



Kunstmatige Intelligentie in Medische Beeldvorming

Associate Prof. Peter van Ooijen, MSc, PhD, CPHIMS
Dept. of Radiation Oncology, Associate Prof. AI in Medical Imaging
Data Science Center in Health, Coordinator Machine Learning Lab

Leerdoelen

- Inzicht krijgen in het proces van Kunstmatige Intelligentie
- Weten wat mogelijke valkuilen van KI zijn
- Besef van relevante wet- en regelgeving binnen de EU
- Kennis van bestaande methoden en checklists te gebruiken bij ontwikkeling en uitrol van KI



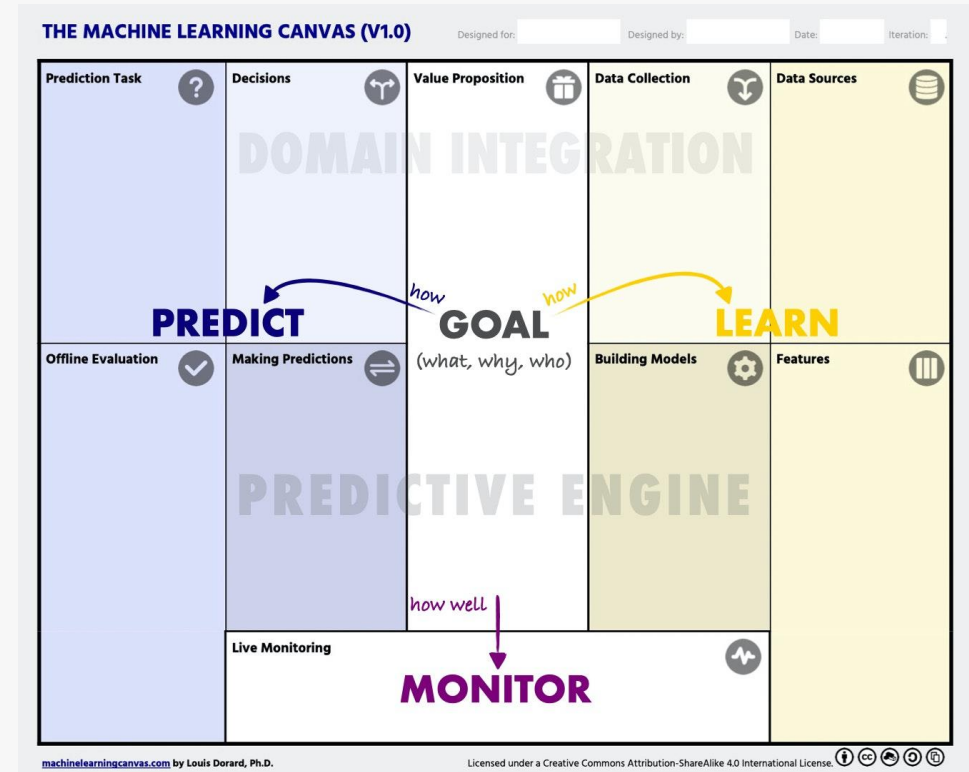
umcg

Stappen in het KI Process



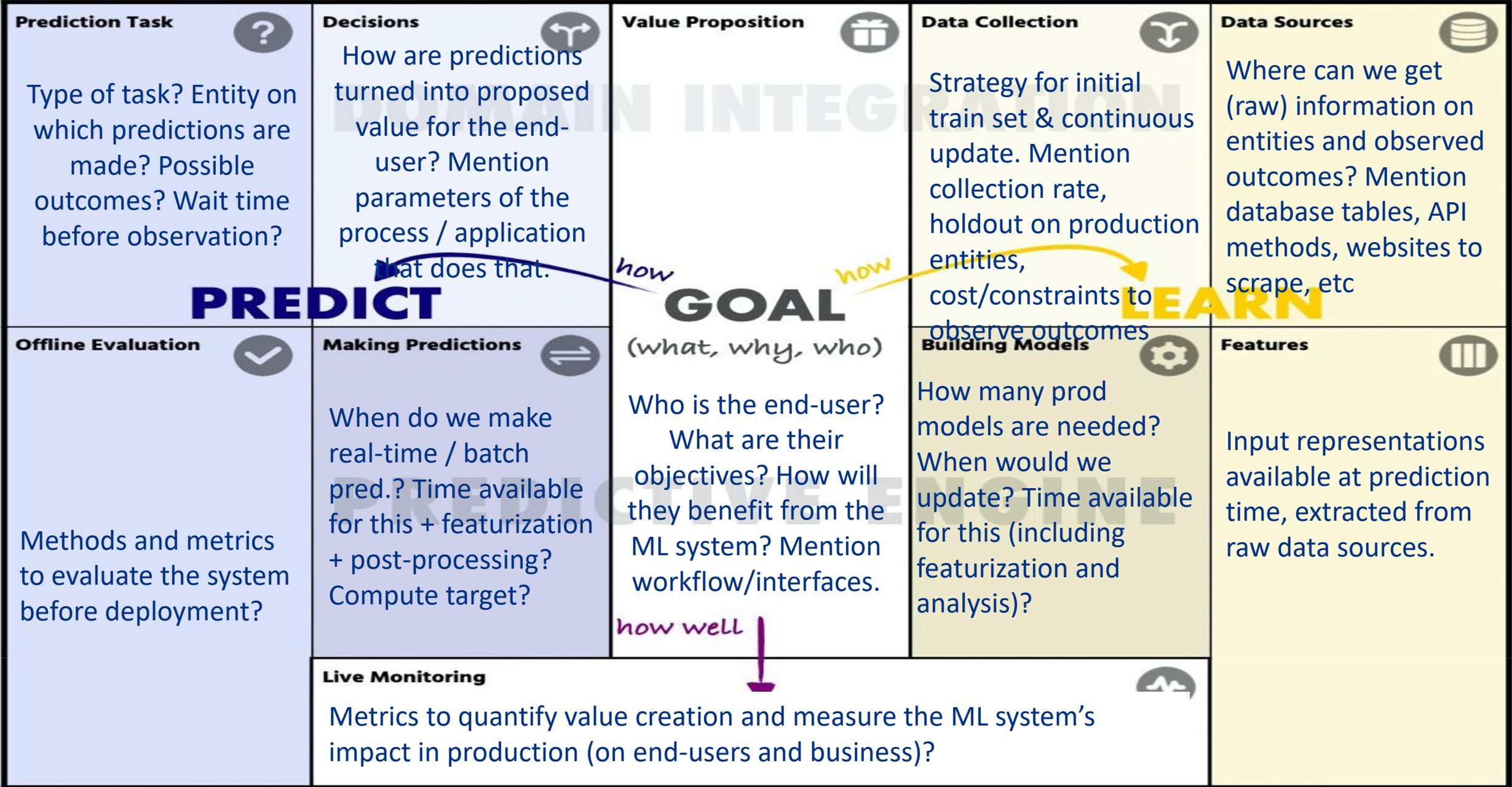
Start van een project Machine Learning Canvas

- Louis Dorard
- <https://www.ownml.co/machine-learning-canvas>



Waarom ML Canvas?

- Describe complex ML systems in a comprehensible and structured way.
- ***Early engagement of stakeholders*** in the ML development process
- Get the key elements of a project
- Assess feasibility
- Detect bottlenecks and technical constraints early on
- ***Collaboration in the team***
- Plan work and choose the right tech



Pneumothorax: definition and diagnosis






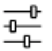
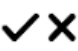


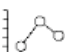
Pneumothorax (PTX)

- Air in the pleural space, causing partial collapse of the lung
- Occurrence of 15% in blunt thorax trauma cases
- Has to be treated prior to mechanical ventilation to avoid tension PTX (can lead to shock or death)

Diagnosis

- CT is gold standard
- Chest X-ray recommended for initial assessment (ATLS guidelines): fast, simple, cheap, always possible



