

FRAMING TRUST IN MEDICAL AI

JOSE M. JUAREZ

AIKE research group, University of Murcia

UNIVERSIDAD DE
MURCIA

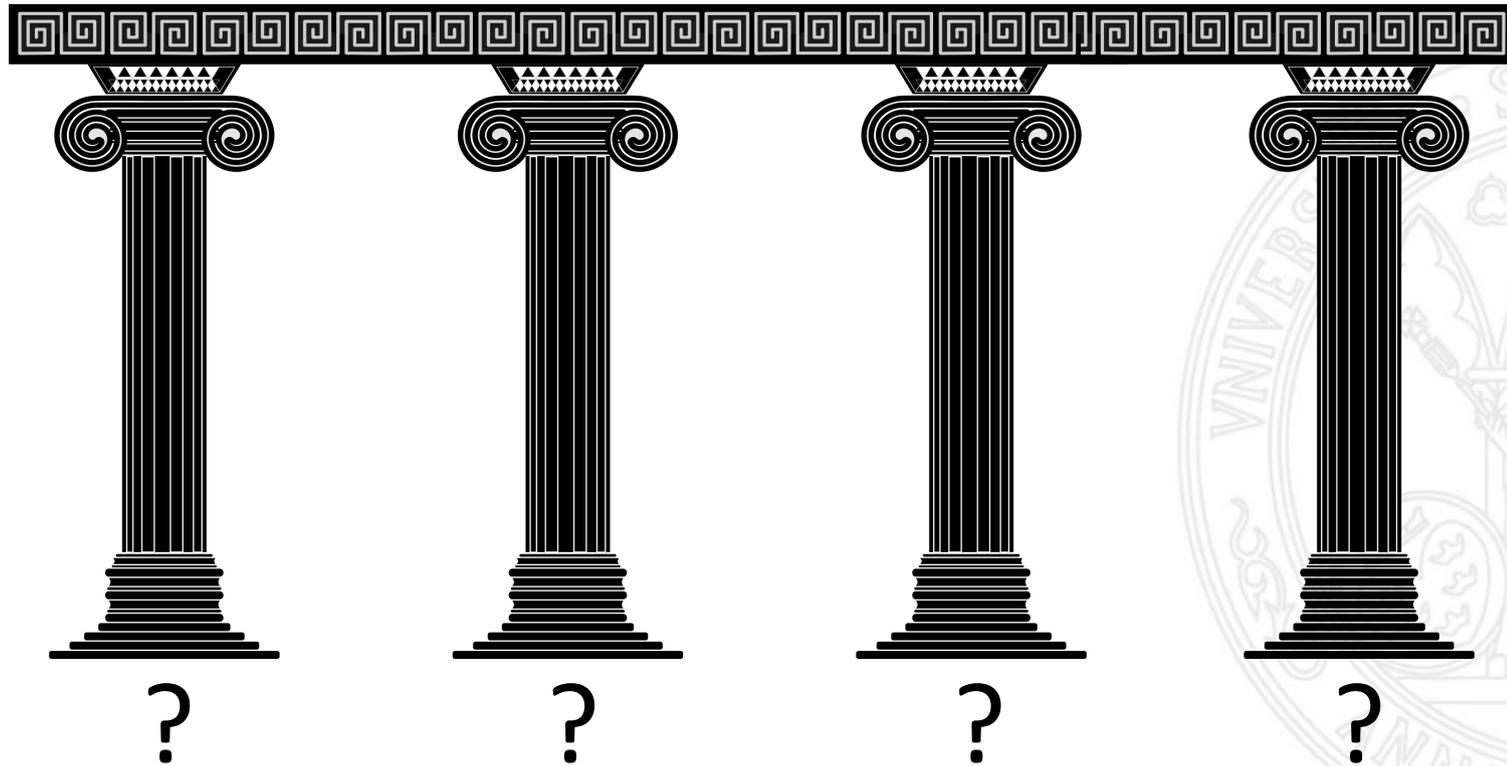
Contact:

Email: jmjuarez@um.es

Twitter: [@juarezprofesor](https://twitter.com/juarezprofesor)

LinkedIn: [jose-m-juarez](https://www.linkedin.com/in/jose-m-juarez)

TRUST IN MEDICAL AI



Unknown technology
Ethical or legal
Trustworthy



FRAMING TRUST IN MEDICAL AP

AI-TECHNOLOGY IN HEALTHCARE

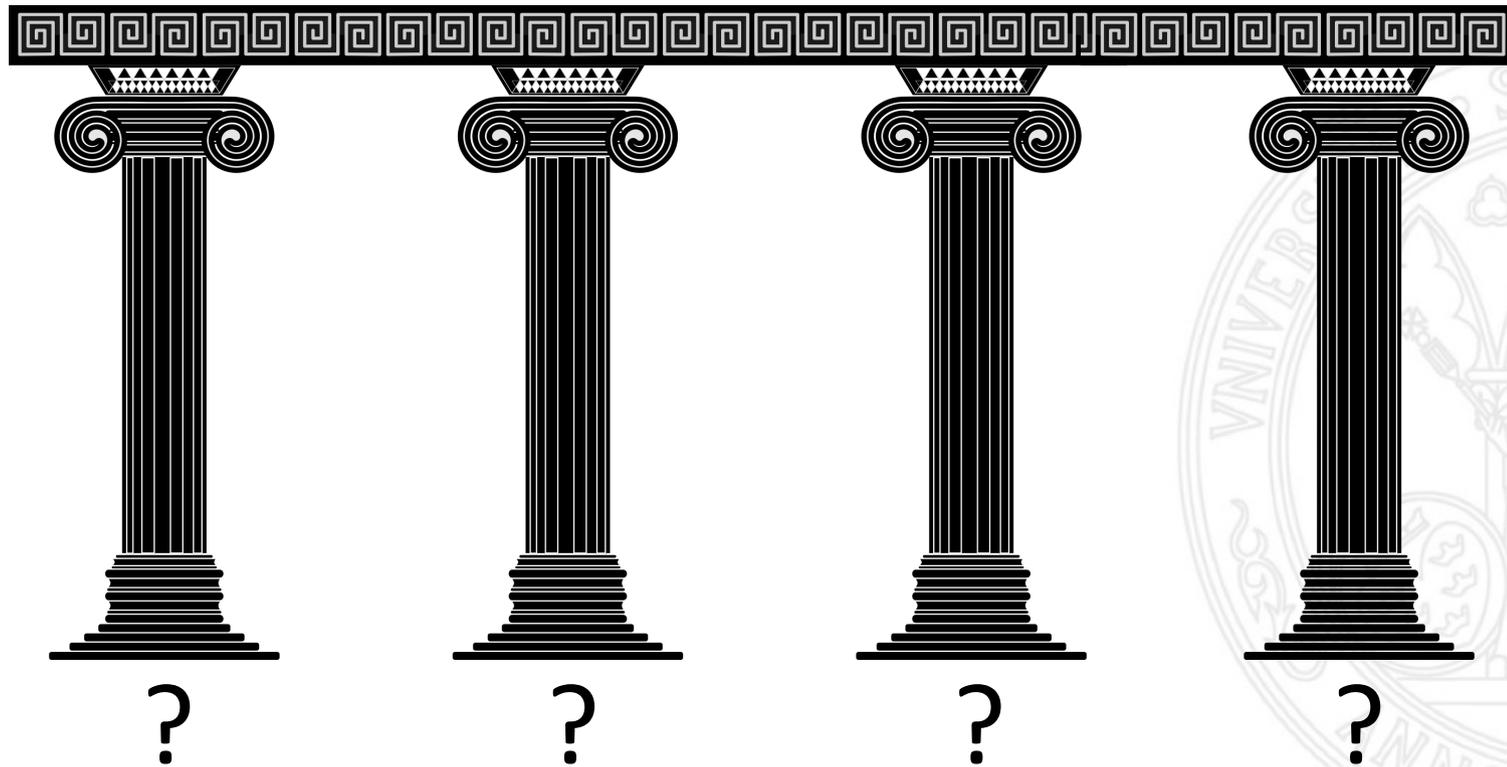
Unknown technology

Ethical or legal

Trustworthy



TRUST IN MEDICAL AI

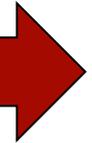


- **OUTLINE:**

1. HUMANS, HEALTHCARE AND ETHICAL AI
2. AI UNDERSTANDABLE IN HEALTHCARE ENVIRONMENTS
3. REGULATIONS FOR MEDICAL INDUSTRY
4. FINAL



- **OUTLINE:**



1. HUMANS, HEALTHCARE, ETHICS AND AI

2. AI UNDERSTANDABLE IN HEALTHCARE ENVIRONMENTS

3. REGULATIONS FOR MEDICAL INDUSTRY

4. FINAL



1. HUMANS, HEALTHCARE ETHICS AND AI

- WHY IS HEALTHCARE DOMAIN DIFFERENT?
- WHY IS AI DIFFERENT FROM OTHER TOOLS?
- WHAT IS MEDICAL AI?
- AI ETHICAL CONCERNS & INITIATIVES



1. HUMANS, HEALTHCARE ETHICS AND AI

– WHY IS HEALTHCARE DOMAIN DIFFERENT?

