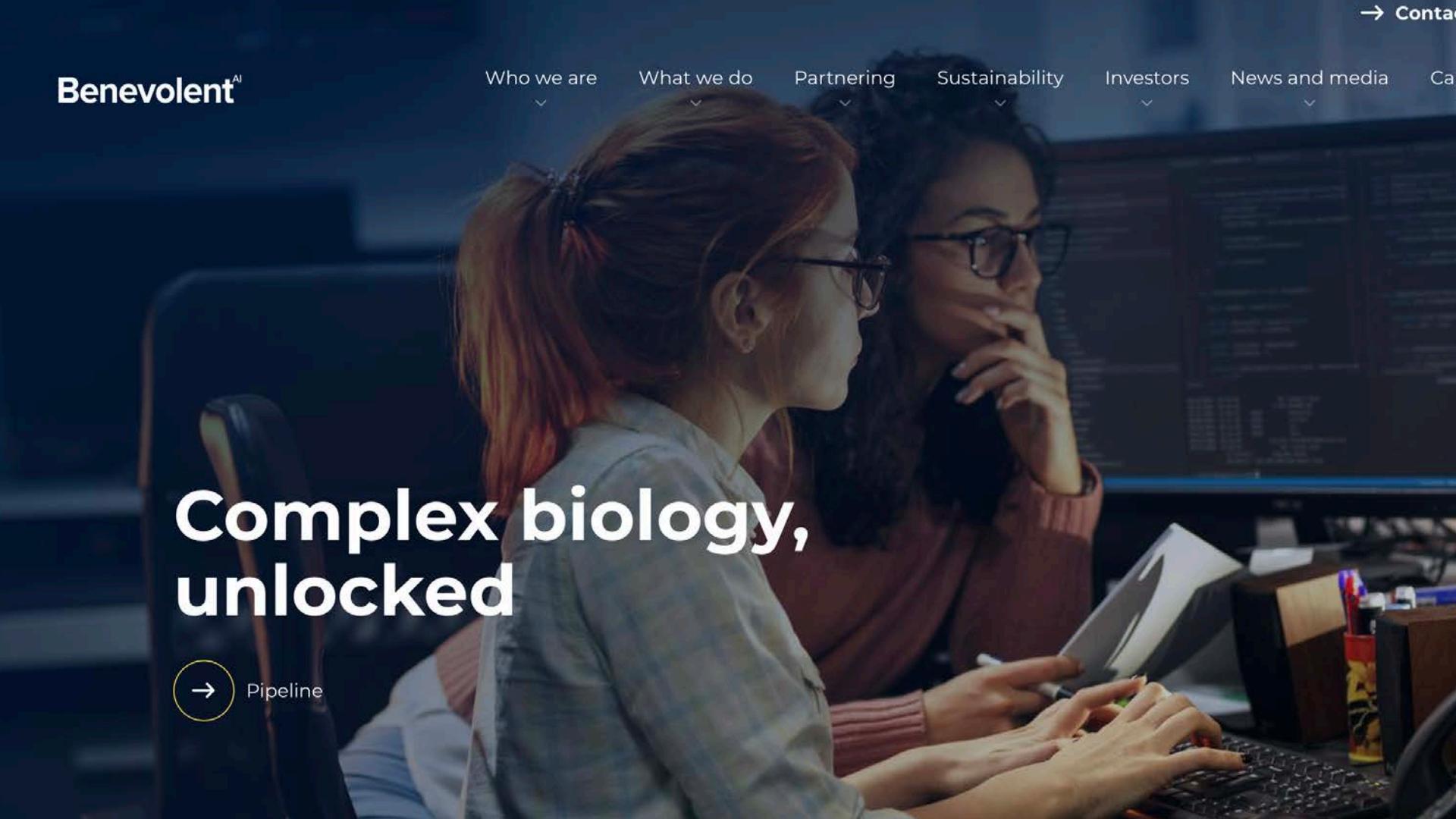
# Hacia una epistemología de algoritmos: fiabilidad computacional y sus límites



T1037 / 6vr4  
90.7 GDT  
(RNA polymerase domain)

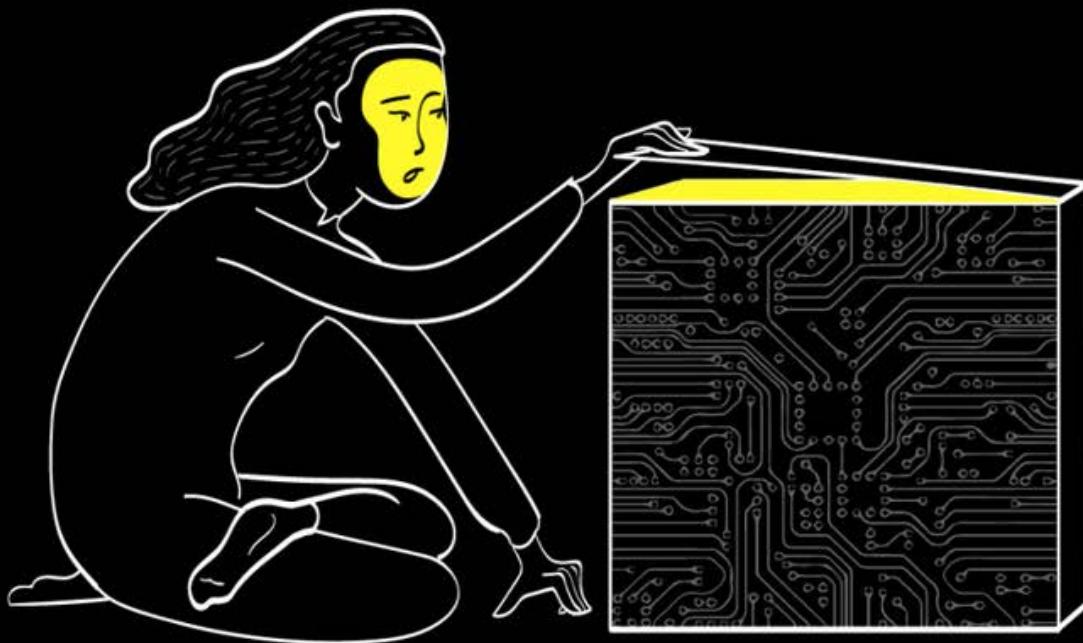
T1049 / 6y4f  
93.3 GDT  
(adhesin tip)

- Experimental result
- Computational prediction

A photograph of two women in an office setting. One woman, with red hair tied back, wears glasses and a light blue plaid shirt, looking down at a laptop keyboard. The other woman, with dark curly hair and glasses, wears a maroon sweater and holds a white document, looking at it thoughtfully with her hand near her chin. They are positioned in front of a computer monitor displaying code or data. The background is slightly blurred.

# Complex biology, unlocked

[Pipeline](#)



# Opacidad epistémica

[A] process is epistemically opaque relative to a cognitive agent  $S$  at time  $t$  just in case  $S$  does not know at  $t$  all of the epistemically relevant elements of the process (Humphreys 2009, 618)

# Opacidad epistémica

[A] process is epistemically opaque relative to a cognitive agent  $S$  at time  $t$  just in case  $S$  does not know at  $t$  all of the epistemically relevant elements of the process (Humphreys 2009, 618)

"For a mathematical proof, one agent may consider a particular step in the proof to be an epistemically relevant part of the justification of the theorem, whereas to another, the step is sufficiently trivial to be eliminable. [...] Within the hybrid scenario [i.e., algorithms], no human can examine and justify every element of the computational processes that produce the output of a computer simulation or other artifacts of computational science [...] Many, perhaps all, of the features that are special to simulations are a result of this inability of human cognitive abilities to know and understand the details of the computational process" (2009, 618)

# Opacidad epistémica

[A] process is epistemically opaque relative to a cognitive agent  $S$  at time  $t$  just in case  $S$  does not know at  $t$  all of the epistemically relevant **elements of the process** (Humphreys 2009, 618)

"For a mathematical proof, one agent may consider **a particular step in the proof** to be an epistemically relevant part of the justification of the theorem, whereas to another, the step is sufficiently trivial to be eliminable. [...] Within the hybrid scenario [i.e., algorithms], no human can examine and justify **every element** of the computational processes that produce the output of a computer simulation or other artifacts of computational science [...] Many, perhaps all, of the features that are special to simulations are a result of this inability of human cognitive abilities to know and understand **the details** of the computational process" (2009, 618)

# Opacidad epistémica

[A] process is epistemically opaque relative to a cognitive agent  $S$  at time  $t$  just in case  $S$  **does not know** at  $t$  all of the epistemically relevant elements of the process (Humphreys 2009, 618)

"For a mathematical proof, one agent may consider a particular step in the proof to be an epistemically relevant part of the justification of the theorem, whereas to another, the step is sufficiently trivial to be eliminable. [...] Within the hybrid scenario [i.e., algorithms], no human can examine and justify every element of the computational processes that produce the output of a computer simulation or other artifacts of computational science [...] Many, perhaps all, of the features that are special to simulations are a result of this inability of human cognitive abilities to know and understand the details of the computational process" (2009, 618)

# Opacidad epistémica

[A] process is epistemically opaque relative to a cognitive agent  $S$  at time  $t$  just in case  $S$  does not know at  $t$  all of **the epistemically relevant elements of the process** (Humphreys 2009, 618)

"For a mathematical proof, one agent may consider a particular step in the proof to be an epistemically relevant part of the justification of the theorem, whereas to another, the step is sufficiently trivial to be eliminable. [...] Within the hybrid scenario [i.e., algorithms], no human can examine and justify every element of the computational processes that produce the output of a computer simulation or other artifacts of computational science [...] Many, perhaps all, of the features that are special to simulations are a result of this inability of human cognitive abilities to know and understand the details of the computational process" (2009, 618)

# Opacidad epistémica

- Naturaleza del proceso y elementos del proceso: instanciación de variables, llamadas a funciones, sentencias condicionales, operaciones aritméticas y lógicas, manejo de errores, estructuras de datos, pero también prácticas, métricas; “one may have excellent reasons for holding that a particular parametric family of models is applicable to the case at hand, yet have only empirical methods available for deciding which parametric values are the right ones” (2004, 150)
- Falta de capacidad de inspección de S: los algoritmos y los procesos computacionales son cajas negras.
- Pertinencia de los elementos epistémicos relevantes son con el propósitos de justificar el resultado

# De opacidad a transparencia

"If we think in terms of such a process [i.e., algorithms] and imagine that its stepwise computation was slowed down to the point where, in principle, a human could examine each step in the process, the computationally irreducible process would become epistemically transparent. What this indicates is that the practical constraints we have previously stressed, primarily the need for computational speed, are the root cause of all epistemic opacity in this area. Because those constraints cannot be circumvented by humans, we must abandon the insistence on epistemic transparency for computational science. What replaces it would require an extended work in itself, but the prospects for success are not hopeless." (Humphreys, 2004, 150)

# Transparencia/interpretabilidad

Synthese (2021) 198:921–9242  
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 characterize 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 (Sugler et al. 2018). These technical achievements have coincided with the increasing ubiquity of ML, which

<sup>1</sup> David S. Watson  
david.watson@oii.ox.ac.uk

<sup>1</sup> Oxford Internet Institute, University of Oxford, 41 Saint Giles, Oxford OX1 3LW, UK  
<sup>2</sup> The Alan Turing Institute, British Library, 96 Euston Road, Kings Cross, London NW1 2DB, UK



## Transparency in Complex Computational Systems

Kathleen A. Crecel<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 have argued that this is not always appropriate. I argue that this is because of a misconception. Instead, I propose an analysis of transparency as having three forms: transparency of the algorithm, the realization of the algorithm in the code, and the way that code is run on particular hardware and data. This targets the transparency most useful 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 go to: University of Pittsburgh, Department of History and Philosophy of Science, 110 Cathedral of Learning, 4200 Fifth Avenue, Pittsburgh, PA 15260; e-mail: kac284@pitt.edu

I am grateful for helpful comments from and discussions with Holly Andemen, Robert Batterman, Nora Mills Boyd, Liam Kofi Bright, Marvita Chirimata, Roger Crecel, Javier Duarte, Muriel Hardcastle, Paul Humphreys, Benjamin Jantzen, Johannes Lenhard, Sabina Leonelli, Jake Levinson, Edmund Machinery, Sandra Mitchell, Elmer Nichols, Kathleen O'Neil, Andrew Potts, Michael R. Williams, and Tracy Peck Williams, as well as Eric Winsberg, and two anonymous reviewers. Thanks also to generous audience 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 Simulation IV in Stuttgart.

Philosophy of Science, 87 (October 2020) pp. 568–589 0031-8248/20/040002510.00  
Copyright 2020 by the Philosophy of Science Association. All rights reserved.

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

ORIGINAL ARTICLE

AI & SOCIETY (2021) 36:985–995  
https://doi.org/10.1007/s00146-020-01966-z

OPEN FORUM

## Artificial intelligence and the value of transparency

Joel Walmsley<sup>1</sup>

Received: 6 January 2020 / Accepted: 25 August 2020 / Published online: 8 September 2020  
© Springer Verlag 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 and ethical dimensions of that concept, as well as surveying and lacuna concerning the variety of ways in which it has been invoked in recent discussions. Whilst “outward” forms of transparency (concerning the relationship between an AI system, its developers, users and the media) may be 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 useful 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 interest in the concept of “transparency” within and about the field. For example, the 2019 report from the European Commission's High-Level Expert Group on AI—entitled Ethical Guidelines for Trustworthy AI—includes the notion of transparency as one of the six fundamental principles. The General Data Protection Regulation (GDPR) includes the stipulation that, when a person is subject to an automated decision based on their personal information, he or she has “the right to obtain human intervention, to express his or her point of view, to obtain an explanation of the logic involved in the automated assessment and to challenge the decision.”<sup>1</sup> In part, these calls respond to an epistemic limitation: machine learning techniques, together with the use of “Big Data” for training purposes, mean that many AI systems are both too complex for a complete understanding, and faster and more powerful than ever before (at least, on the relatively narrow set of tasks for which AI is designed). Of course, in many cases,

<sup>1</sup> Some also discussed under the heading of “explainability,” “interpretability,” “accountability” (e.g., by Rohr 2019) or with reference, also, to “accountability,” “intelligibility” and “interpretability” (e.g., in Floridi et al. 2018).

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

<sup>3</sup> See Dosen (1971).