Decisions  How are predictions used to make decisions that provide the proposed value to the end-user? <i>Provide the score of the Xray to the emergency physician. Based on the score the physician will make a decision on how to further treat the patient or if additional diagnosis (e.g. performing Computed Tomography) is needed.</i>	ML task  Input, output to predict, type of problem. We want to be able to answer the question <i>"Is there a pneumothorax in the presented thoracic X Ray?"</i> right after performing the X Ray. Input: X Ray image Output: "Pneumothorax" or "no Pneumothorax" => Binary Classification Identification of the pneumothorax region => image segmentation	Value Propositions  What are we trying to do for the end-user(s) of the predictive system? What objectives are we serving? Provide automated diagnosis of pneumothorax in thoracic X-ray images in the emergency department as input for the emergency physician. Aim to facilitate the process and avoid delays in treatment of patients by faster diagnosis of the -ray images. In case of doubt, expert review by a radiologist is still required.	Data Sources  Which raw data sources can we use (internal and external)? X-rays acquired at emergency department Diagnosis description of the images. Weak or strong labels?	Collecting Data  How do we get new data to learn from (inputs and outputs)? PACS Export images + Structured Reports Manual tagging of the images by a human expert. Alternative could be to automatically retrieve the tagging using NLP from the structured report by the radiologist.
Making Predictions  When do we make predictions on new inputs? How long do we have to featurize a new input and make a prediction? <i>Predictions have to be made right after the thoracic X Ray is required. The prediction should be provided in seconds because of the emergency setting in which waiting is not possible.</i>	Offline Evaluation  Methods and metrics to evaluate the system before deployment. <i>False negative predictions should be avoided. System should be fine-tuned for this.</i>	Features  Input representations extracted from raw data sources. <i>Image features of the thoracic Xrays</i> <i>Diagnosis information (yes or no pneumothorax) from free text reports using NLP</i>	Building Models  When do we create/update models with new training data? How long do we have to featurize training inputs and create a model? <i>One model is built for general use based on the retrospectively collected data.</i> <i>feedback from the expert users will be collected during operation.</i> <i>New model will be trained after collecting feedback for 6 months. Only 'frozen' models can be used in the clinic, no continuous learning.</i>	Live Evaluation and Monitoring Methods and metrics to evaluate the system after deployment, and to quantify value creation. <i>Changes of acquisition setup should be tracked (could influence performance of the trained network).</i> <i>Comparison of processing time of thoracic X-ray before and after AI implementation.</i> <i>Recording the agreement and disagreement of the emergency physician with the decision as made by AI</i> 

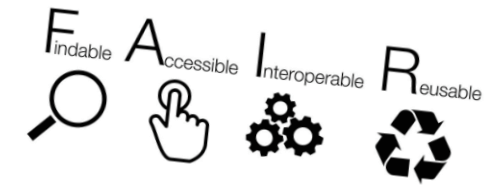
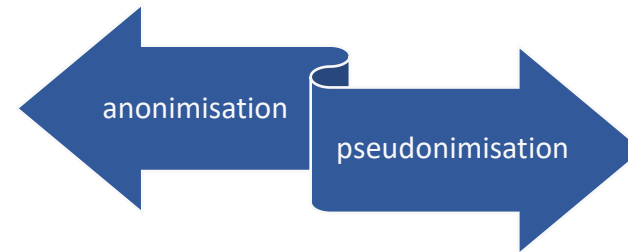


Data Collection



Data Collection

- Probably the most important part of the process
- Very often the limiting factor
- Data Quality
- Standardization
- Small data set
- Bias



https://www.emedicinehealth.com/informed_consent/article_em.htm

Datasheets for Datasets

arXiv.org > cs > arXiv:1803.09010

Search...

Help | Advanced

Computer Science > Databases

[Submitted on 23 Mar 2018 (v1), last revised 19 Mar 2020 (this version, v7)]

Datasheets for Datasets

Timnit Gebru, Jamie Morgenstern, Briana Vecchione, Jennifer Wortman Vaughan, Hanna Wallach, Hal Daumé III, Kate Crawford

The machine learning community currently has no standardized process for documenting datasets, which can lead to severe consequences in high-stakes domains. To address this gap, we propose datasheets for datasets. In the electronics industry, every component, no matter how simple or complex, is accompanied with a datasheet that describes its operating characteristics, test results, recommended uses, and other information. By analogy, we propose that every dataset be accompanied with a datasheet that documents its motivation, composition, collection process, recommended uses, and so on. Datasheets for datasets will facilitate better communication between dataset creators and dataset consumers, and encourage the machine learning community to prioritize transparency and accountability.