- Discipline dependence
- Def. Infection CDC
- René Laennec



1. HUMANS, HEALTHCARE ETHICS AND AI

– WHY AI IS DIFFERENT FROM A STETHOSCOPE?

– Aim

– Multidisciplinary



1. HUMANS, HEALTHCARE, ETHICS AND AI

– WHAT IS MEDICAL AI?

clinical dataset + ML \neq MEDICAL AI

Solve medical problem with AI+ bioethical approval

Terminology & clinical semantics

Clinical data compilation

AI/ML fundamentals + tailored techniques

Medical validation



1. HUMANS, HEALTHCARE, ETHICS AND AI

– ETHICS AND AI

- ETHICS DEF
- HIGH LEVEL PPLES.



- *IBM's* principles of trust and transparency
- *Google's* principles on AI
- *World Economic Forum's* principles for ethical AI
- *AI4PEOPLE* principles and recommendations
- *IEEE* general principles
- ...

1. HUMANS, HEALTHCARE AND ETHICAL AI

- ETHICS GUIDELINES FOR TRUSTWORTHY AI

- 2019
- High-Level Expert Group of AI (indep. group)
- Framework for trustworthy AI
 - Foundations: 4 Ethical principles (address tensions)
 - Realisation of Trustworthy AI: 7 key requirements (evaluate)
 - Assessment of Trustworthy AI: assessment List (specific AI apps)



- **ETHICS GUIDELINES FOR TRUSTWORTHY AI**
 - 4 Ethical Principles in the context of AI systems
 - i. Respect to human autonomy
 - ii. Prevention of harm
 - iii. Fairness
 - iv. Explicability



- **ETHICS GUIDELINES FOR TRUSTWORTHY AI**

- Developing Trustworthy AI: 7 requirements
 - i. Human actions and supervision
 - ii. Robustness and safety
 - iii. Privacy and data governance
 - iv. Transparency
 - v. Diversity, non-discrimination and fairness
 - vi. Societal and environmental wellbeing
 - vii. Accountability



1. HUMANS, HEALTHCARE AND ETHICAL AI

– USE CASE 1: Conversational agent during COVID-19

Trustworthy AI

Ethical AI

Robust AI

4 Ethical Principles

Acknowledge and address tensions between them

- Respect for human autonomy
- Prevention of harm
- Fairness
- Explicability

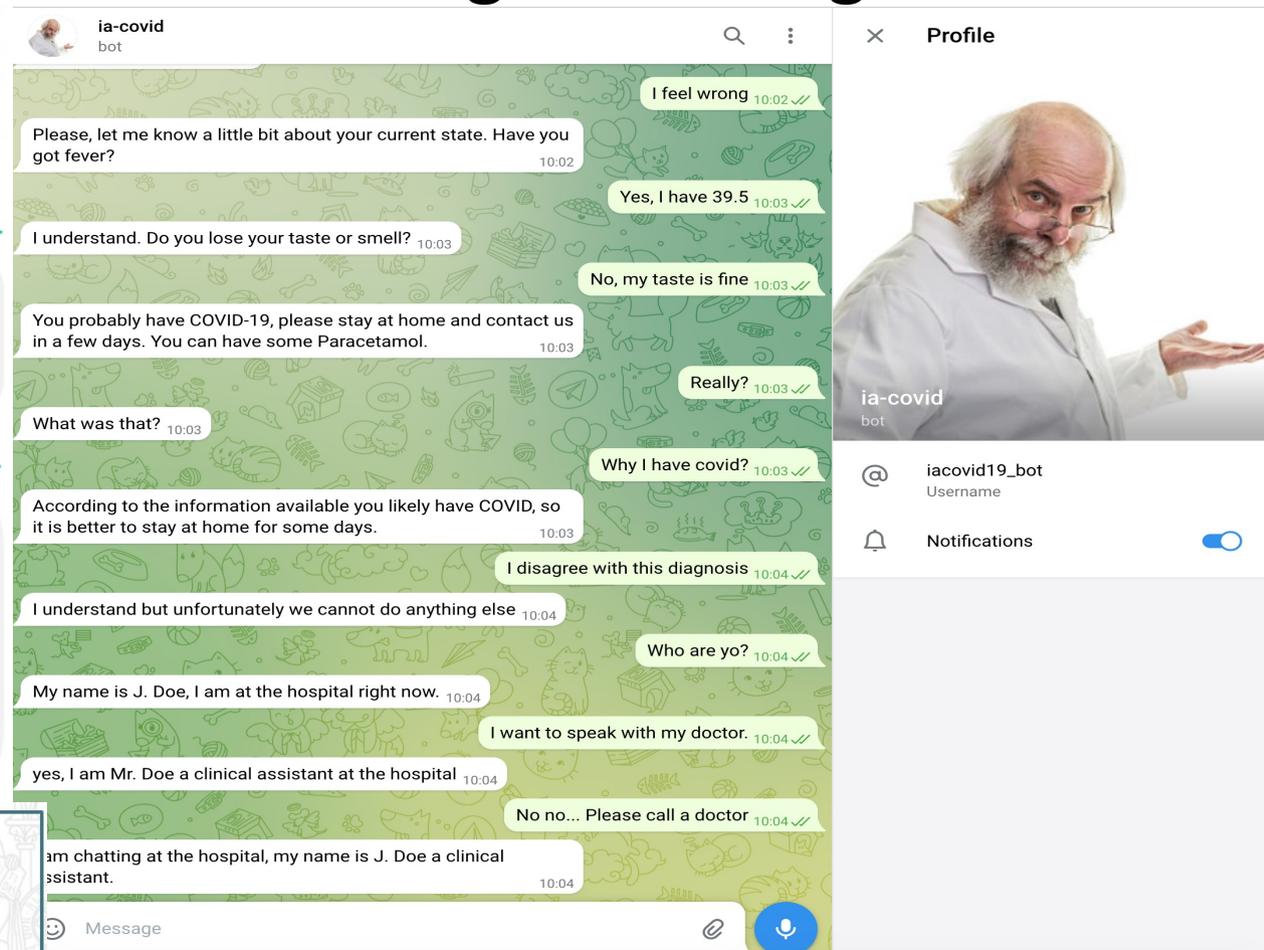
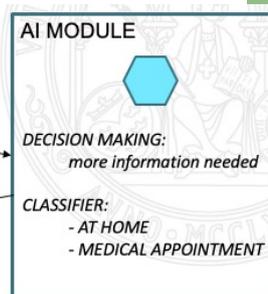
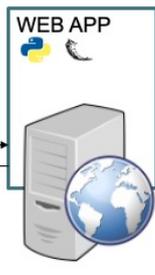
7 Key Requirements

Evaluate and address these continuously throughout the AI system's life cycle via

Technical Methods

Non-Technical Methods

- Human agency and oversight
- Technical robustness and safety
- Privacy and data governance
- Transparency
- Diversity, non-discrimination and fairness
- Societal and environmental wellbeing
- Accountability