# ¿Cómo obtenemos transparencia?

Table 2. Summary of Methods for Opening Black Boxes Solving the *Model Explanation Problem*

Name	Ref.	Authors	Year	Explainer	Black Box	Data Type	General	Random	Examples	Code	Dataset
Trepan	[22]	Craven et al.	1996	DT	NN	TAB	✓				✓
—	[57]	Krishnan et al.	1999	DT	NN	TAB	✓		✓		
DecText	[12]	Boz	2002	DT	NN	TAB	✓	✓			
GPDT	[46]	Johansson et al.	2009	DT	NN	TAB	✓	✓			
Tree Metrics	[17]	Chapman et al.	1998	DT	TE	TAB					
CCM	[26]	Domingos et al.	1998	DT	TE	TAB	✓	✓			
—	[34]	Gibbons et al.	2013	DT	TE	TAB	✓	✓			
STA	[140]	Zhou et al.	2016	DT	TE	TAB	✓				
CDT	[104]	Schetinin et al.	2007	DT	TE	TAB					
—	[38]	Hara et al.	2016	DT	TE	TAB	✓				
TSP	[117]	Tan et al.	2016	DT	TE	TAB					
Conj Rules	[21]	Craven et al.	1994	DR	NN	TAB	✓				
G-REX	[44]	Johansson et al.	2003	DR	NN	TAB	✓	✓			
REFINE	[141]	Zhou et al.	2003	DR	NN	TAB	✓	✓			
RxREN	[6]	Augusta et al.	2012	DR	NN	TAB	✓				
SVM+P	[82]	Nunez et al.	2002	DR	SVM	TAB					
—	[33]	Fung et al.	2005	DR	SVM	TAB					
inTrees	[25]	Deng	2014	DR	TE	TAB					
—	[70]	Lou et al.	2013	FI	AGN	TAB	✓				
GoldenEye	[40]	Henelius et al.	2014	FI	AGN	TAB	✓	✓			
PALM	[58]	Krishnan et al.	2017	DT	AGN	ANY	✓				
FIRM	[142]	Zien et al.	2009	FI	AGN	TAB	✓	✓			
MFI	[124]	Vidovic et al.	2016	FI	AGN	TAB	✓	✓			
—	[121]	Tolomei et al.	2017	FI	TE	TAB					
POIMs	[111]	Sonnenburg et al.	2007	FI	SVM	TAB					

Guidotti, Monreale, Ruggieri,  
Turini, Giannotti, Pedreschi, (2018)

Table 3. Summary of Methods for Opening Black Boxes Solving the *Outcome Explanation Problem*

Name	Ref.	Authors	Year	Explainer	Black Box	Data Type	General	Random	Examples	Code	Dataset
—	[134]	Xu et al.	2015	SM	DNN	IMG			✓	✓	✓
—	[30]	Fong et al.	2017	SM	DNN	IMG			✓		
CAM	[139]	Zhou et al.	2016	SM	DNN	IMG					
Grad-CAM	[106]	Selvaraju et al.	2016	SM	DNN	IMG					
—	[109]	Simonian et al.	2013	SM	DNN	IMG					
PWD	[7]	Bach et al.	2015	SM	DNN	IMG					
—	[113]	Sturm et al.	2016	SM	DNN	IMG					
DTD	[78]	Montavon et al.	2017	SM	DNN	IMG					
DeepLIFT	[107]	Shrikumar et al.	2017	FI	DNN	ANY					
CP	[64]	Landecker et al.	2013	SM	NN	IMG					
—	[143]	Zintgraf et al.	2017	SM	DNN	IMG					
VBP	[11]	Bojarski et al.	2016	SM	DNN	IMG					
—	[65]	Lei et al.	2016	SM	DNN	TXT					
ExplainD	[89]	Poulin et al.	2006	FI	SVM	TAB					
—	[29]	Strumbelj et al.	2010	FI	AGN	TAB	✓				
LIME	[98]	Ribeiro et al.	2016	FI	AGN	ANY	✓				
MES	[122]	Turner et al.	2016	DR	AGN	ANY	✓				
Anchors	[99]	Ribeiro et al.	2018	DR	AGN	ANY	✓				
—	[110]	Singh et al.	2016	DT	AGN	TAB	✓				
LORE	[37]	Guidotti et al.	2018	DR	AGN	TAB	✓				
MFI	[124]	Vidovic et al.	2016	FI	AGN	TAB	✓				
—	[39]	Haufe et al.	2014	FI	NLM	TAB					

Table 4. Summary of Methods for Opening Black Boxes Solving the *Model Inspection Problem*

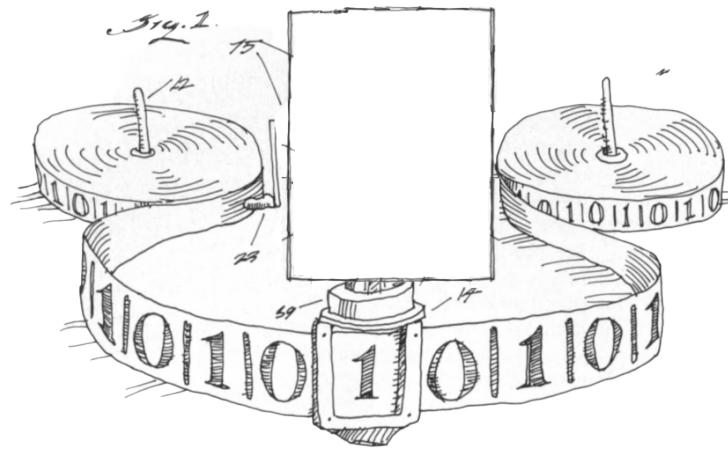
Name	Ref.	Authors	Year	Explainer	Black Box	Data Type	General	Random	Examples	Code	Dataset
NID	[83]	Olden et al.	2002	SA	NN	TAB			✓		
GDP	[8]	Baehrens	2010	SA	AGN	TAB	✓		✓		
QII	[24]	Datta et al	2016	SA	AGN	TAB	✓		✓		
IG	[115]	Sundararajan	2017	SA	DNN	ANY			✓		
VEC	[18]	Cortez et al.	2011	SA	AGN	TAB	✓		✓		
VIN	[42]	Hooker	2004	PDP	AGN	TAB	✓		✓		
ICE	[35]	Goldstein et al.	2015	PDP	AGN	TAB	✓		✓	✓	
Prospector	[55]	Krause et al.	2016	PDP	AGN	TAB	✓		✓		
Auditing	[2]	Adler et al.	2016	PDP	AGN	TAB	✓		✓	✓	
OPIA	[1]	Adebayo et al.	2016	PDP	AGN	TAB	✓		✓		
—	[136]	Yosinski et al.	2015	AM	DNN	IMG			✓		
IP	[108]	Shwartz et al.	2017	AM	DNN	TAB			✓		
—	[137]	Zeiler et al.	2014	AM	DNN	IMG	✓		✓		
—	[112]	Springenberg et al.	2014	AM	DNN	IMG			✓		
DGN-AM	[80]	Nguyen et al.	2016	AM	DNN	IMG			✓	✓	
—	[72]	Mahendran et al.	2016	AM	DNN	IMG			✓	✓	
—	[95]	Radford	2017	AM	DNN	TXT			✓		
—	[143]	Zintgraf et al.	2017	SM	DNN	IMG			✓	✓	
VBP	[11]	Bojarski et al.	2016	SM	DNN	IMG			✓	✓	
TreeView	[119]	Thiagarajan et al.	2016	DT	DNN	TAB			✓	✓	

# Transparencia/interpretabilidad

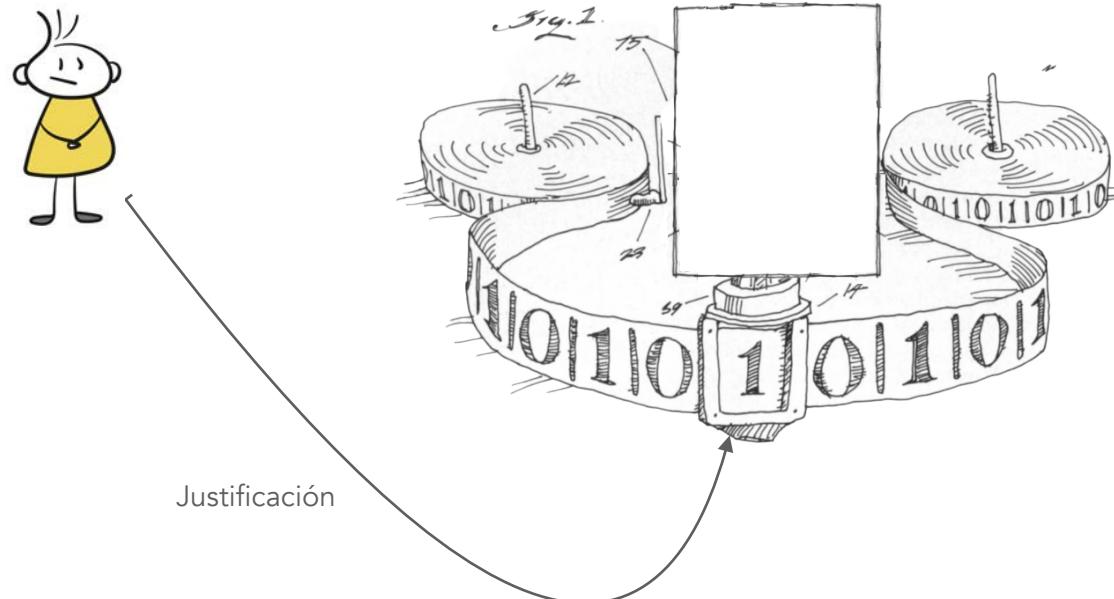
- Requiere “abrir” el algoritmo
  - i.e., rastrear el “path-dependency” de un resultado (Durán, 2021)
- La justification es asegurada mediante un third-party algorithm:
  - Interpretable predictors
  - Algunas formas de XAI (e.g., post-hoc explanation)
  - Transparency reports (e.g., Qualitative Input Influence)