Comments: Working Paper, comments are encouraged

Subjects: **Databases (cs.DB)**; Artificial Intelligence (cs.AI); Machine Learning (cs.LG)

Cite as: [arXiv:1803.09010](#) [cs.DB]

(or [arXiv:1803.09010v7](#) [cs.DB] for this version)



Data sheets for datasets

- Motivation
- Composition
- Collection process
- Preprocessing/cleaning/labeling
- Uses
- Distribution
- Maintenance



<https://arxiv.org/abs/1803.09010>

<https://www.microsoft.com/en-us/research/project/datasheets-for-datasets/>

Data Collection

European Radiology
<https://doi.org/10.1007/s00330-020-07453-w>

LETTER TO THE EDITOR

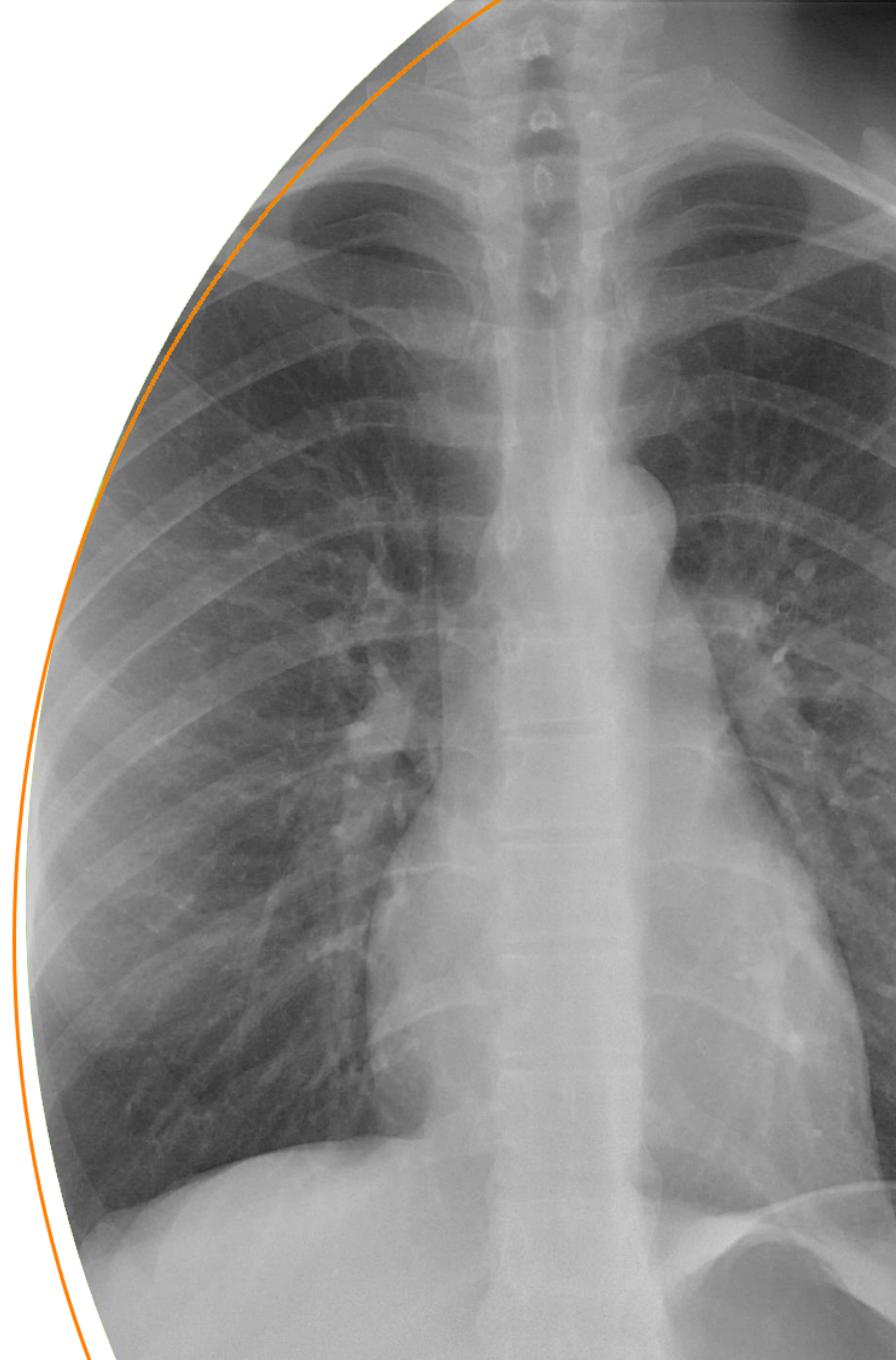
COVID-19, AI enthusiasts, and toy datasets: radiology without radiologists