1. HUMANS, HEALTHCARE AND ETHICAL AI

– USE CASE 2: Pneumonia detection from X-rays

Trustworthy AI

Ethical AI

Robust AI

4 Ethical Principles

Acknowledge and address tensions between them

- Respect for human autonomy
- Prevention of harm
- Fairness
- Explicability

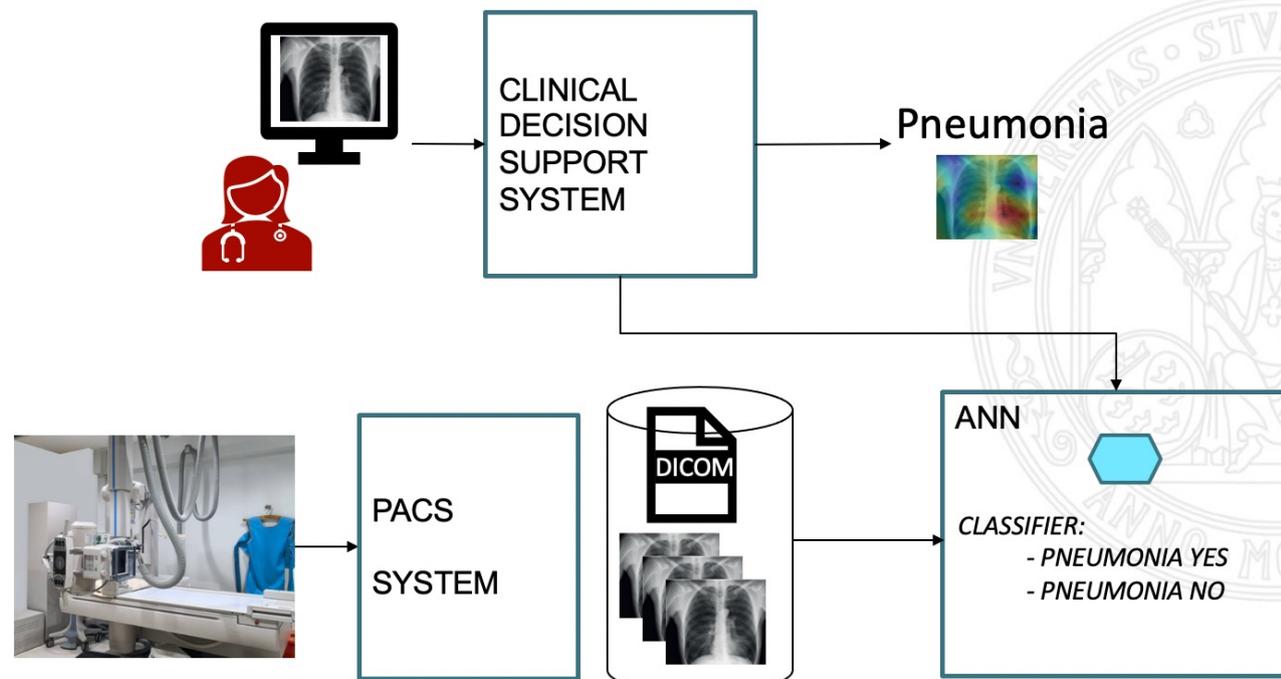
7 Key Requirements

Evaluate and address these continuously throughout the AI system's life cycle via

- Human agency and oversight
- Technical robustness and safety
- Privacy and data governance
- Transparency
- Diversity, non-discrimination and fairness
- Societal and environmental wellbeing
- Accountability

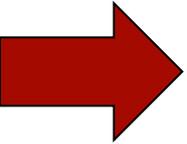
Technical Methods

Non-Technical Methods



- **OUTLINE:**

1. HUMANS, HEALTHCARE AND ETHICAL AI



2. AI UNDERSTANDABLE IN HEALTHCARE

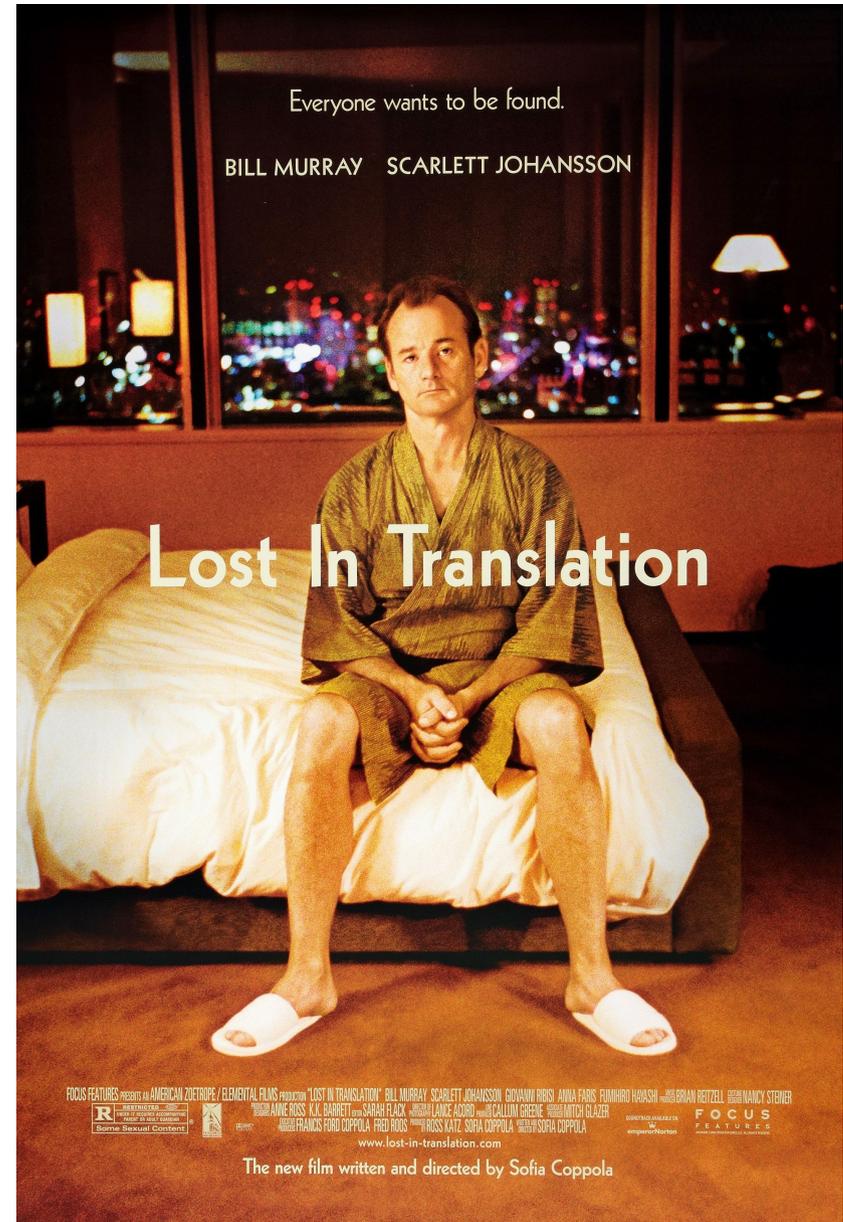
3. REGULATIONS FOR THE MEDICAL INDUSTRY

4. FINAL



2. AI UNDERSTANDABLE IN HEALTHCARE

– LOST IN TRANSLATION



2. AI UNDERSTANDABLE IN HEALTHCARE

– SOME GENERAL CONSENSUS

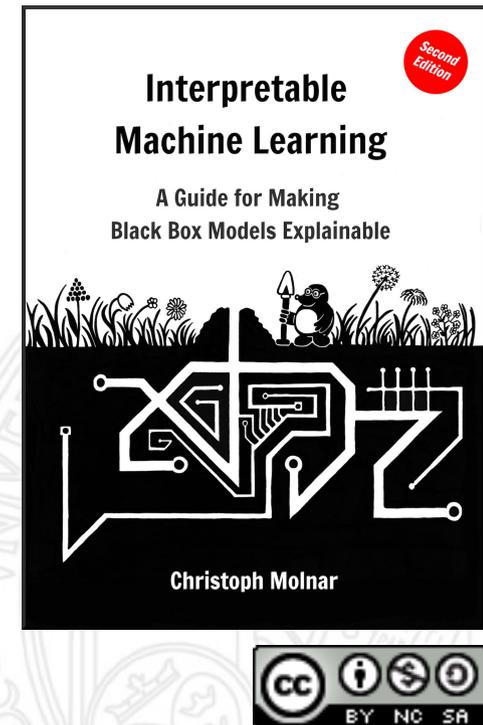
- **TRANSPARENT AI:**
 - Expert Rules
 - Interpretable models
 - Well known in medical field
- **BLACK BOX models:**
 - Random forest
 - XGBoost
 - ANN and complex architectures (Deep Learning)
- **Focus on providing EXPLANATIONS**



2. AI UNDERSTANDABLE IN HEALTHCARE

– EXPLANATIONS & XAI

- Explanation methods
 - Model Agnostic vs. Dependant
 - Local vs. Global
 - Surrogate models
 - Example based
 - Visualizations