# ¿Cómo se justifica via transparencia?

# Justificación via Transparencia (a sketch)



# Justificación via Transparencia (a sketch)

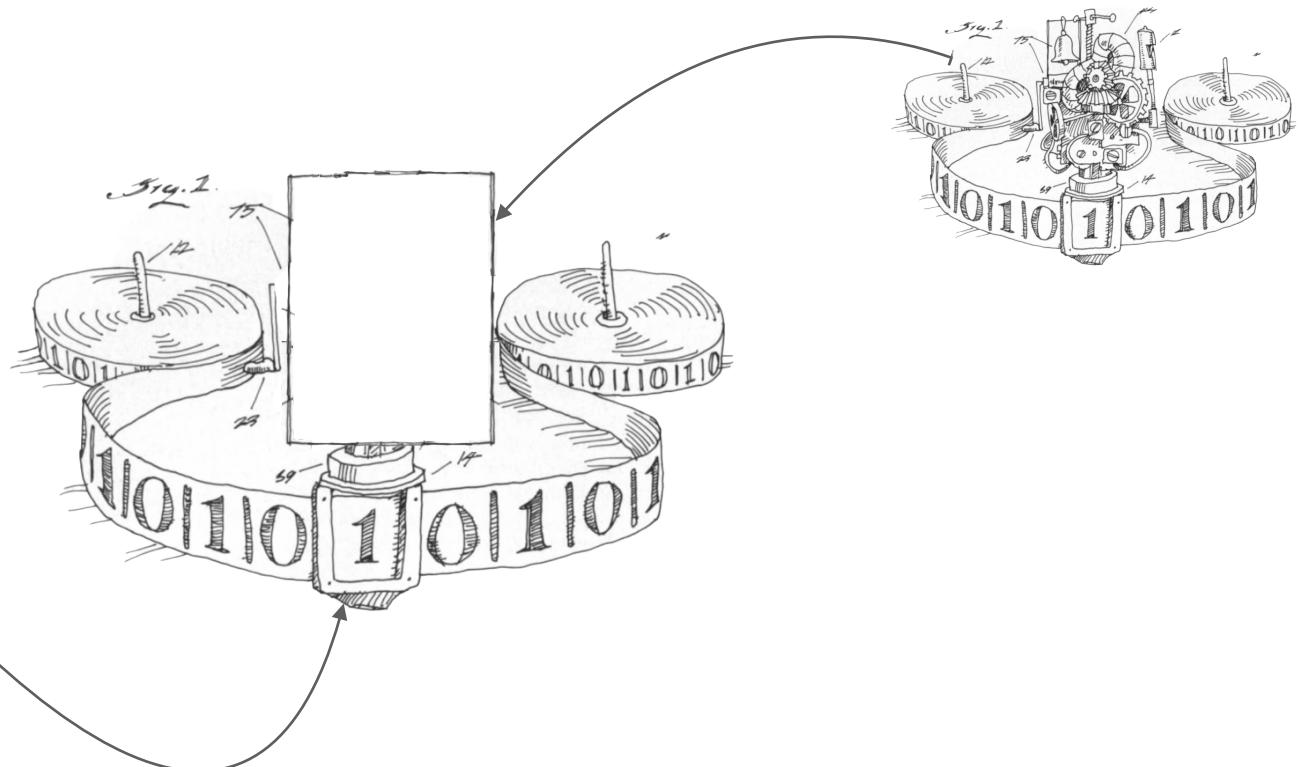


Justificación

# Justificación via Transparencia (a sketch)



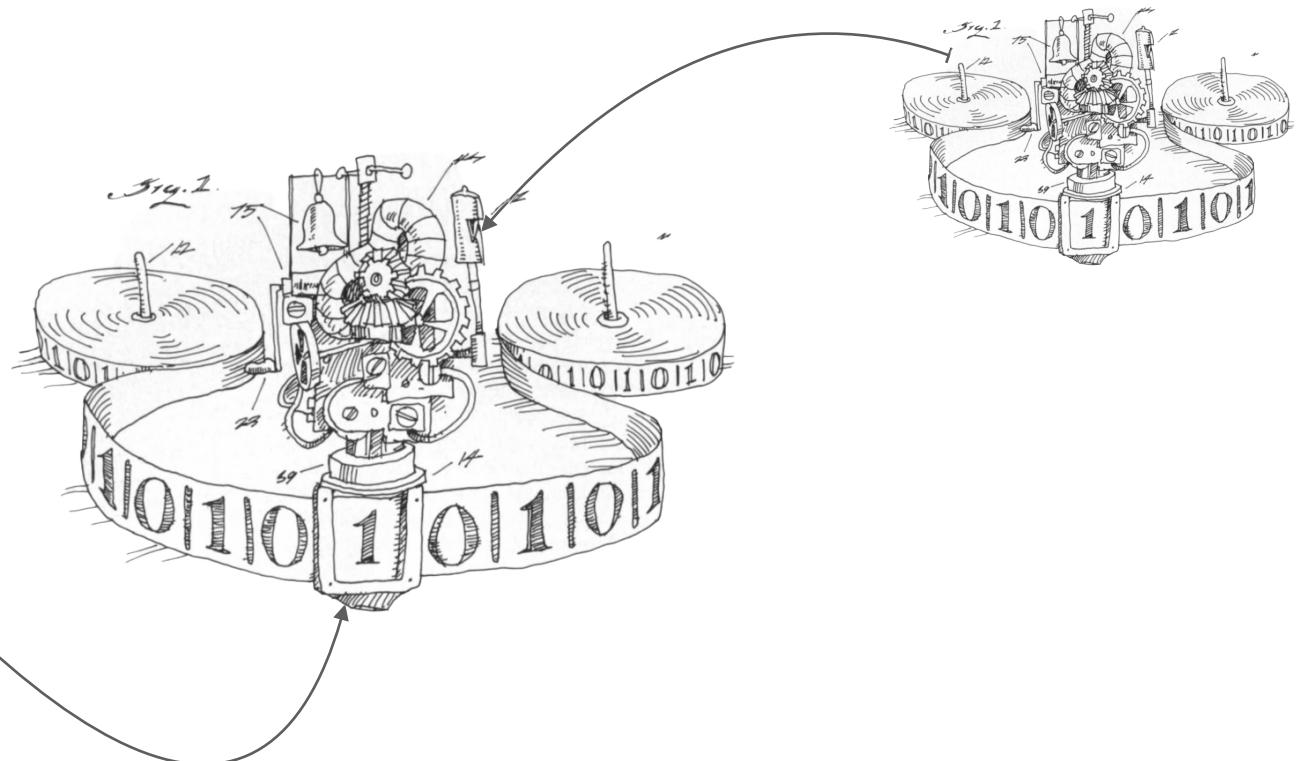
Justificación



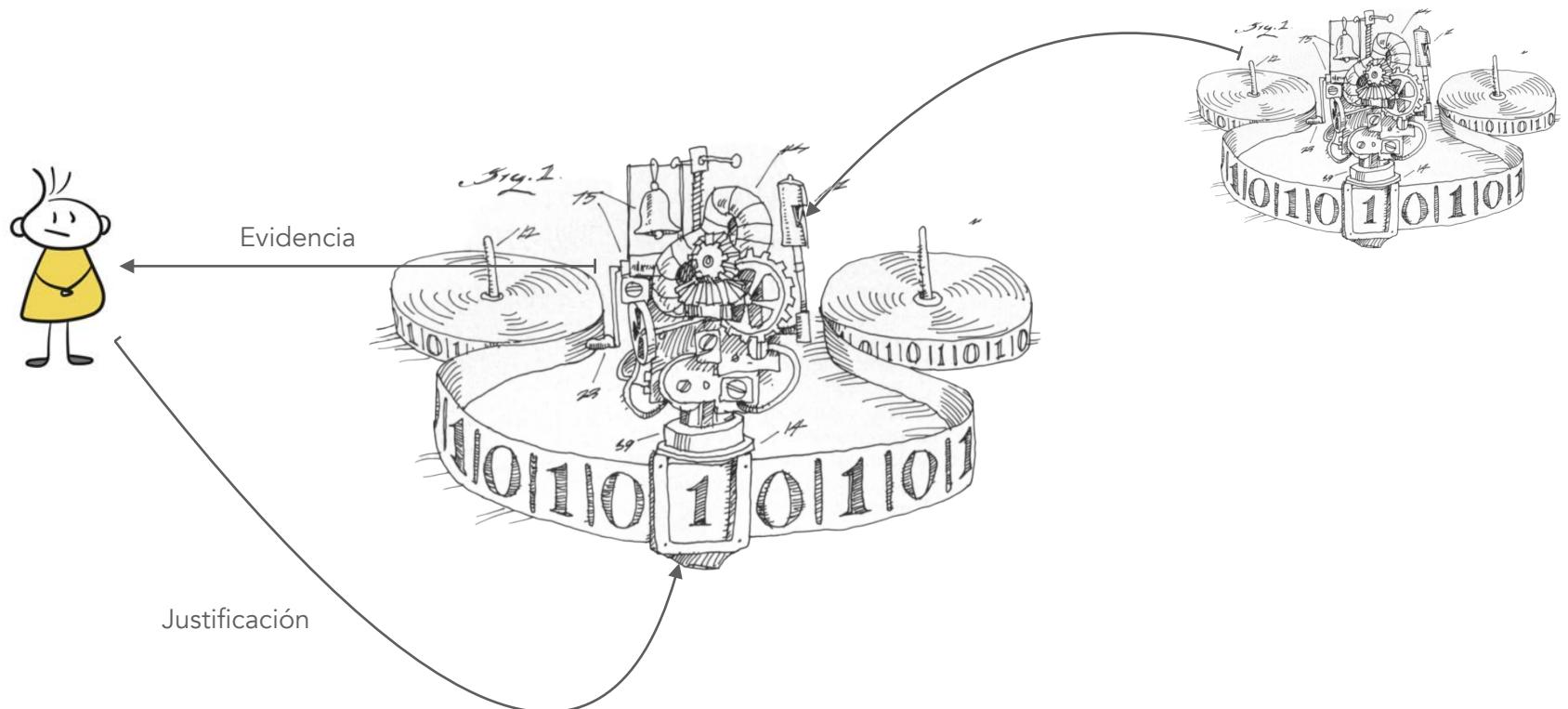
# Justificación via Transparencia (a sketch)



Justificación



# Justificación via Transparencia (a sketch)



# Against epistemic transparency in algorithmic science

Juan M. Durán



Against epistemic transparency of algorithms

Juan M. Durán

## Abstract

Epistemic transparency is often proposed as a solution to algorithmic opacity, wherein revealing the inner logic of an algorithm provides reasons or supporting evidence for the justification of its outputs. I argue that transparency is defective in belief formation and thus inadequate as an epistemology of algorithms. Two objections are developed: *transparency regress* and *bootstrapping*, both rooted in transparency's status as a time-sliced, 'outsourced' inferentialist epistemology. While weaker forms—such as contextual transparency—are briefly considered, the central argument remains that transparency fails to provide justification and should be abandoned as an epistemology of algorithms.

## 1 Introduction

It is widely accepted that algorithms are epistemically opaque. Opacity is here understood in line with Humphreys' definition, whereby the cognitive limitations of human agents prevent them from knowing the epistemically relevant elements of an algorithm that would justify belief in a given output [22, 618]. Now, there would be no special concern about opaque algorithms if it were not for the fact that they are central to many scientific endeavors, progressively displacing humans from the center of knowledge production. The opacity of algorithms, then, triggers significant epistemic anxiety.

To address opacity, many have sought to counteract this lack of knowledge. It is in this context that *transparency* gains traction. Creel, for instance, takes this idea in its strictest form. According to this author, opacity and transparency are “two sides of the same coin: opacity is a lack of transparency and vice versa” [28, 569, footnote 1]. As we will see, many have folded under this view. And while transparency encompasses a range of methodologies and strategies, there is a common justificatory aim, that is, to provide reasons or supporting evidence for believing the algorithm’s output. From where is this justification drawn? From identifying the functions, variables, values, and similar epistemically relevant elements within the algorithm that generate the output. This justificatory aim has been collectively referred to as “showing the inner logic of the algorithm” [7, 843]. For example, BenevolentAI identified *baricitinib* as an effective drug to combat COVID-19 symptoms. For a researcher to be justified in this belief, transparency shows how BenevolentAI internally operates: how it identifies drug-signal blocking, maps causal relations between viruses and drugs,

# De transparencia a fiabilismo (reliabilism)

"If we think in terms of such a process [i.e., algorithms] and imagine that its stepwise computation was slowed down to the point where, in principle, a human could examine each step in the process, the computationally irreducible process would become epistemically transparent. What this indicates is that the practical constraints we have previously stressed, primarily the need for computational speed, are the root cause of all epistemic opacity in this area. Because those constraints cannot be circumvented by humans, we must abandon the insistence on epistemic transparency for computational science. What replaces it would require an extended work in itself, but the prospects for success are not hopeless." (Humphreys, 2004, 150)

# ¿Qué es process reliabilism? ( $\approx$ Goldman 1979)

**Process reliabilism:** la creencia de Diego está justificada en caso que sea producida por un proceso (o secuencia de procesos) de formación de creencia fiable

Un proceso de formación de creencia fiable tiene la tendencia de producir creencias que son verdaderas antes que *falsas*



# ¿Qué es process reliabilism? ( $\approx$ Goldman 1979)

**Process reliabilism:** la creencia de Diego está justificada en caso que sea producida por un proceso (o secuencia de procesos) de formación de creencia fiable

Un proceso de formación de creencia fiable tiene la tendencia de producir creencias que son verdaderas antes que *falsas*



# CR + Instrumental reliabilism

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

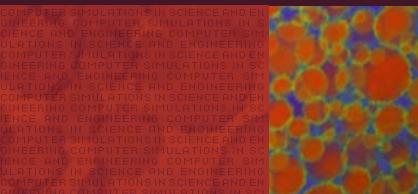


## Grounds for Trust: Essential Epistemic Opacity and Computational Reliabilism

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

### THE FRONTIERS COLLECTION



Juan Manuel Durán

## COMPUTER SIMULATIONS IN SCIENCE AND ENGINEERING

Concepts—Practices—Perspectives

Technology, Jaffalaan  
5, Nobelstraat 19,  
2628 BX Delft

## Beyond transparency: computational reliabilism as an externalist epistemology of algorithms

forthcoming in

*Philosophy of Science for Machine Learning: Core Issues and New Perspectives* – Juan M. Durán and Giorgia Pozzi (eds.)

Synthese Library

Juan M. Durán

**Abstract** This chapter examines the epistemology of algorithms, framing the discussion as a question of epistemic justification. Current approaches emphasize algorithmic transparency, which involves elucidating internal mechanisms—such as functions and variables—and demonstrating how (or that) these compute outputs. Thus, the mode of justification through transparency is contingent upon what is to be shown about the algorithm and, in this sense, is *internal* to it. In contrast, I propose an *externalist* epistemology of algorithms called *reliabilism* (CR). While I have previously developed CR in the context of computer simulations (160, 74, 112), this chapter extends the framework to algorithms used across scientific disciplines, particularly in machine learning and deep neural networks. At its core, CR posits that an algorithm's reliability is determined by its output, where reliability is defined by reliability indicators. These indicators range from formal standards, algorithmic expert competencies, research cultures, and other scientific practices to primary objectives to delineate the foundations of CR, explain mechanisms, and outline its potential as an externalist epistemology.

### 1 Introduction

The use of algorithms for scientific purposes is delivering remarkable examples of what suffices to illustrate this. In molecular biolog-

### Epistemic Opacity and Epistemic Inaccessibility

In this paper I shall revisit the concepts of epistemic opacity and essential epistemic opacity with the hope of clarifying and elaborating those concepts. I shall begin with a clarification of the definitions and relate the concept of epistemic opacity to that of epistemic inaccessibility. I then introduce and describe representational opacity, including three distinctions between types of representation and argue that this is an important source of opacity in some deep neural networks. Next, I proceed to an examination of other sources of epistemic opacity, some of which follow from differences between applied and pure mathematics. Then, after looking at some related concepts, I conclude with some ways in which opacity can be ameliorated.

#### 1. Introduction

As a reminder, here are the two definitions of epistemic opacity as formulated in Humphreys 2009:

A process is *epistemically opaque* relative to a cognitive agent X at time t just in case X does not know at all of the epistemically relevant elements of the process.

A process is *essentially epistemically opaque* to X if and only if it is impossible, given the nature of X, for X to know all of the epistemically relevant elements of the process.

European Journal for Philosophy of Science (2025) 15:37  
https://doi.org/10.1007/s13194-025-00664-2

### PAPER IN GENERAL PHILOSOPHY OF SCIENCE

## In defense of reliabilist epistemology of algorithms

Juan M. Durán<sup>1</sup>

Received: 19 August 2024 / Accepted: 3 June 2025  
© The Author(s) 2025

### Abstract

In a reliabilist epistemology of algorithms, a high frequency of accurate output representations is indicative of the algorithm's reliability. Recently, Humphreys challenged this assumption, arguing that reliability depends not only on frequency but also on the quality of output. Specifically, he contends that radical and egregious misrepresentations have a distinct epistemic impact on our assessment of an algorithm's reliability, regardless of the frequency of their occurrence. He terms these *statistically insignificant but serious errors* (SIS-Errors) and maintains that their occurrence warrants revoking our epistemic attitude towards the algorithm's reliability. This article argues against this challenge by presenting a series of counterexamples of algorithms against the challenge posed by Humphreys. It also argues that *instrumental reliabilism* as a foundational conditions designed to prevent SIS-Errors

algorithms - Computational reliabilism -

## Inductive Failures in Neural Nets: Why Reliabilism is an inappropriate epistemology for them



PAUL HUMPHREYS

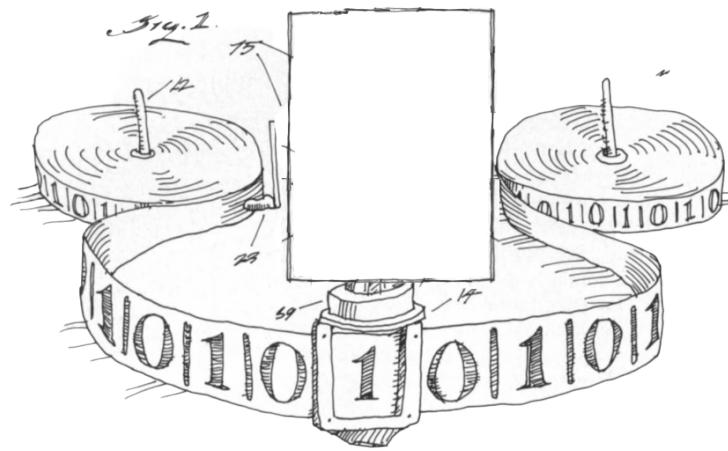
University of Virginia

at the alleged novelty of computer simulations argued that, although technologically routine, it is a philosophical novelty not "[a] revolution". Philosophers were worried about in the past" (Humphreys 2009).