H. R. Tizhoosh^{1,2}  • Jennifer Fratesi³

Received: 4 September 2020 / Revised: 23 September 2020 / Accepted: 2 November 2020

Pneumothorax detection

- Dataset already available in public domain
 - Con: No influence on data collection/selection
 - Pro: easy to start with
- ChestXray14¹
- 112.120 Xrays



Data Annotation



umcg

Data Annotation/Labeling

- Annotation/Label quality!
 - Appropriate for the task at hand? (global/local labels)
 - Class balance? (the more imbalance, the more difficult the task, the more data likely needed to solve it)
 - Consistency? (everyone used the same labeling guidelines/procedure?)
 - Label format? (computer readable? Or hidden in e.g. radiology report?)

Data Annotation/Labeling



Heart segmentation:
Local, pixel level labels

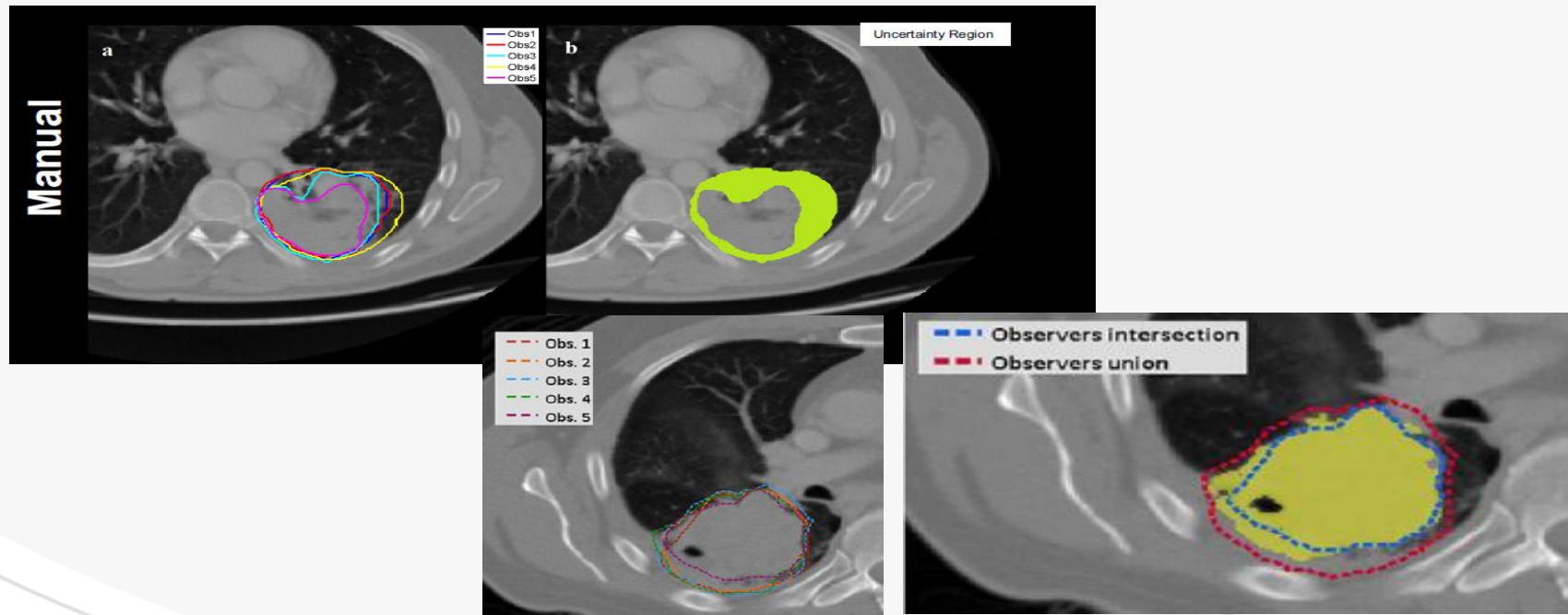


Fracture detection:
Global or local labels?

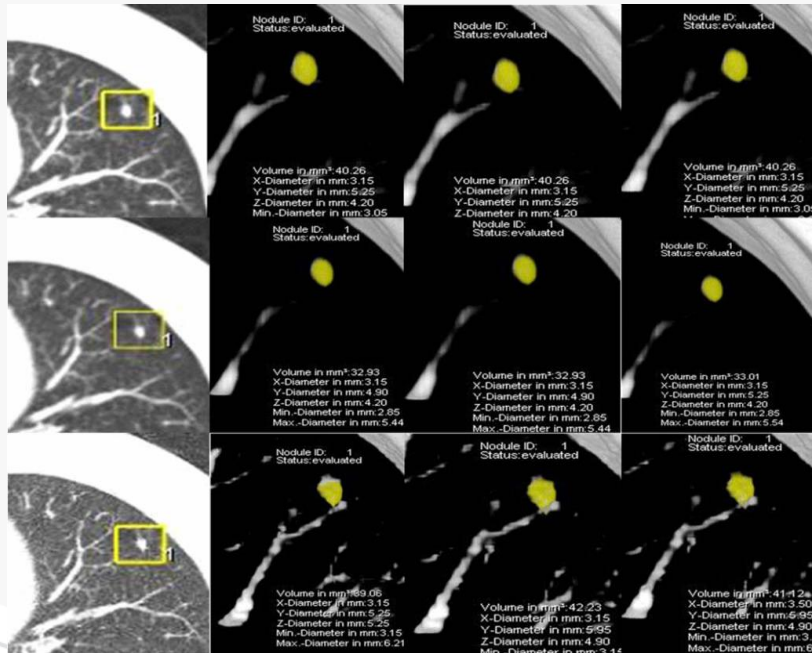
Data Annotation/Labeling

- Manual delineation of regions of interest (ROI)
 - Time consuming
 - inter- en intra-observer variability
- Semi-automatic: some user interaction
 - Dependent on starting ROI
 - Variability of seedpoint selection
 - Less user dependent than manual