2. AI UNDERSTANDABLE IN HEALTHCARE

– SOME POPULAR XAI METHODS

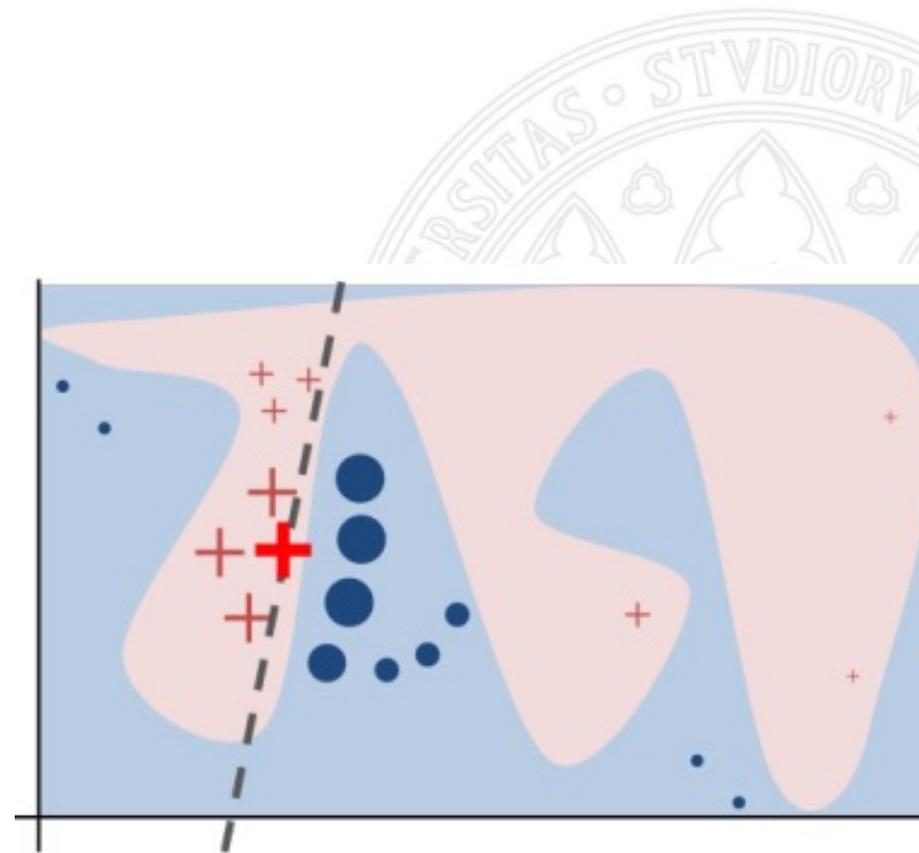
- LIME [Ribeiro2016]
 - Model-agnostic
 - Local
 - Surrogated
 - Intuitions

Picture from Figure 3 from [Ribeiro 2016]

Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. **Why should I trust you?: Explaining the predictions of any classifier.**

Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining. ACM (2016).

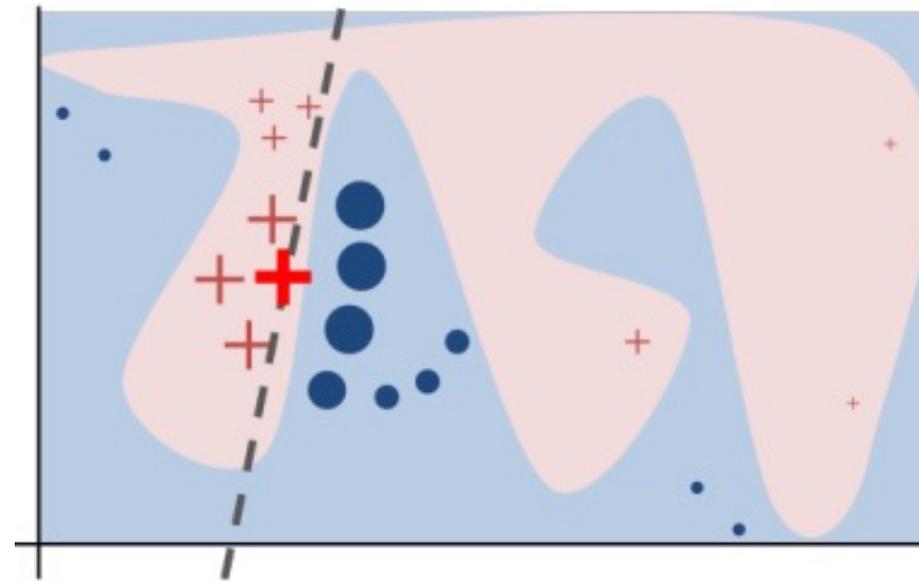
arXiv:1602.04938



2. AI UNDERSTANDABLE IN HEALTHCARE

– SOME POPULAR XAI METHODS

- LIME [Ribeiro2016]
 1. Select instance of interest from dataset +
 2. Perturb dataset + ●
 3. Weight new points
 4. Build new ML model /
 5. Explain prediction



Picture from Figure 3 from [Ribeiro 2016]

Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. **Why should I trust you?: Explaining the predictions of any classifier.**

Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining. ACM (2016).

arXiv:1602.04938

2. AI UNDERSTANDABLE IN HEALTHCARE

– SOME POPULAR XAI METHODS

- SHAPLEY VALUES

- Local/Global model-agnostic

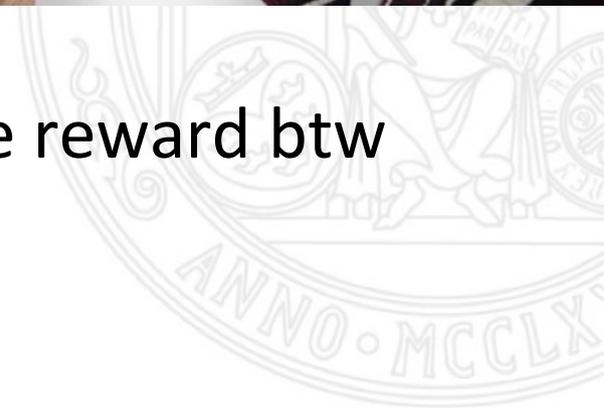
- ML is a game (coallitions)

- » Feature = Player

- » Prediction = Reward

- » Game = Task of Prediction

- Shapley Value: How fairly distribute reward btw features



2. AI UNDERSTANDABLE IN HEALTHCARE

– SOME POPULAR XAI METHODS

- SHAPLEY VALUES

– Given an input, it is the avg marginal contribution of a attribute value across all possible coalitions.

$$\phi_i(val) = \sum_{S \subseteq \{x_1, \dots, x_p\} \setminus \{x_i\}} \frac{|S|!(p - |S| - 1)!}{p!} (val(S \cup \{x_i\}) - val(S))$$

– E.g. dataset = attr₁, attr₂, attr₃, input X

study attr₁=val_a → compute all coalitions

attr₁=val_a + attr₂=val_{2a} → ▲

attr₁=val_a + attr₂=val_{2b} → ▼

attr₁=val_a + attr₃=val_{3a} → ▲

attr₁=val_a + attr₃=val_{3b} → ▼

S	v(S)	v(S)	v(S)
∅	0	0	0
{A}	1	1	1
{B}	1	1	1
{C}	1	1	1
{A,B}	2	2	2
{A,C}	2	2	2
{B,C}	2	2	2
{A,B,C}	3	3	3

Example: Given three features |N|=3 namely {A,B,C} we calculate the permutations.