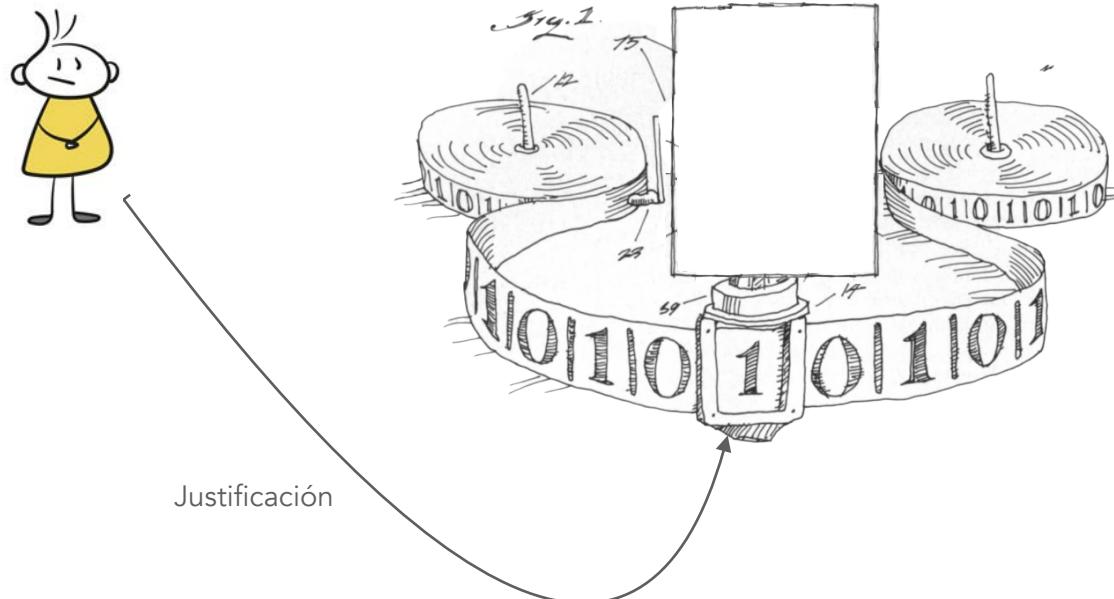
Technology, Jaffalaan 5, 2628 BX Delft,

# ¿Cómo se justifica via CR?

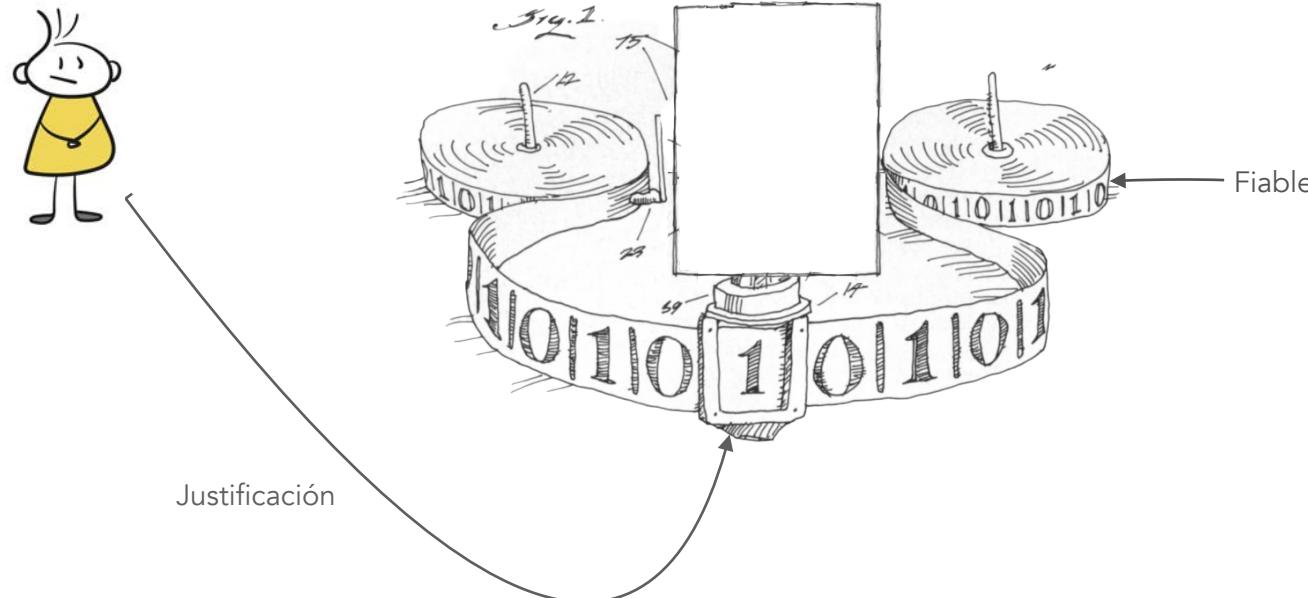
# Justificación via CR (a sketch)



# Justificación via CR (a sketch)



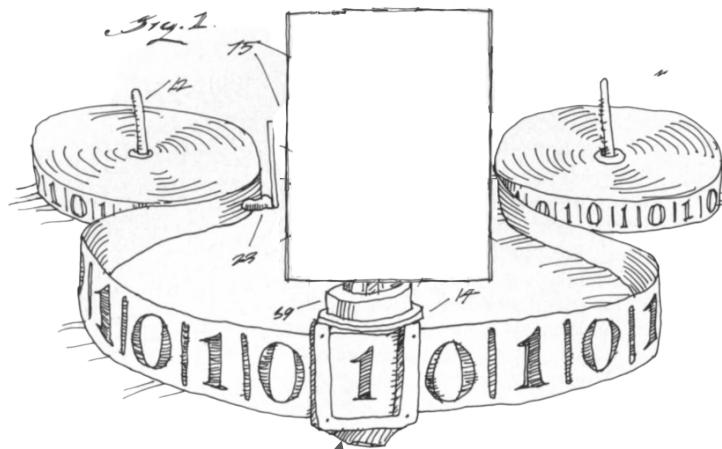
# Justificación via CR (a sketch)



# Justificación via CR (a sketch)



Justificación



Fiable

Indicadores  
de  
fiabilidad

# Fiabilidad computacional

- Aceptamos algoritmos de “caja negra”
- La justificación viene de asegurarse la fiabilidad del algoritmo a través de indicadores de fiabilidad

# Fiabilidad computacional

- Aceptamos algoritmos de “caja negra”
- La justificación viene de asegurarse la fiabilidad del algoritmo a través de indicadores de fiabilidad
  - RI<sub>1</sub> Performance de los algoritmos
  - RI<sub>2</sub> Práctica científica basada en algoritmos
  - RI<sub>3</sub> Construcción social de la fiabilidad

# Somes obvious sources of reliabilty

- RI<sub>1</sub> Performance de los algoritmos
- RI<sub>2</sub> Práctica científica basada en algoritmos
- RI<sub>3</sub> Construcción social de la fiabilidad

# Somes obvious sources of reliabilty

- RI<sub>1</sub> Performance de los algoritmos

# Somes obvious sources of reliabilty



## Rl<sub>1</sub> Performance de los algoritmos

- Métodos de verificación y validación (→ accuracy)
- Elecciones de datos y bases de datos
- Parametrización e hyper-parametrización
- Prácticas de entrenamiento de modelos
- Una historia de éxitos y fracasos en las implementaciones
- Tratamiento de errores (e.g., uncertainty quantification)

(Durán 2018; Durán & Formanek 2018)

# Somes obvious sources of reliabilty

- RI<sub>1</sub> Performance de los algoritmos
- RI<sub>2</sub> Práctica científica basada en algoritmos

# Somes obvious sources of reliability

— RI<sub>1</sub> Performance de los algoritmos

— RI<sub>2</sub> Práctica científica basada en algoritmos

- Causalidad y otros mecanismos implementados
- Valores epistémicos, morales, etc.
- Representación, idealizaciones, abstracciones, etc
- Teorías, leyes, principios, etc
- Algún grado de coherencia teórica, conceptual, etc
  - Conocimiento de fondo (background knowledge, e.g., knowledge of the learned weights)
  - Algún morfismo entre hipótesis, modelos, etc y la sintaxis del algoritmo (e.g., al usar natural kinds, cut-off-values que son aceptados, definiciones — e.g., cuándo un melanoma es cancerígeno)

# Somes obvious sources of reliabilty

- RI<sub>1</sub> Performance de los algoritmos
- RI<sub>2</sub> Práctica científica basada en algoritmos
- RI<sub>3</sub> Construcción social de la fiabilidad

# Somes obvious sources of reliabilty

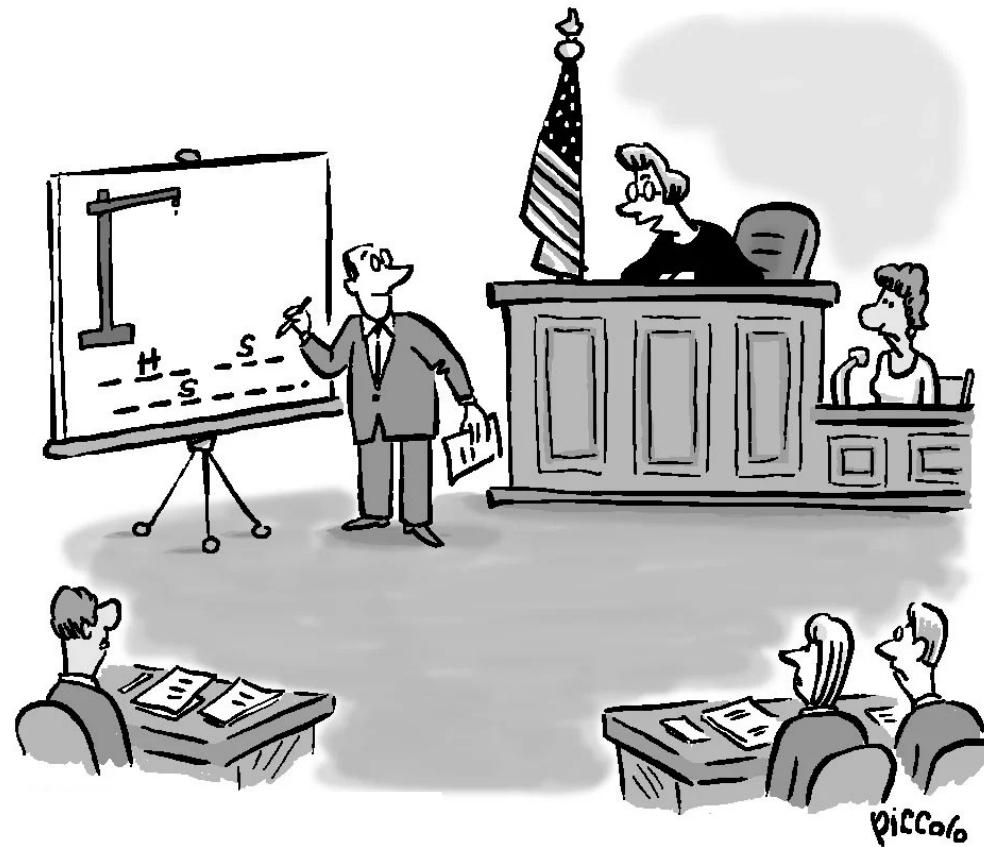
- RI<sub>1</sub> Performance de los algoritmos
- RI<sub>2</sub> Práctica científica basada en algoritmos
- RI<sub>3</sub> Construcción social de la fiabilidad
  - Aceptabilidad social del resultado del algoritmo
  - Favorable, exitoso/fracaso, usos (in)válidos del algoritmo y sus resultados (e.g., cuando se desarrollan nuevas hipótesis o se expande en nuevas estrategias)
  - Coherencia con un cuerpo de conocimiento estable
  - Coherencia con los compromisos epistémicos, éticos, etc de los científicos y sus comunidades
  - Realización de los valores epistémicos, éticos, etc
  - Culturas de investigación (científico, modelaje, comunidades...)

# ¿Cómo funciona CR? Un caso en ciencia forense Parte I

# El experto, el abogado, y el juez

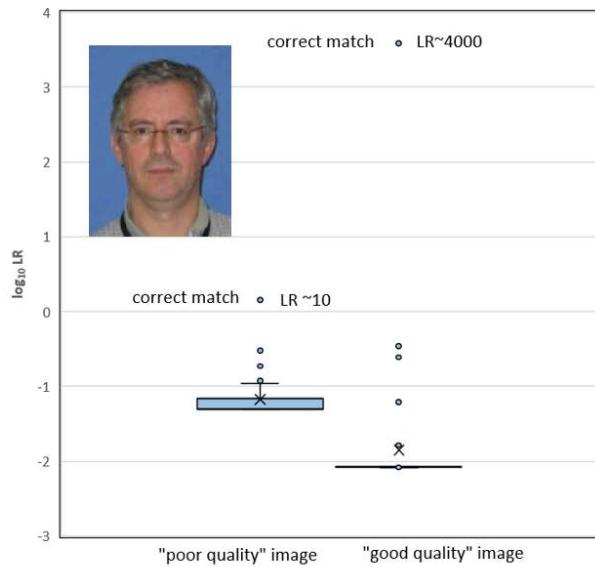


Nederlands Forensisch Instituut  
*Ministerie van Justitie en Veiligheid*



# Comparación de rostros en ciencia forense

"Poor quality" image



"Good quality" image





# Nederlands Forensisch Instituut

## Ministerie van Justitie en Veiligheid

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

	Zaak_01 (2013)	Zaak_04 (2012)	Zaak_05 (2013)	Zaak_07 (2013)	Zaak_13 (2013)	Zaak_14 (2013)	Zaak_15 (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 gezichtbeelddvergelijkingen. Groen: geselecteerd voor gedetailleerde vergelijking. Rood: op basis van voorselectie uitgesloten.

	Zaak_01 (2013)	Zaak_04 (2012)	Zaak_05 (2013)	Zaak_07 (2013)	Zaak_13 (2013)	Zaak_14 (2013)	Zaak_15 (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

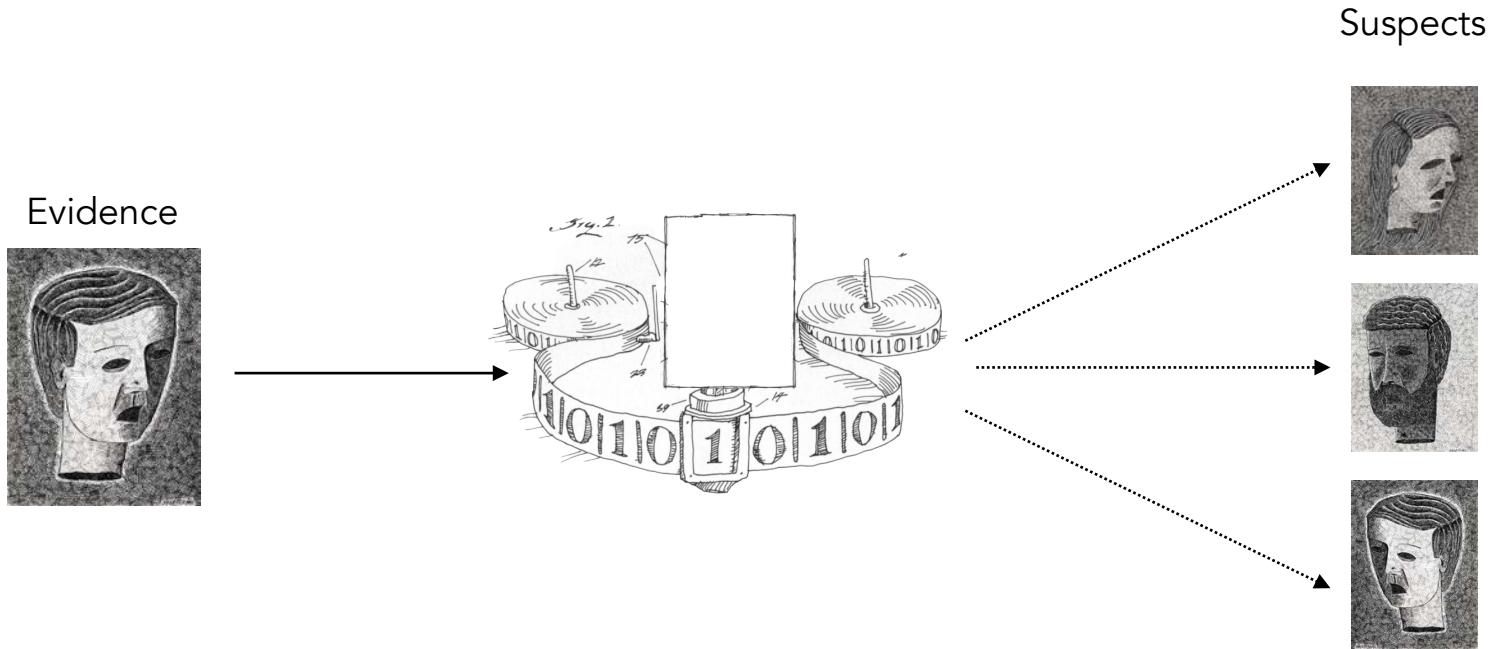
Zaaknr.	Rapportage aanvraag 002 (25 maart 2016)	Huidige aanvraag 003
Zaak 01	Waarschijnlijker	Waarschijnlijker
Zaak 04	Zeer veel waarschijnlijker	Zeer veel waarschijnlijker
Zaak 05	Zeer veel waarschijnlijker	Zeer veel waarschijnlijker
Zaak 07	--	Iets waarschijnlijker
Zaak 13	Veel waarschijnlijker	Veel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

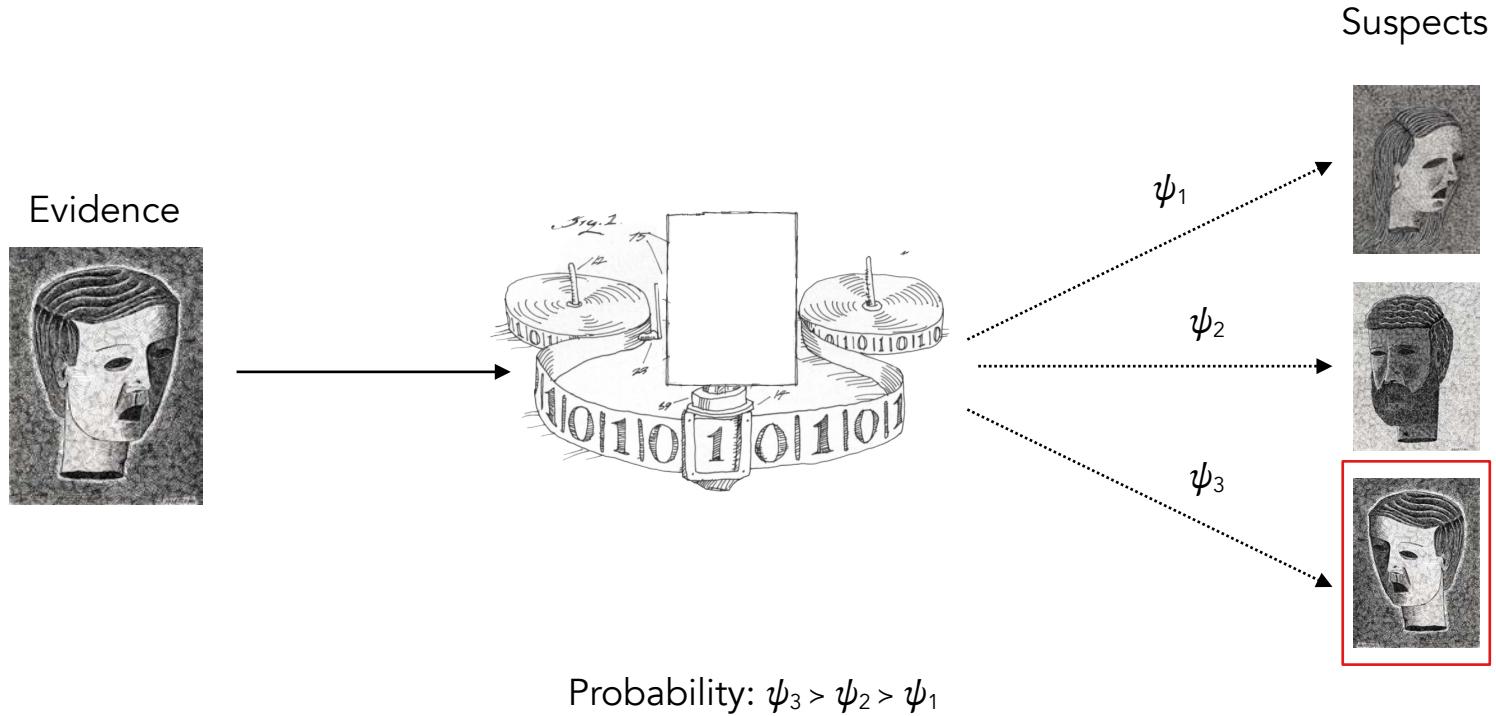
Zaak	Verdachte	a	b	c	Consensus	Zaak	Verdachte	a	b	c	Consensus
01	Sofia	2	2	2	2	13	Sofia	3	3	2	3
01	Diana	-3	-3	-3	-3	13	Diana	-3	-3	-3	-3
01	Djanija	-4	-3	-4	-4	13	Djanija	-3	-2	-4	-3
01	Vahida (jonger)	-2	-2	-4	-3	13	Vahida (ouder)	-3	-2	-3	-3
01	Vahida (ouder)	-3	1	-4	-2	14	Sofia	1	2	2	2
01	Vera	-4	-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									

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	extrem veel waarschijnlijker wanneer H2 waar is dan wanneer H1 waar is
H2 waar is		

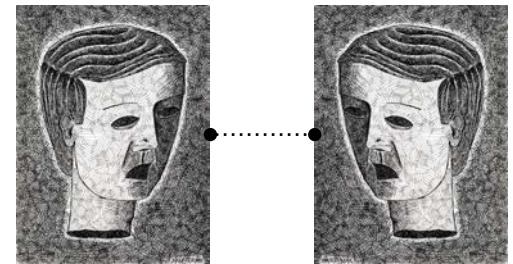
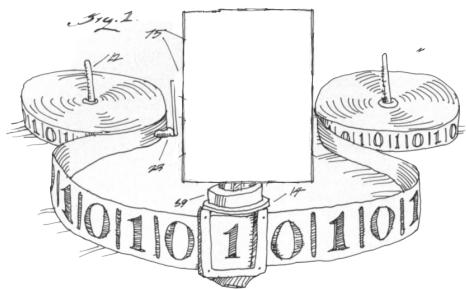
# Ciencia forense y AI



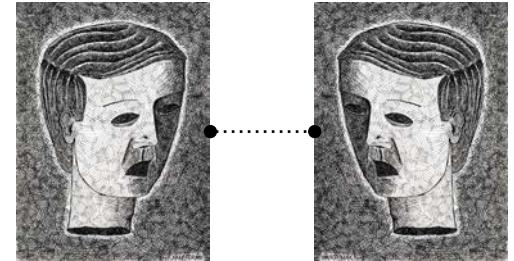
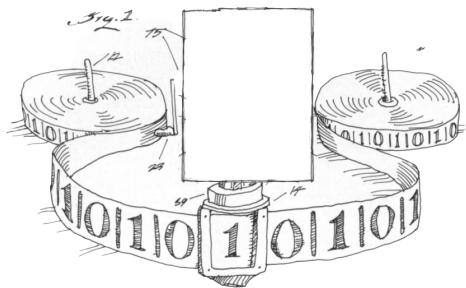
# Ciencia forense y AI



# Ciencia forense y CR



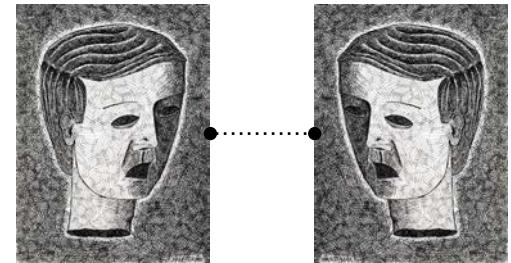
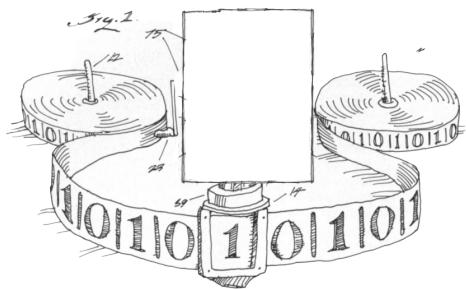
# Ciencia forense y CR



## RI<sub>1</sub> Performance de los algoritmos

- Verificación y Validación
- Análisis de robustez
- Redundancia
- "Calidad" de los datos

# Ciencia forense y CR



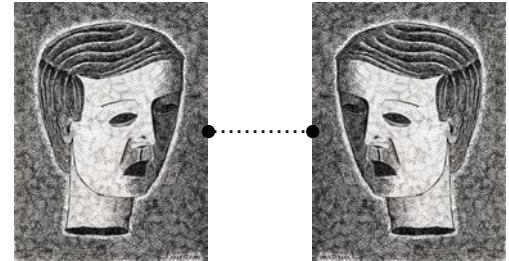
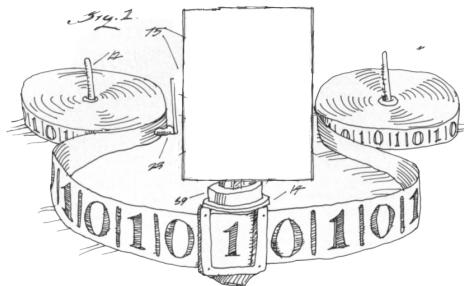
## RI<sub>1</sub> Performance de los algoritmos

- Verificación y Validación
- Análisis de robustez
- Redundancia
- “Calidad” de los datos

## RI<sub>2</sub> Práctica científica basada en algoritmos

- Implementación de “quality assurance methods”
- “We implement well-known and tested libraries”
- Técnicas biométricas

# Ciencia forense y CR



## R<sub>I</sub><sub>1</sub> Performance de los algoritmos

- Verificación y Validación
- Análisis de robustez
- Redundancia
- "Calidad" de los datos

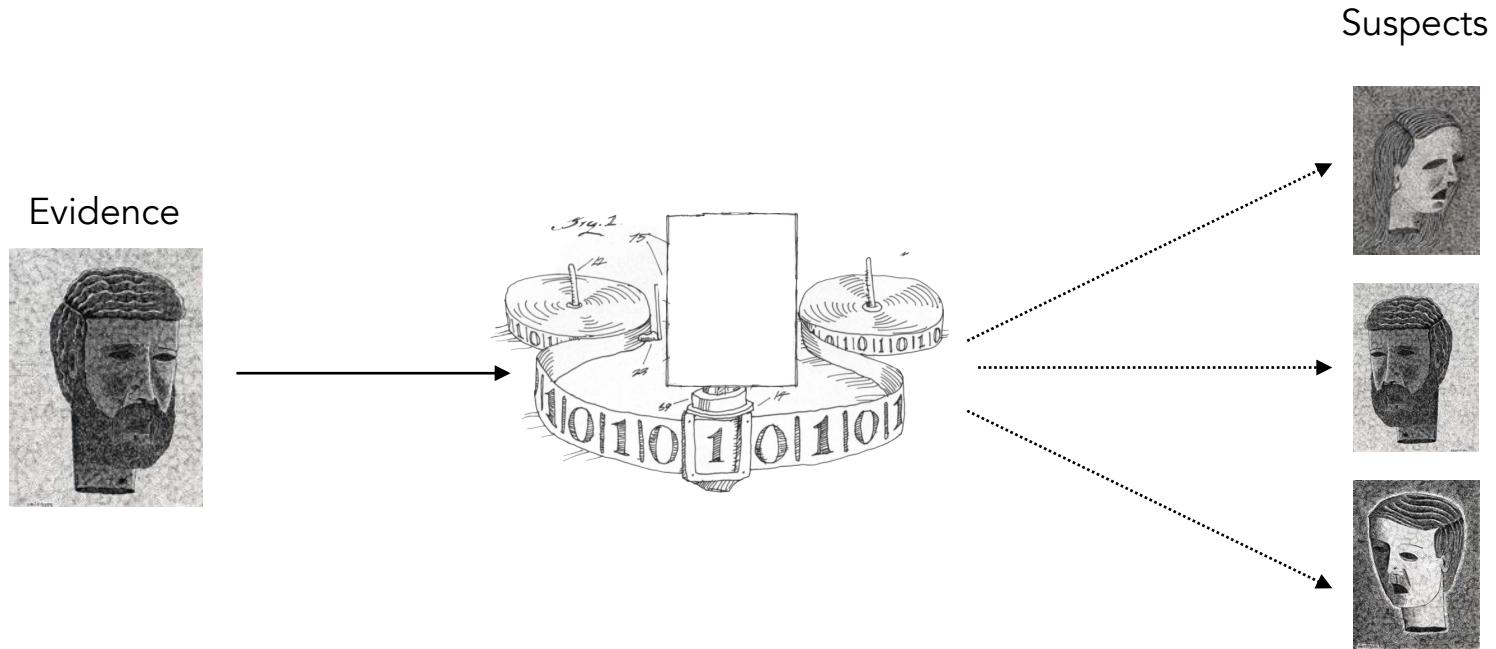
## R<sub>I</sub><sub>2</sub> Práctica científica basada en algoritmos