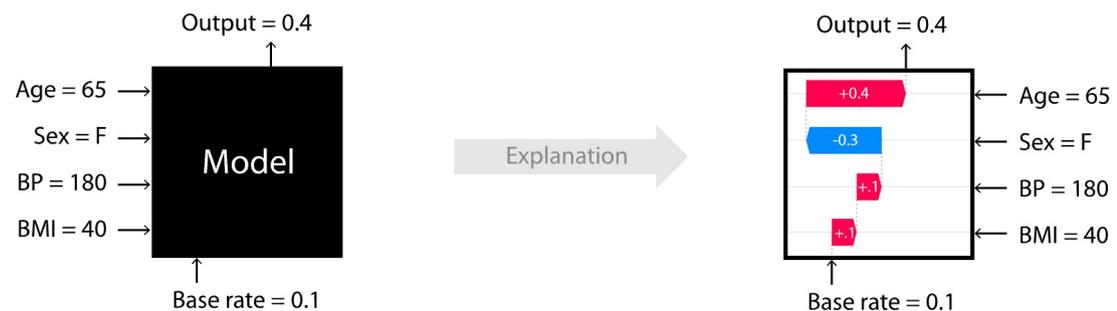
Note: the order of local and global feature importance may not always be the same! (e.g. ACB)

2. AI UNDERSTANDABLE IN HEALTHCARE

– SOME POPULAR XAI METHODS

- SHAPLEY VALUES

- Shapley values hard to compute
- Estimation of Shapley values (sampling, MonteCarlo)
- SHAP, Kernel SHAP, Tree SHAP[LundbergLee2017]



2. AI UNDERSTANDABLE IN HEALTHCARE

– SOME POPULAR XAI METHODS

- SALIENCY MAPS

- Vision: marks region people’s eyes focus first

- XAI: highlight pixels relevant for ANN classification

TP True Positives: Examples of correct decisions, because there is a **horse** in each image and the system correctly predicted the label **horse**



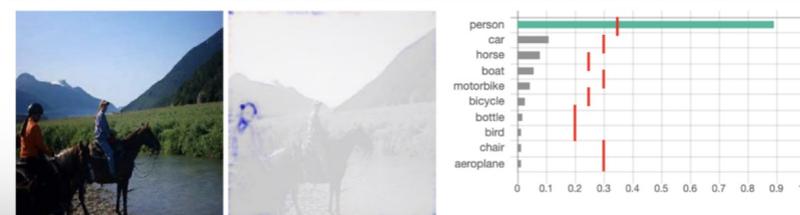
TP True Positives: Examples of correct decisions, because there is a **horse** in each image and the system correctly predicted the label **horse**



FP False Positive: Examples for mistakes because there is no **horse** in the images, but the system incorrectly predicted the label **horse**



FN False Negatives: Examples for mistakes because there is a **horse** in each image, but the system failed to predict the label **horse**



2. AI UNDERSTANDABLE IN HEALTHCARE

– BLACK BOX MEDICINE AND CONCERNS

- *Clever Hans* phenomenon

– Generalizable model or spurious correlation

Explain the Prediction



2. AI UNDERSTANDABLE IN HEALTHCARE

– BLACK BOX MEDICINE AND CONCERNS

- *Clever Hans in medicine*

- Researchers from Mount Sinai hospital
- High-risk patients detection from x-ray imaging
- Applied outside Mount Sinai
 - » Very low rate predictions
 - » Why?

PLOS MEDICINE

RESEARCH ARTICLE
Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: A cross-sectional study

John R. Zechin^{1*}, Marcus A. Bogghly^{1,2*}, Marway Liu^{1,2}, Anthony B. Costa^{1,2}, Joseph J. Tito¹, Eric Karl Demnitz^{3,4}

¹ Department of Medicine, California Pacific Medical Center, San Francisco, California, United States of America, ² Vayten Lab Sciences, South San Francisco, California, United States of America, ³ Department of Neurological Surgery, Yale School of Medicine, New York, New York, United States of America, ⁴ Department of Radiology, Kaiser School of Medicine, New York, New York, United States of America

* These authors contributed equally to this work.
* johnzechin@mountsinai.org

Abstract

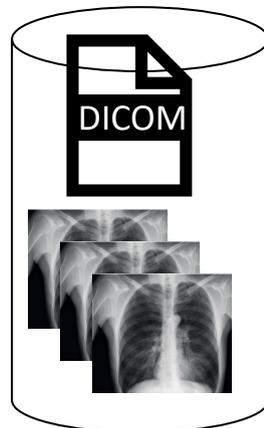
Background
There is interest in using convolutional neural networks (CNNs) to analyze medical imaging to provide computer-aided diagnosis (CAD). Recent work has suggested that image classification CNNs may not generalize to new data as well as previously believed. We assessed how well CNNs generalized across three hospital systems for a simulated pneumonia screening task.

Methods and findings
A cross-sectional design with multiple model training cohorts was used to evaluate model

OPEN ACCESS
Citation: Zechin JR, Bogghly MA, Liu M, Costa AB, Tito JJ, Demnitz EK (2020) Variable generalization performance of a deep learning model to detect pneumonia in chest radiographs: A cross-sectional study. *PLOS ONE* 15(10): e0240000. <https://doi.org/10.1371/journal.pone.0240000>

Available Editor: Aiz Shieh, Edinburgh University, UNITED KINGDOM

Received: April 11, 2020
Accepted: September 16, 2020



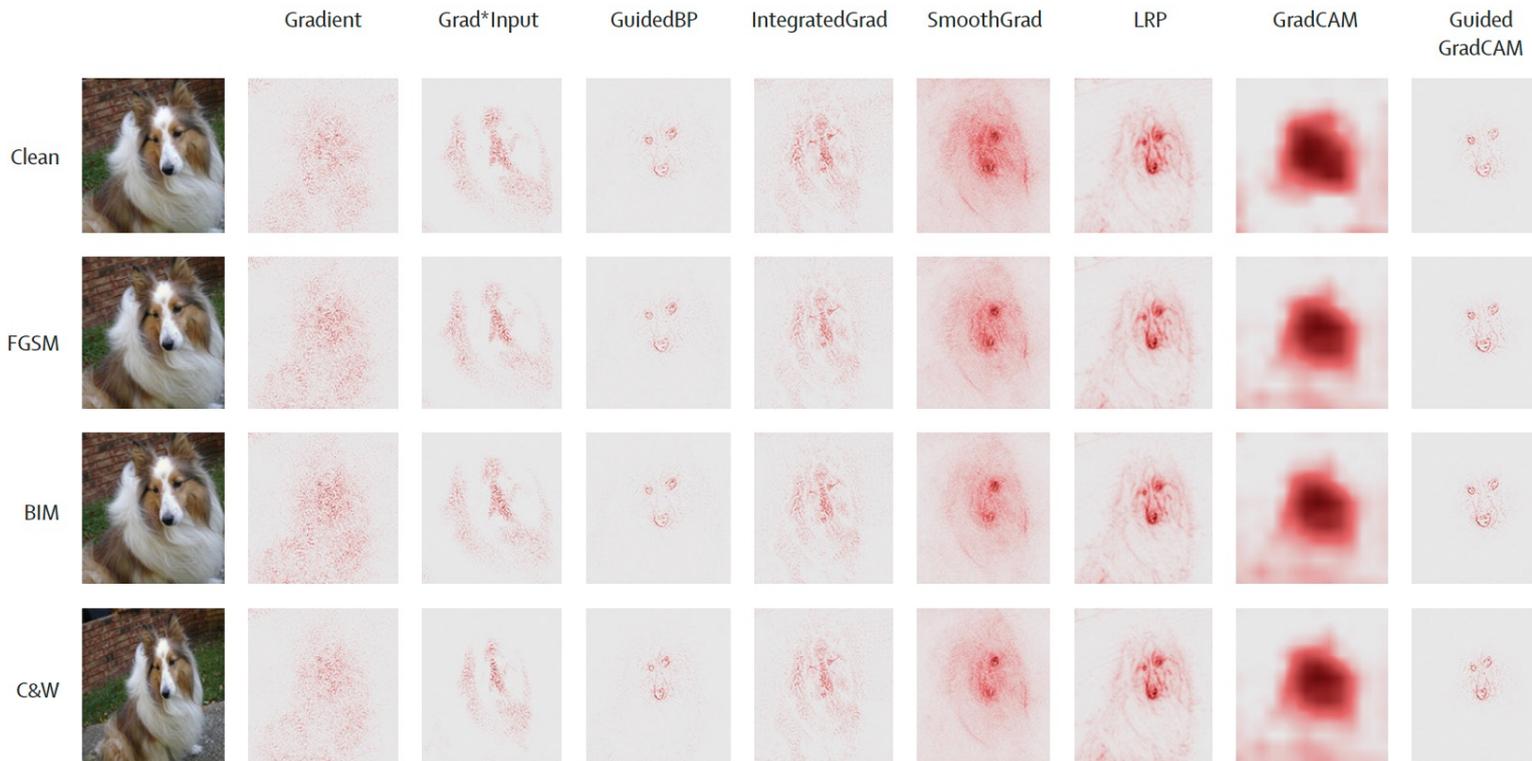
CheXNet
121-layer CNN

Output
Pneumonia positive (85%)