- Implementación de "quality assurance methods"
- "We implement well-known and tested libraries"
- Técnicas biométricas

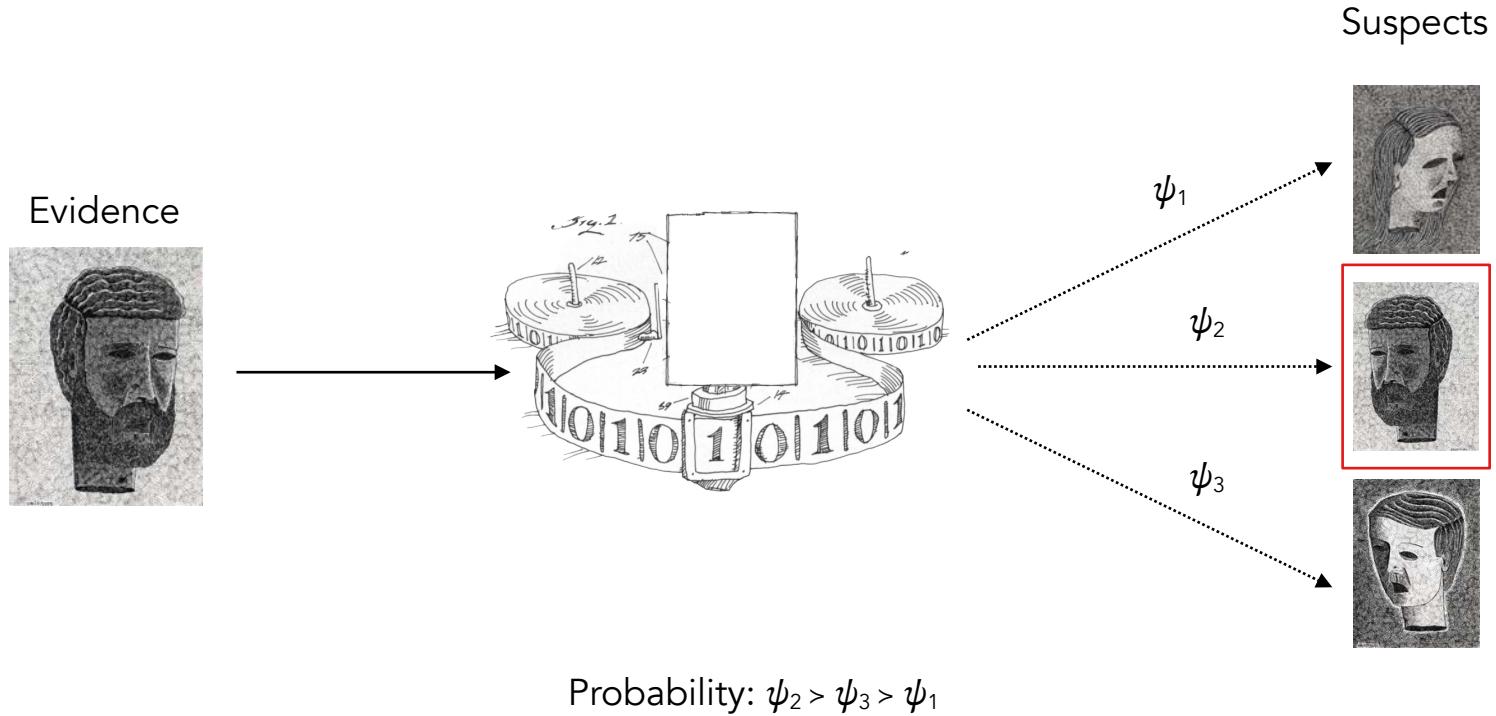
## R<sub>I</sub><sub>3</sub> Construcción social de la fiabilidad

- Debates sobre la aplicabilidad de los resultados (e.g., soundness, legal, etc.)
- Contraste con otros métodos no-computacionales (e.g., uso de revisiones de código, reportes de expertos)

# Ciencia forense y CR



# Ciencia forense y CR



# ¿Cómo funciona CR? Un caso en ciencia forense Parte II

# 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.

# Automated Inference on Criminality Using Face Images

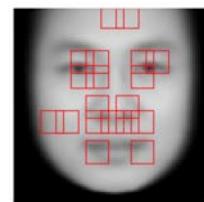
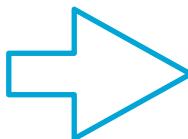


(a) Three samples in criminal ID photo set  $S_c$ .

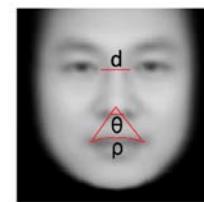


(b) Three samples in non-criminal ID photo set  $S_n$

Figure 1. Sample ID photos in our data set.



(a)



(b)

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

# Automated Inference on Criminality Using Face Images

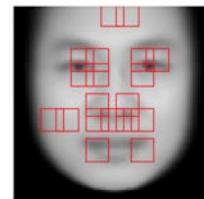
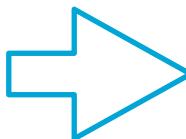


(a) Three samples in criminal ID photo set  $S_c$ .

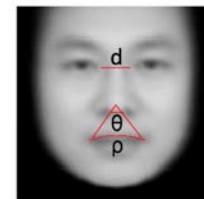


(b) Three samples in non-criminal ID photo set  $S_n$

Figure 1. Sample ID photos in our data set.



(a)



(b)

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

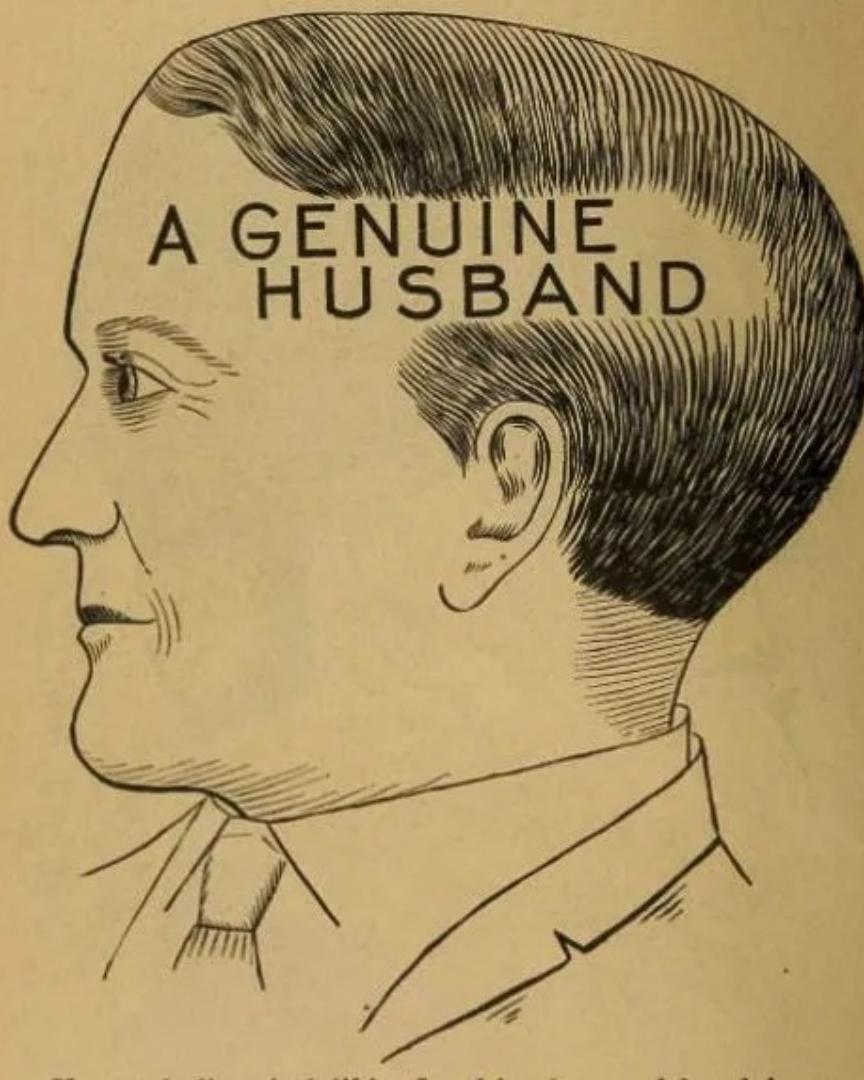


(a)



(b)

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



A GENUINE  
HUSBAND



AN UNRELIABLE  
HUSBAND

# 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.

"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)

# 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.

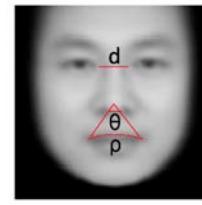
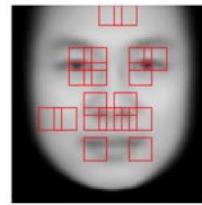
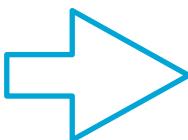


Figure 4. (a) FGM results; (b) Three discriminative features  $\rho$ ,  $d$  and  $\theta$ .



(a)



(b)

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

# Automated Inference on Criminality Using Face Images

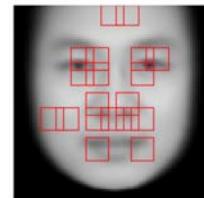
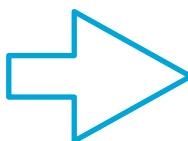


(a) Three samples in criminal ID photo set  $S_c$ .

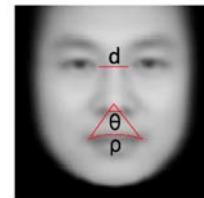


(b) Three samples in non-criminal ID photo set  $S_n$

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.

## RI<sub>1</sub> Performance de los algoritmos

- Verificación y Validación (?)
- “Calidad” de los datos
- El sistema es forzado a elegir entre dos clases: {criminal}/{no-criminal}

# Automated Inference on Criminality Using Face Images

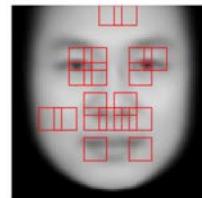
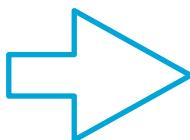


(a) Three samples in criminal ID photo set  $S_c$ .

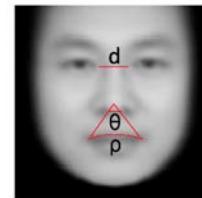


(b) Three samples in non-criminal ID photo set  $S_n$

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.

## RI<sub>1</sub> Performance de los algoritmos

- Verificación y Validación (?)
- “Calidad” de los datos
- El sistema es forzado a elegir entre dos clases: {criminal}/{no-criminal}

## RI<sub>2</sub> Práctica científica basada en algoritmos

- Fossiliza conceptos tales como “criminal”
- Está desconectada de un cuerpo de conocimiento:
  - La construcción social de la criminalidad
  - Las bases socio-económicas
  - La psicología de los autores de crímenes

# Automated Inference on Criminality Using Face Images

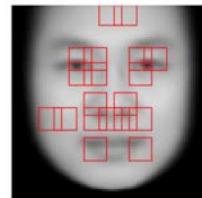
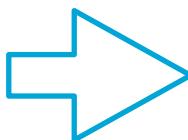


(a) Three samples in criminal ID photo set  $S_c$ .

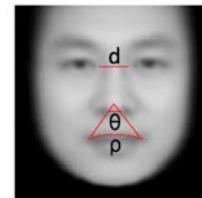


(b) Three samples in non-criminal ID photo set  $S_n$

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.

## RI<sub>1</sub> Performance de los algoritmos

- Verificación y Validación (?)
- “Calidad” de los datos
- El sistema es forzado a elegir entre dos clases: {criminal}/{no-criminal}

## RI<sub>2</sub> Práctica científica basada en algoritmos

- Fossiliza conceptos tales como “criminal”
- Está desconectada de un cuerpo de conocimiento:
  - La construcción social de la criminalidad
  - Las bases socio-económicas
  - La psicología de los autores de crímenes

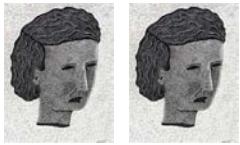
## RI<sub>3</sub> Construcción social de la fiabilidad

- Falla en dar cuenta de, o incluir valores epistémicos, éticos, etc (e.g., la presunción de inocencia)

# Un límite para CR

# SISE-errors

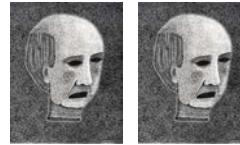
A[1]



A[2]



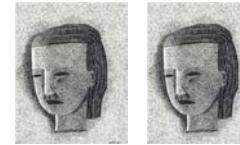
A[3]



A[4]



A[5]



...

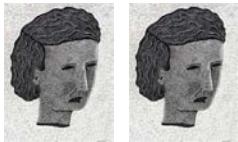
A[n]



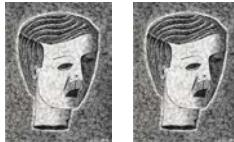
# SISE-errors

SIS-Errors: Statistically insignificant, but serious errors pueden debilitar la fidelidad del algoritmo  
(warrant revoking the reliability of the algorithm)

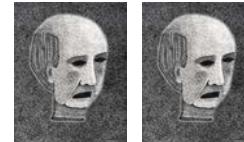
A[1]



A[2]



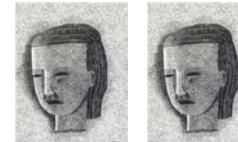
A[3]



A[4]

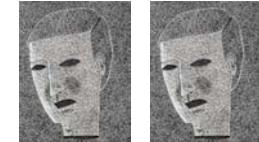


A[5]



...