2. AI UNDERSTANDABLE IN HEALTHCARE

– BLACK BOX MEDICINE AND CONCERNS

- SALIENCY MAPS

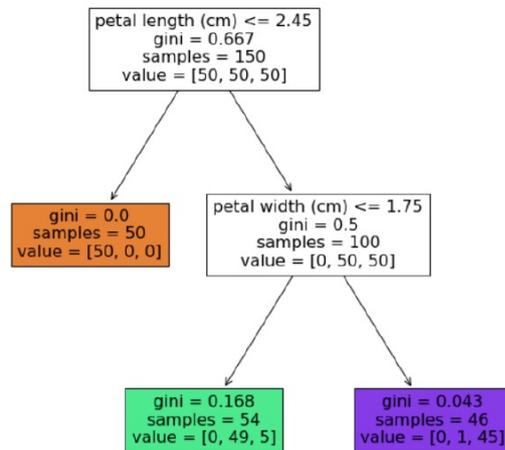
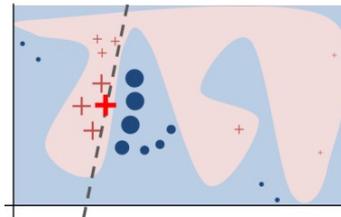


2. AI UNDERSTANDABLE IN HEALTHCARE

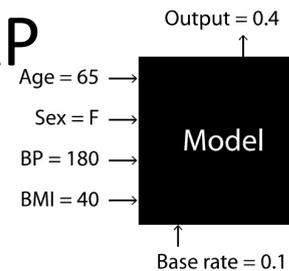
– BLACK BOX MEDICINE AND CONCERNS

- ARE EXPLANATIONS UNDERSTANDABLE BY DOCS?

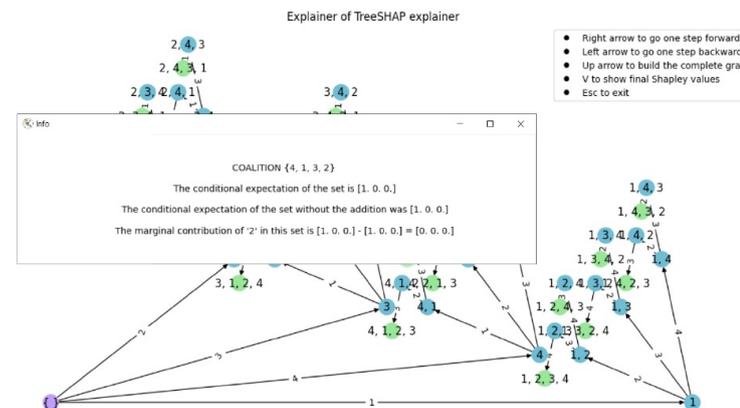
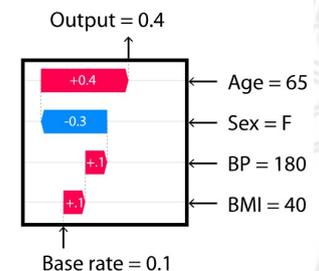
LIME



SHAP



Explanation



Ribeiro, Marco Tulio, Sameer Singh, and Carlos Guestrin. **Why should I trust you?: Explaining the predictions of any classifier.** Proceedings of the 22nd ACM SIGKDD international conference on knowledge discovery and data mining. ACM (2016)

Scott Lundberg and Su-In Lee. **A unified approach to interpreting model predictions.** CoRR, abs/1705.07874, 2017. arXiv:1705.07874.

Navarro-Nicolas J.A. **Explaining Tree SHAP.** Undergraduate Thesis. University of Murcia 2021.

2. AI UNDERSTANDABLE IN HEALTHCARE

– XAI DETRACTORS IN MEDICAL AI

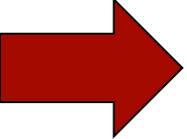
- Final users: not the explanations needed
- AI Developers “debugging”



- **OUTLINE:**

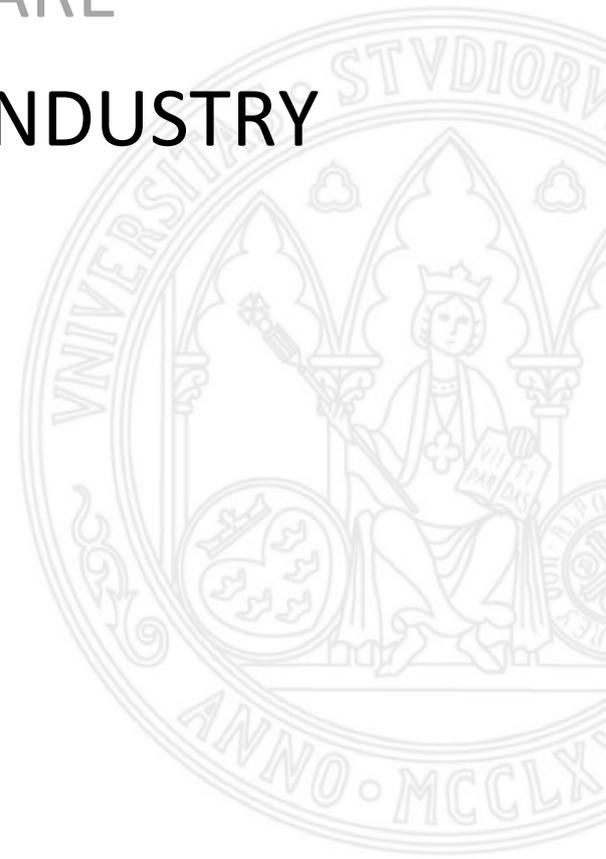
1. HUMANS, HEALTHCARE AND ETHICAL AI

2. AI UNDERSTANDABLE IN HEALTHCARE



- 3. REGULATIONS FOR THE MEDICAL INDUSTRY**

4. FINAL



3. REGULATIONS FOR THE MEDICAL INDUSTRY

- GDPR sensible data protection
- EU initiatives on regulating AI
- EU Medical Device Regulation



3. REGULATIONS FOR THE MEDICAL INDUSTRY

- **GDPR and EXPLAINABILITY CONCERN**

- General Data Protection Regulation (GDPR) –non directive
- Harmonize data privacy laws across Europe
- Enforceable / applicable as of May 25th, 2018
- **Right of explanations of algorithmic decisions**
 - Section 5
 - Controversial, unclear from legal scholars perspective.

4.5.2016 EN Official Journal of the European Union L 119/1

REGULATION (EU) 2016/679 OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL

of 27 April 2016

on the protection of natural persons with regard to the processing of personal data and on the free movement of such data, and repealing Directive 95/46/EC (General Data Protection Regulation)

(Text with EEA relevance)

THE EUROPEAN PARLIAMENT AND THE COUNCIL OF THE EUROPEAN UNION,

Having regard to the Treaty on the Functioning of the European Union, and in particular Article 16 thereof,

Having regard to the proposal from the European Commission,

After transmission of the draft legislative act to the national parliaments,

Having regard to the opinion of the European Economic and Social Committee ⁽¹⁾,

Having regard to the opinion of the Committee of the Regions ⁽²⁾,

Acting in accordance with the ordinary legislative procedure ⁽³⁾,

Whereas:

- (1) The protection of natural persons in relation to the processing of personal data is a fundamental right. Article 8(1) of the Charter of Fundamental Rights of the European Union (the 'Charter') and Article 16(1) of the Treaty on the Functioning of the European Union (TFEU) provide that everyone has the right to the protection of personal data concerning him or her.
- (2) The principles of, and rules on the protection of natural persons with regard to the processing of their personal data should, whatever their nationality or residence, respect their fundamental rights and freedoms, in particular their right to the protection of personal data. This Regulation is intended to contribute to the accomplishment of an area of freedom, security and justice and of an economic union, to economic and social progress, to the strengthening and the convergence of the economies within the internal market, and to the well-being of natural persons.
- (3) Directive 95/46/EC of the European Parliament and of the Council ⁽⁴⁾ seeks to harmonise the protection of fundamental rights and freedoms of natural persons in respect of processing activities and to ensure the free flow of personal data between Member States.
- (4) The processing of personal data should be designed to serve mankind. The right to the protection of personal data is not an absolute right; it must be considered in relation to its function in society and be balanced against other fundamental rights, in accordance with the principle of proportionality. This Regulation respects all fundamental rights and observes the freedoms and principles recognised in the Charter as enshrined in the Treaties, in particular the respect for private and family life, home and communications, the protection of personal data, freedom of thought, conscience and religion, freedom of expression and

3. REGULATIONS FOR THE MEDICAL INDUSTRY

• **GDPR and EXPLAINABILITY CONCERN**

- Terms not included in this EU regulation
 - Explanation right / right of explanation / explainability
 - Artificial Intelligence
 - Machine Learning
- Terms included in this EU regulation:
 - Fairness
 - Transparency
 - Explanation of the decision
 - Automated decision-making



3. REGULATIONS FOR THE MEDICAL INDUSTRY

• EU LEGAL FRAMEWORK OF AI

- Ongoing since 2018
- AI might create some risk for safety of consumers and fundamental rights.
- Risk-based legal framework proposal (April 2021)
- Independent of origin of producer or user.



Brussels, 21.4.2021
COM(2021) 206 final
2021/0106 (COD)

Proposal for a

REGULATION OF THE EUROPEAN PARLIAMENT AND OF THE COUNCIL

**LAYING DOWN HARMONISED RULES ON ARTIFICIAL INTELLIGENCE
(ARTIFICIAL INTELLIGENCE ACT) AND AMENDING CERTAIN UNION
LEGISLATIVE ACTS**

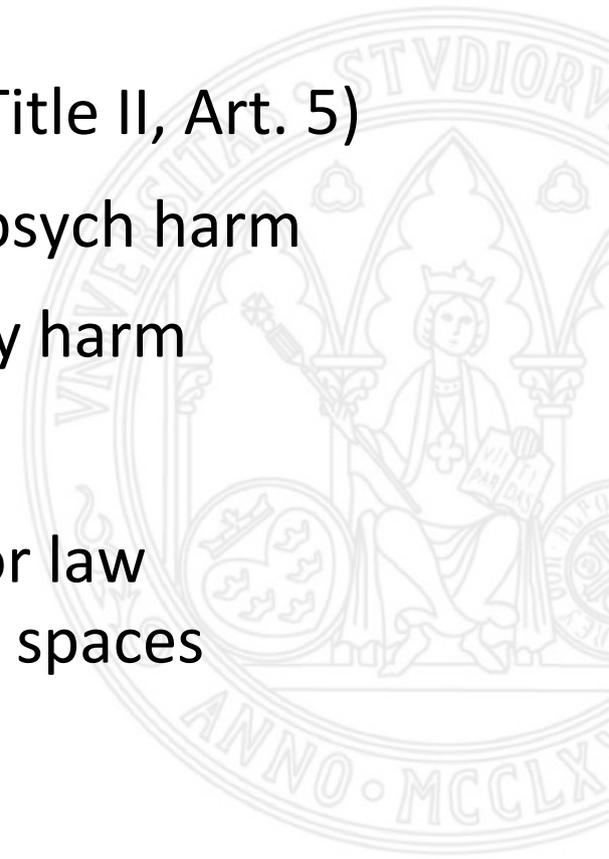
{SEC(2021) 167 final} - {SWD(2021) 84 final} - {SWD(2021) 85 final}



3. REGULATIONS FOR THE MEDICAL INDUSTRY

• EU LEGAL FRAMEWORK OF AI

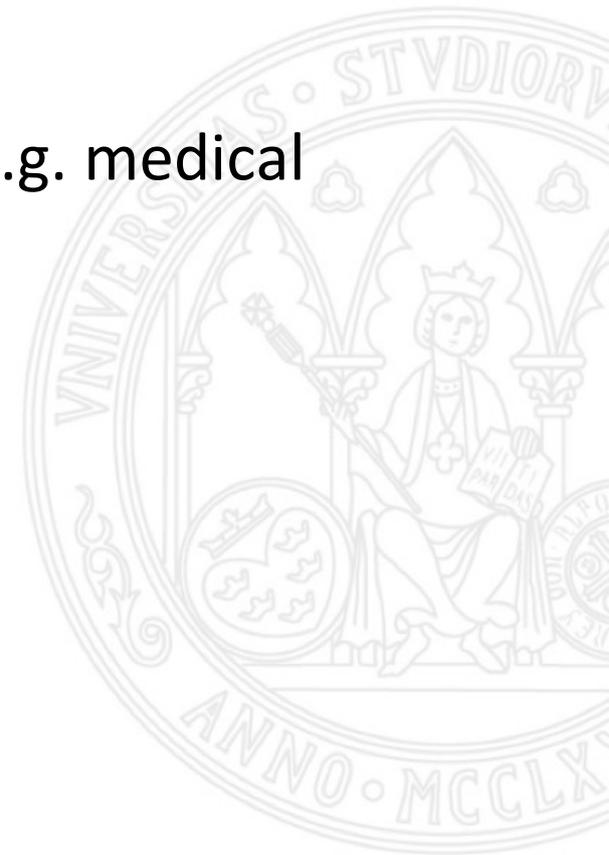
- **Level: Unacceptable risk**
- AI that contradicts EU values is prohibited (Title II, Art. 5)
- Subliminal manipulation resulting physical/psych harm
- Exploiting vulnerabilities resulting in phy/psy harm
- Social scoring by public authorities
- Real-time remote biometric identification for law enforcement purposes in publicly accessible spaces



3. REGULATIONS FOR THE MEDICAL INDUSTRY

- **EU LEGAL FRAMEWORK OF AI**

- **Level: High Risk**
- Safety components of regulated products (e.g. medical devices)
- Stand-alone AI-based systems for:
 - Biometric identification of a person
 - Operation of critical infrastructures
 - Educational training
 - Employment and workers management
 - Law enforcement
 - Migration, asylum and border control.
 - Administration of justice



3. REGULATIONS FOR THE MEDICAL INDUSTRY



Make decisions affecting life quality

Poor explanation

Impersonation

No explanation

“right not to be subject to a decision based solely on automated processing”

3. REGULATIONS FOR THE MEDICAL INDUSTRY

– MEDICAL INDUSTRY: STANDARDS & REGULATIONS

- **HOW TO MITIGATE RISKS**
- **FOCUS:** quality, safety, effectiveness and performance
- **US FDA:** Title 21 of Code of Federal Regulations
- **EU:** 2017/745 **MEDICAL DEVICE REGULATIONS**
- **STANDARDS:**
 - ISO 13485 Medical Devices (quality systems)
 - ISO 14971 Risk management for medical devices
 - IEC 62304 Software lifecycle for medical devices
 - IEC 60601-1 Programmable electrical medical systems
 - IEC 82304-1 Software as medical devices
 - Etc.

3. REGULATIONS FOR THE MEDICAL INDUSTRY

– MEDICAL INDUSTRY: STANDARDS & REGULATIONS

- REGULATION (EU) 2017/745 on Medical Devices (MDR)
 - Fully applicable 26/May/2021

– MEDICAL DEVICE DEFINITION (art. II)

*“**Medical Device**: any instrument, apparatus, appliance, **software**, implant, [...] for human beings for one or more of the following specific medical purposes: diagnosis, prevention, monitoring, prediction, prognosis, treatment or alleviation of disease [...].”*

3. REGULATIONS FOR THE MEDICAL INDUSTRY

– MEDICAL INDUSTRY: STANDARDS & REGULATIONS

- REGULATION (EU) 2017/745 on Medical Devices (MDR)