A[n]



SIS-Errors: the algorithm's output does not match, to the degree permissible by the relevant community, a given representation (Humphreys, 2021)

# Tipos de Errores (y de cómo los aborda el fiabilista)

# SIS-Errors: Errores aleatorios

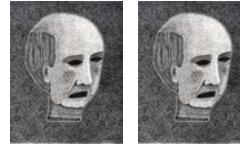
A[1]



A[2]



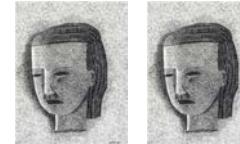
A[3]



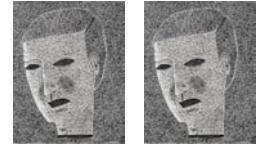
A[4]



A[5]



A[n]



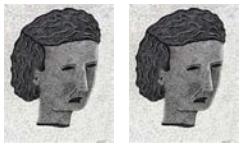
...

Errores aleatorios:

Son errores que no pueden anticiparse, que son impredecibles, y que son usualmente ejecuciones fallidas del algoritmo o los datos que se procesan. Estos errores no son sistemáticos ni consistentes, sino que ocurren esporádicamente y son muy difíciles de reproducir (e.g., bit flips in memory due to cosmic rays; randomised initiations  $\rightarrow$  non-deterministic ML outputs)

# SIS-Errors: Errores aleatorios

A[1]



A[2]



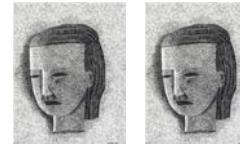
A[3]



A[4]

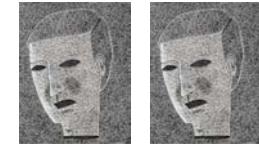


A[5]



...

A[n]



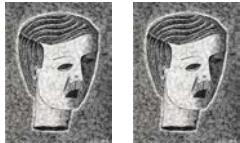
**Anti “epistemic bad luck” condition:** Permitase que un algoritmo se ejecute al menos dos veces bajo condiciones similares (i.e., variables de entrada, parámetros, configuraciones de datos, metaparámetros, etc.). Si el algoritmo produce SIS-Errors idénticos (o considerados lo suficientemente idénticos) en ejecuciones independientes, entonces la probabilidad de que estos resultados sean casos de “epistemic bad luck” es insignificante. Por lo tanto, tal recurrencia proporciona un fuerte respaldo a que el error es sistemático y no aleatorio.

# SIS-Errors: Errores sistemáticos

A[1]



A[2]



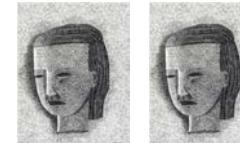
A[3]



A[4]

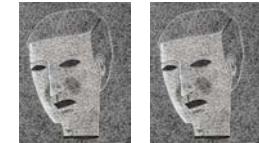


A[5]



...

A[n]



Errores sistemáticos:

Son errores que no son arbitrarios, y que se reproducen a lo largo de distintas ejecuciones del mismo algoritmo. Con las herramientas apropiadas, errores sistemáticos se pueden identificar y medir (e.g., rounding off errors)

# La estrategia

- I. Hacer una estrategia epistemológica, entendida como condiciones de posibilidad para la fiabilidad. No necesariamente tiene correlato metodológico o aplicable en la práctica.

# La estrategia

- I. Hacer una estrategia epistemológica, entendida como condiciones de posibilidad para la fiabilidad. No necesariamente tiene correlato metodológico o aplicable en la práctica.
- II. Conectar errores de representación con tipos de indicadores de fidelidad.
  - i. Para poder evitar tener que revocar la fiabilidad de un algoritmo, nuestra mejor estrategia es identificar IF que han fallado (inadecuados, incorrectos, o ausentes)

# Clases de errores sistemáticos

"Once you unleash it on large data, deep learning has its own dynamics, it does its own repair and its own optimisations, and it gives you the right results most of the time. But when it doesn't, you don't have a clue about what went wrong and what should be fixed. In particular, you'd not know if the fault is in the program, in the method, or because things have changed in the environment."

(Pearl, 2019, 15)

# Clases de errores sistemáticos

"Once you unleash it on large data, deep learning has its own dynamics, it does its own repair and its own optimisations, and it gives you the right results most of the time. But when it doesn't, you don't have a clue about what went wrong and what should be fixed. In particular, you don't know if the fault is in the program, in the method, or because things have changed in the environment."

(Pearl, 2019, 15)

- Errores de clase 1: errores que ocurren cuando el algoritmo produce resultados incorrectos dado calcular mal (e.g., rounding off errors, overflow, etc)

# Clases de errores sistemáticos

"Once you unleash it on large data, deep learning has its own dynamics, it does its own repair and its own optimisations, and it gives you the right results most of the time. But when it doesn't, you don't have a clue about what went wrong and what should be fixed. In particular, you don't know if the fault is in the program, in the method, or because things have changed in the environment."

(Pearl, 2019, 15)

- Errores de clase<sub>1</sub>: errores que ocurren cuando el algoritmo produce resultados incorrectos dado calcular mal (e.g., rounding off errors, overflow, etc)
- Errores de clase<sub>2</sub>: errores que surgen cuando los métodos o técnicas implementadas en el algoritmo son inapropiadas para un objetivo en particular (e.g., usar el "sorting algorithm" incorrecto, o la implementación de conceptos núcleo {criminal/no-criminal})

# Clases de errores sistemáticos

"Once you unleash it on large data, deep learning has its own dynamics, it does its own repair and its own optimisations, and it gives you the right results most of the time. But when it doesn't, you don't have a clue about what went wrong and what should be fixed. In particular, you don't know if the fault is in the program, in the method, or because things have changed in the environment."

(Pearl, 2019, 15)

- Errores de clase<sub>1</sub>: errores que ocurren cuando el algoritmo produce resultados incorrectos dado calcular mal (e.g., rounding off errors, overflow, etc)
- Errores de clase<sub>2</sub>: errores que surgen cuando los métodos o técnicas implementadas en el algoritmo son inapropiadas para un objetivo en particular (e.g., usar el "sorting algorithm" incorrecto, o la implementación de conceptos núcleo {criminal/no-criminal})
- Errores de clase<sub>3</sub>: errores que ocurren cuando el entorno hace el algoritmo y sus resultados irrelevantes (e.g., detección facial en los tiempos del COVID-19)

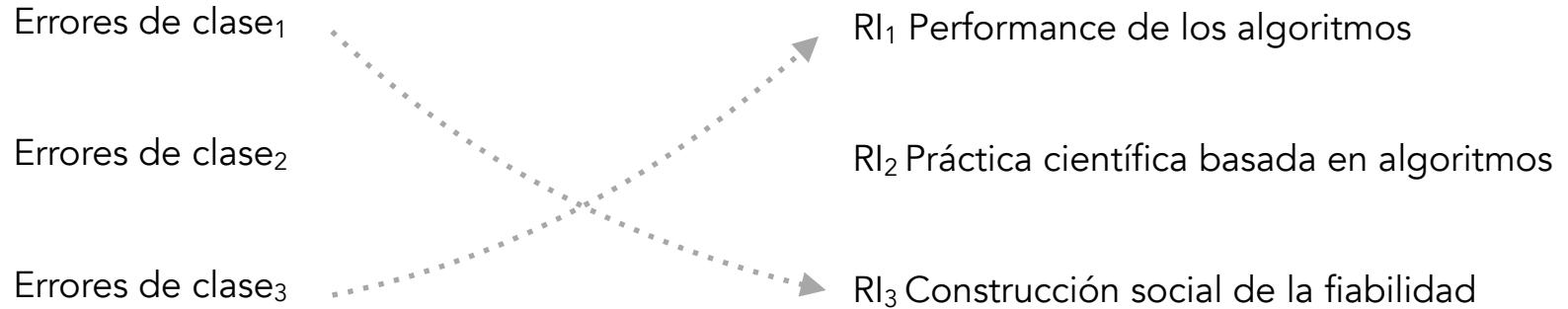
# Errores e Indicadores de Fiabilidad

Errores de clase<sub>1</sub> ..... ➔ RI<sub>1</sub> Performance de los algoritmos

Errores de clase<sub>2</sub> ..... ➔ RI<sub>2</sub> Práctica científica basada en algoritmos

Errores de clase<sub>3</sub> ..... ➔ RI<sub>3</sub> Construcción social de la fiabilidad

# Errores e Indicadores de Fiabilidad



# Errores e Indicadores de Fiabilidad

Errores de clase<sub>1</sub> ..... inadecuados ⤵ RI<sub>1</sub> Performance de los algoritmos

Errores de clase<sub>2</sub> ..... incorrectos ⤵ RI<sub>2</sub> Práctica científica basada en algoritmos

Errores de clase<sub>3</sub> ..... ausentes ⤵ RI<sub>3</sub> Construcción social de la fiabilidad

# IF inadecuados

**SPRINGER NATURE** Link

[Find a journal](#)   [Publish with us](#)   [Track your research](#)

 Search

[Home](#) > [Advances in Biometrics](#) > Conference paper

## Why Is Facial Occlusion a Challenging Problem?

Conference paper

pp 299–308 | [Cite this conference paper](#)

[Download book PDF](#) 



**Advances in Biometrics**  
(ICB 2009)

# IF incorrectos

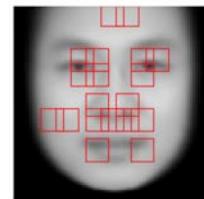
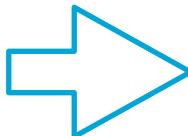


(a) Three samples in criminal ID photo set  $S_c$ .

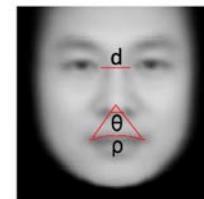


(b) Three samples in non-criminal ID photo set  $S_n$ .

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)

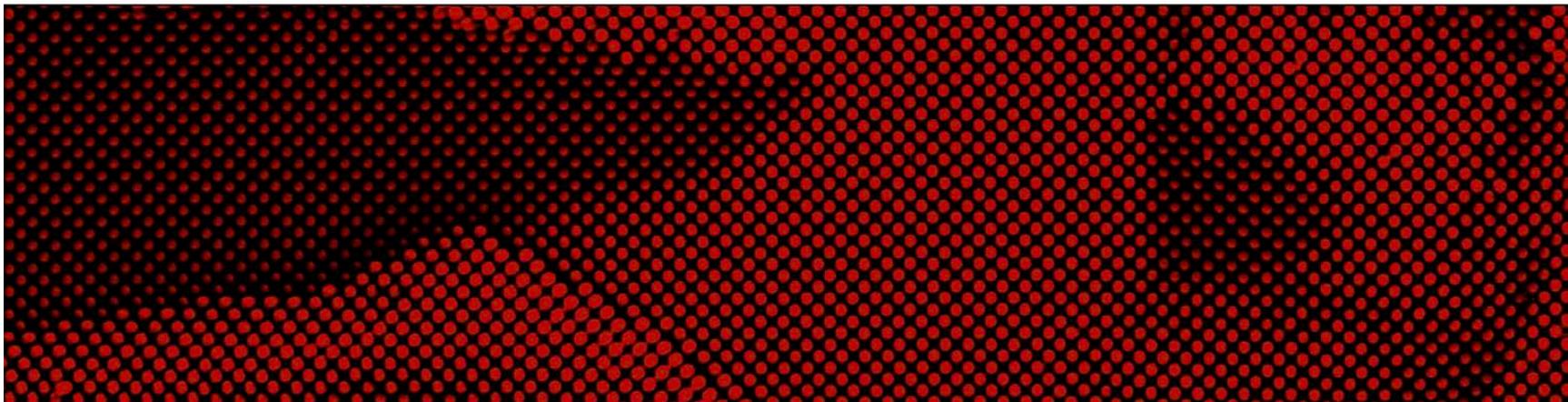
Figure 9. (a) The four subtypes of criminal faces; (b) The three subtypes of non-criminal faces.

# IF ausentes

≡ WIRED

SECURITY POLITICS THE BIG STORY BUSINESS SCIENCE CULTURE REVIEWS

👤 NEWSLETTERS ⚙



VIDEO: SAM CANNON

MAIA SZALAVITZ THE BIG STORY AUG 11, 2021 6:00 AM

## The Pain Was Unbearable. So Why Did Doctors Turn Her Away?

A sweeping drug addiction risk algorithm has become central to how the US handles the opioid crisis. It may only be making the crisis worse.

# La estrategia

- I. Hacer una estrategia epistemológica, entendida como condiciones de posibilidad para la fiabilidad. No necesariamente tiene correlato metodológico o aplicable en la práctica.
- II. Conectar errores de representación con tipos de indicadores de fidelidad.
  - i. Para poder evitar tener que revocar la fiabilidad de un algoritmo, nuestra mejor estrategia es identificar IF que han fallado (inadecuados, incorrectos, o ausentes)
- III. Establezco cláusulas de fiabilidad que, de violarse, nos obligan a suspender nuestras creencias —y así resguardar la fiabilidad (i.e., no tener que revocarla)

## SIS-Errors

Errores de clase<sub>1</sub>

SIS-Errors ..... ➔ Errores de clase<sub>2</sub>

Errores de clase<sub>3</sub>

Errores de clase<sub>1</sub> ..... ➔ RI<sub>1</sub> Robustez técnica de los algoritmos

SIS-Errors ..... ➔ Errores de clase<sub>2</sub> ..... ➔ RI<sub>2</sub> Práctica científica basada en algoritmos

Errores de clase<sub>3</sub> ..... ➔ RI<sub>3</sub> Construcción social de la fiabilidad

Errores de clase<sub>1</sub> ..... inadecuados ..... → RI<sub>1</sub> Robustez técnica de los algoritmos

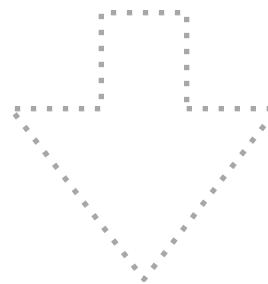
SIS-Errors ..... → Errores de clase<sub>2</sub> ..... incorrectos ..... → RI<sub>2</sub> Práctica científica basada en algoritmos

Errores de clase<sub>3</sub> ..... ausentes ..... → RI<sub>3</sub> Construcción social de la fiabilidad

Errores de clase<sub>1</sub> ..... inadecuados ..... RI<sub>1</sub> Robustez técnica de los algoritmos

SIS-Errors ..... ➔ Errores de clase<sub>2</sub> ..... incorrectos ..... ➔ RI<sub>2</sub> Práctica científica basada en algoritmos

Errores de clase<sub>3</sub> ..... ausentes ..... ➔ RI<sub>3</sub> Construcción social de la fiabilidad

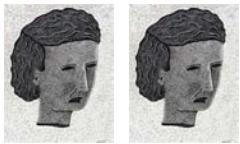


Mantener fiabilidad?

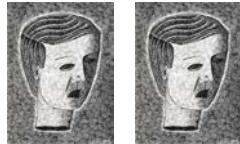
condiciones para suspender, revisar o anular la formación de creencia

# SIS-Errors: Errores sistemáticos clase<sub>1</sub>

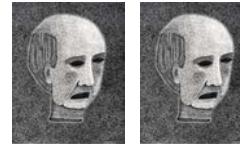
A[1]



A[2]



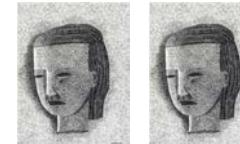
A[3]



A[4]

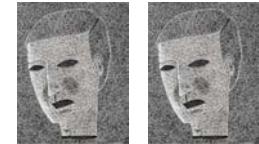


A[5]



...

A[n]



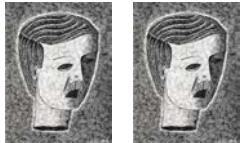
- **Conditional reliable token-RI:** is one that confers reliability [to the new context] if the methods, standards, and breadth of application are based on are also reliable [for the new context]

# SIS-Errors: Errores sistemáticos clase<sub>1</sub>

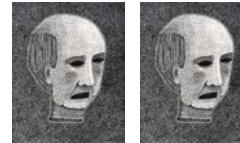
A[1]



A[2]



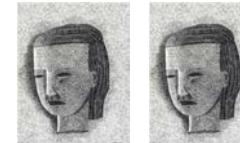
A[3]



A[4]

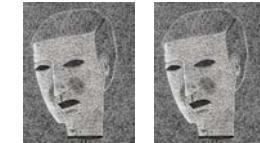


A[5]



...

A[n]



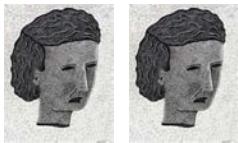
- **Conditional reliable token-RI:** is one that confers reliability [to the new context] if the methods, standards, and breadth of application are based on are also reliable [for the new context]

**Idea:** Realiza un seguimiento de la representación evaluando la probabilidad de que el token-RI mantenga la confiabilidad en todos los contextos.

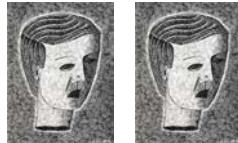
Nos permite mantener nuestras creencias en la fiabilidad del algoritmo en tanto y en cuanto las condiciones iniciales que dieron fiabilidad todavía se aplican en contextos nuevos

# SIS-Errors: Errores sistemáticos clase<sub>1</sub>

A[1]



A[2]



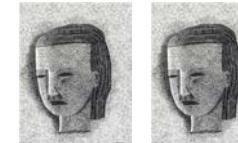
A[3]



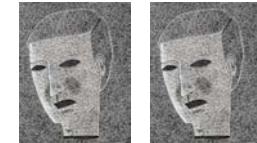
A[4]



A[5]



A[n]



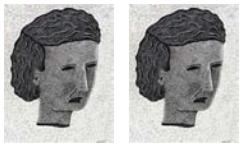
...

- **Conditional reliable token-RI:** is one that confers reliability [to the new context] if the methods, standards, and breadth of application are based on are also reliable [for the new context]

**Ejemplo (IF inadecuado - suspendo creencia):** un algoritmo que fue diseñado para detección de rostros aplicado en un nuevo contexto (e.g., durante el uso de COVID-19) necesita una actualización de la base de datos de entrenamiento, un nuevo módulo que trate con oclusión facial, etc

# SIS-Errors: Errores sistemáticos clase<sub>2</sub>

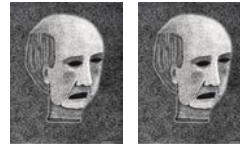
A[1]



A[2]



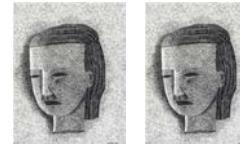
A[3]



A[4]

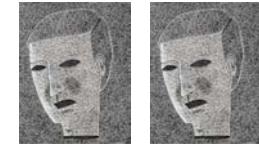


A[5]



...

A[n]



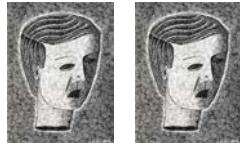
- **Anti-defeat clause:**  $S$  is epistemically warranted in maintaining token-RI as a suitable reliability indicator unless there is a defeater token-RI\* epistemically available to  $S$  such that, if  $S$  were to use toke-RI\*,  $S$  would no longer hold the belief that the output represents a fact  $F$

# SIS-Errors: Errores sistemáticos clase<sub>2</sub>

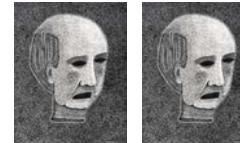
A[1]



A[2]



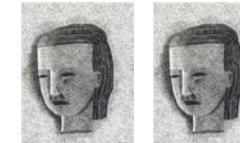
A[3]



A[4]



A[5]



...

A[n]



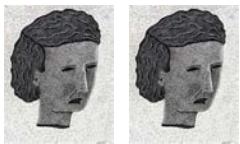
- **Anti-defeat clause:**  $S$  is epistemically warranted in maintaining token-RI as a suitable reliability indicator unless there is a defeater token-RI\* epistemically available to  $S$  such that, if  $S$  were to use token-RI\*,  $S$  would no longer hold the belief that the output represents a fact  $F$

**Idea:** El estatus justificado de una creencia se preserva mientras que la IF siga valiendo frente a alternativas nuevas y que potencialmente remueven credibilidad.

Mantenemos la fiabilidad en tanto y en cuanto no encontramos un defeater. En ese caso, estamos epistémicamente obligados a adoptarlo

# SIS-Errors: Errores sistemáticos clase<sub>2</sub>

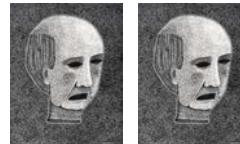
A[1]



A[2]



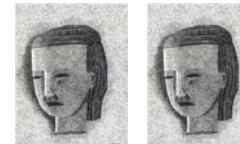
A[3]



A[4]

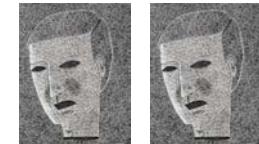


A[5]



...

A[n]



- **Anti-defeat clause:**  $S$  is epistemically warranted in maintaining token-RI as a suitable reliability indicator unless there is a defeater token-RI\* epistemically available to  $S$  such that, if  $S$  were to use token-RI\*,  $S$  would no longer hold the belief that the output represents a fact  $F$

**Ejemplo (IF incorrecto - reviso de la creencia):** Un algoritmo que selecciona {criminal, no criminal} basado en características faciales: token-IF = distancia entre los ojos, curvatura boca, etc.

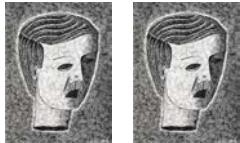
Defeater, más robusto: token-IF\* = implementación de conceptos socio-económicos de criminalidad

# SIS-Errors: Errores sistemáticos clase 3

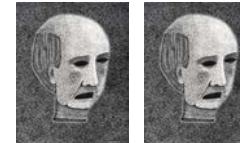
A[1]



A[2]



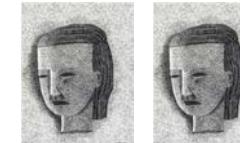
A[3]



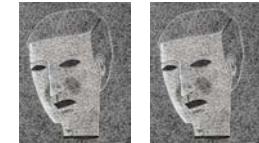
A[4]



A[5]



A[n]

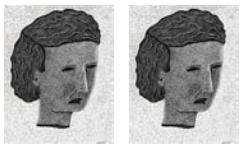


...

- **Supercharging type-RI<sub>3</sub> indicator:** Maintaining algorithmic reliability depends on subjecting their outputs to debate and other forms of scientific engagement. In this sense, the social construction of belief plays a crucial role in determining reliability and can, at times, take precedence over other indicators

# SIS-Errors: Errores sistemáticos clase 3

A[1]



A[2]



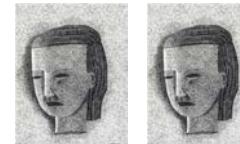
A[3]



A[4]

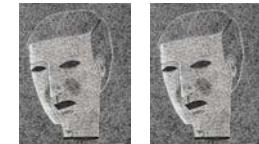


A[5]



A[n]

...



- **Supercharging type-RI<sub>3</sub> indicator:** Maintaining algorithmic reliability depends on subjecting their outputs to debate and other forms of scientific engagement. In this sense, the social construction of belief plays a crucial role in determining reliability and can, at times, take precedence over other indicators

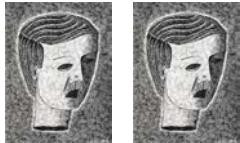
**Idea:** principio preventivo epistémico, donde la comunidad relevante evalúa los méritos de los resultados del algoritmo (tiene poder de veto!)

# SIS-Errors: Errores sistemáticos clase 3

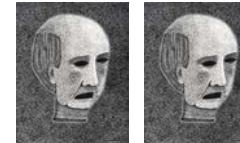
A[1]



A[2]



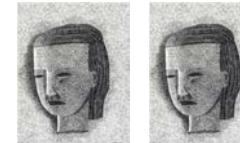
A[3]



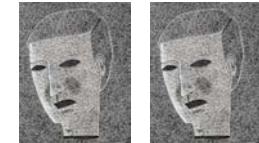
A[4]



A[5]



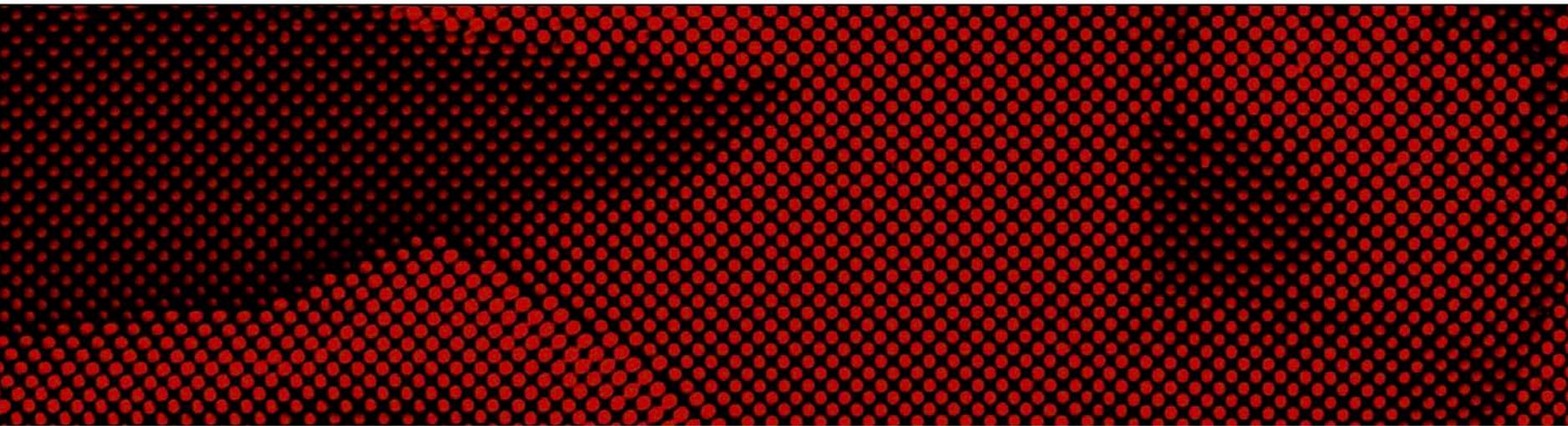
A[n]



...

- **Supercharging type-RI<sub>3</sub> indicator:** Maintaining algorithmic reliability depends on subjecting their outputs to debate and other forms of scientific engagement. In this sense, the social construction of belief plays a crucial role in determining reliability and can, at times, take precedence over other indicators

Ejemplo (ausencia de IF - anulo creencia): Kathryn!



VIDEO: SAM CANNON

MAIA SZALAVITZ

THE BIG STORY AUG 11, 2021 6:00 AM

# The Pain Was Unbearable. So Why Did Doctors Turn Her Away?

A sweeping drug addiction risk algorithm has become central to how the US handles the opioid crisis. It may only be making the crisis worse.

# Recapitulando

# Recapitulando

- Externalista:
  - La justificación no depende de un “third-party algorithms” y es independiente de nuestras capacidad por entender algoritmos
  - Cara indicador de fiabilidad está justificado independientemente
  - Hay una variedad de RI (e.g., es una epistemología “no-centralizada”)
  - En principio se pueden detectar indicadores de resultados no creíbles/falsos
  - Se acomoda a diferentes necesidades epistémicas/éticas: en algunos sistemas, los métodos de verificación y validación son más importantes, mientras que en otros es la coherencia teórica, mientras que en otros es el alineamiento con leyes vigentes

# Recapitulando

- SIS-Errors constituyen un problema para el fiabilista (para toda epistemología)
  - Condiciones para suspender nuestras creencias—y mantener fiabilidad?
  - Problemas de orden práctico
- Más problemas:
  - No es claro la precedencia, el orden, el peso, etc de cada indicador.
  - No es claro cómo se procede cuando hay conflicto entre los indicadores
  - *La tiranía de unos pocos indicadores:* unos pocos indicadores pueden ser suficiente para desbalancear la fiabilidad/no-fiabilidad de un algoritmo



Muchas gracias!

j.m.duran@tudelft.nl