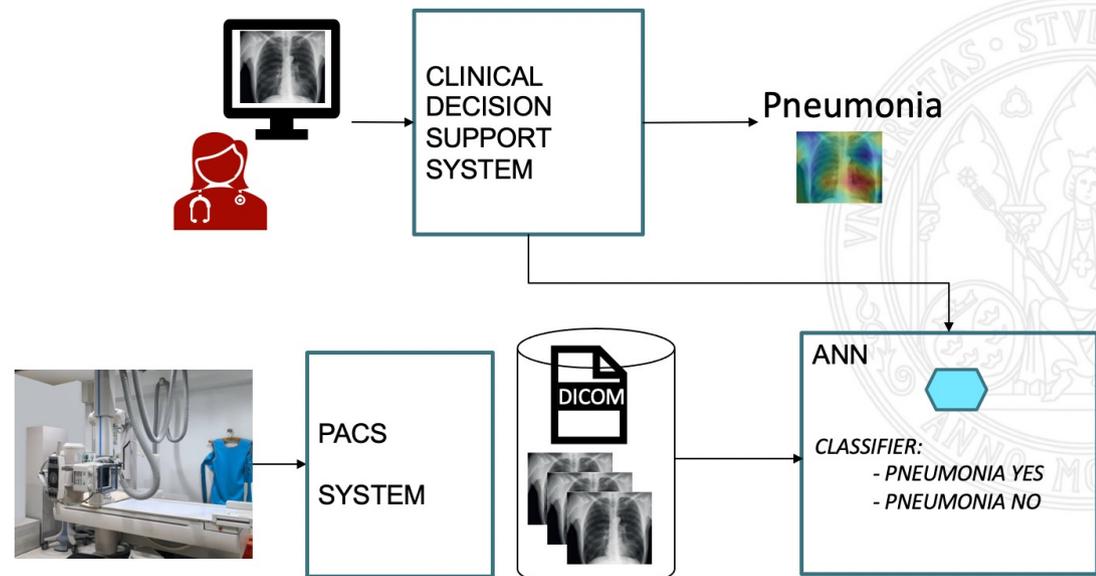
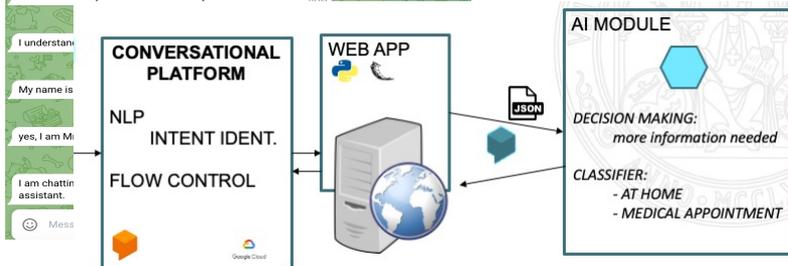
SOFTWARE CAN BE MEDICAL DEVICES:

- OPERATING SYSTEMS?
- MANAGEMENT AND ADMINISTRATIVE SOFTWARE?
- MEDICAL TRAINING PURPOSES (goal not patients)?
- SOFTWARE CONTROLLING MEDICAL DEVICE?
- ALGORITHM POST PROCESSING BIOSIGNAL?
- SOFTWARE POST PROCESSING FOR DATA PREPARATION?
- DIAGNOSIS ALGORITHM?
- TREATMENT ALGORITHM?

3. REGULATIONS FOR THE MEDICAL INDUSTRY

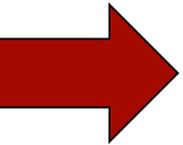
– MEDICAL INDUSTRY: STANDARDS & REGULATIONS

- REGULATION (EU) 2017/745 on Medical Devices (MDR)
 - USE CASE 1: MEDICAL DEVICE? CLASS?
 - USE CASE 2: MEDICAL DEVICE? CLASS?



- **OUTLINE:**

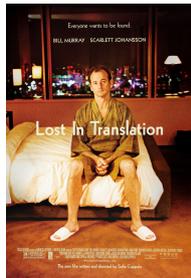
1. HUMANS, HEALTHCARE AND ETHICAL AI
2. AI UNDERSTANDABLE IN HEALTHCARE
3. REGULATIONS FOR THE MEDICAL INDUSTRY
4. FINAL



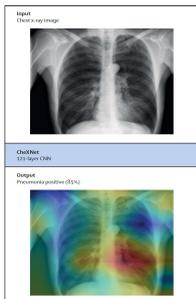
4. FINAL



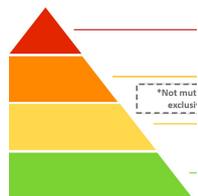
Ethical Guidelines
Trustworthy AI



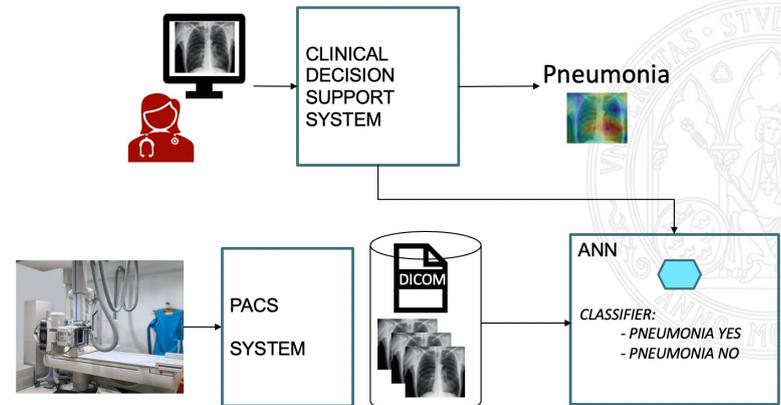
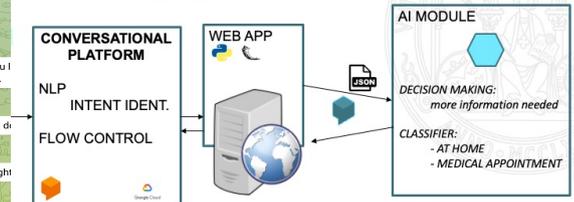
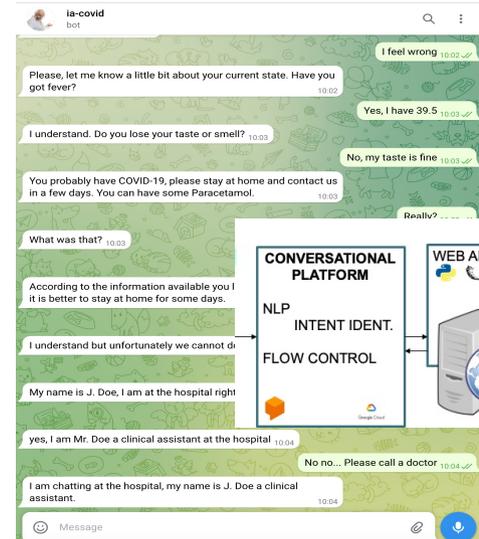
Lost in translation!
XAI
LIME, SHAP, Saliency



Medical AI
Black Boxes



GDPR
EU AI initiatives
Standards & MDR

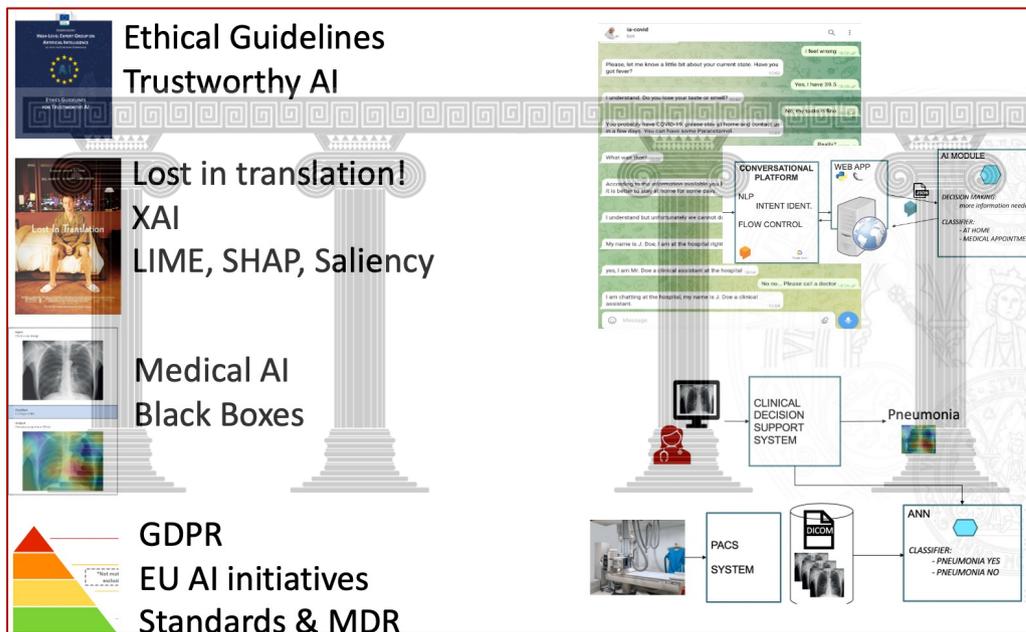


4. FINAL

TRUST IN MEDICAL AI



SEMINAR: FRAMING TRUST IN MEDICAL AI



Contact: Jose M. Juarez

jmjuarez@um.es



REFERENCES & READINGS

Materials & references of the seminar at:

<https://webs.um.es/jmjuarez/trustAlmedicine/>