Data Annotation/Labeling



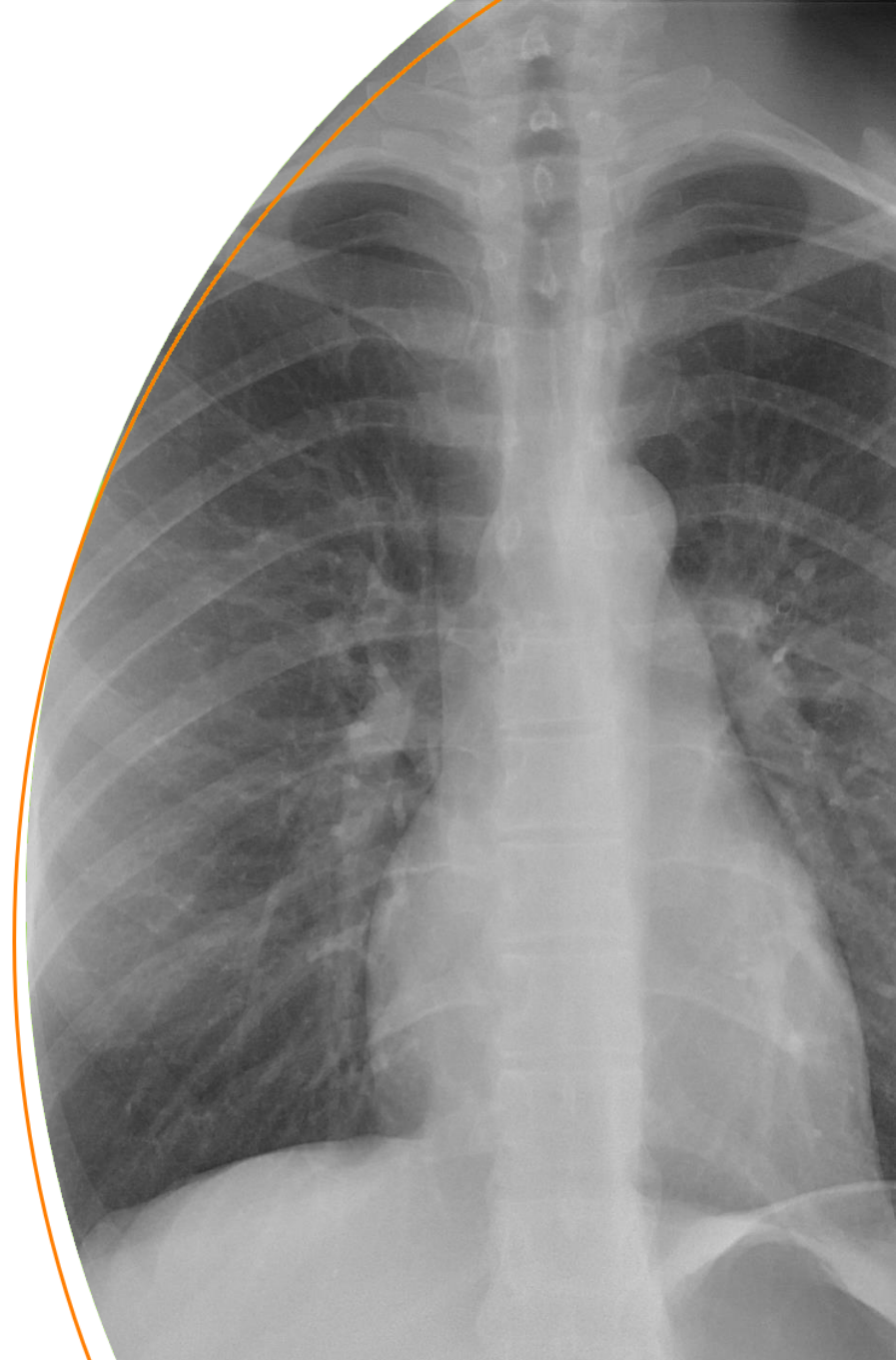
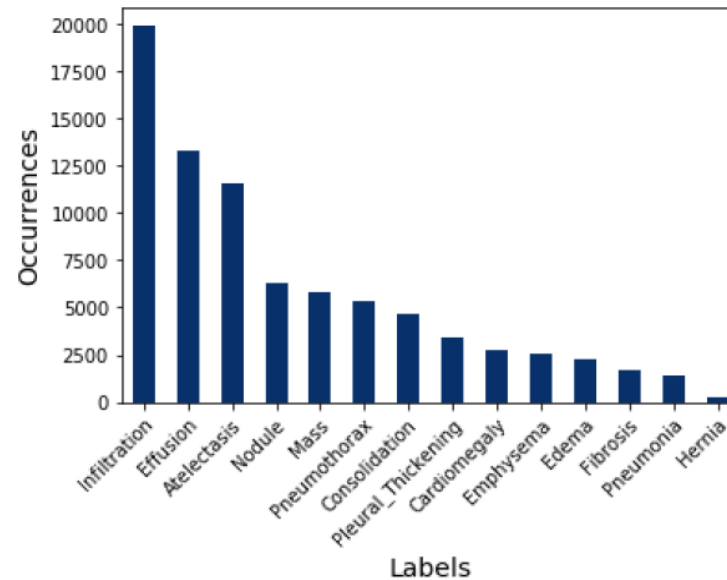
Data Annotation/Labeling



- Variation in CT reconstruction algorithm
- 1 mm soft kernel, 2 mm soft kernel, and 2 mm sharp kernel

Pneumothorax detection

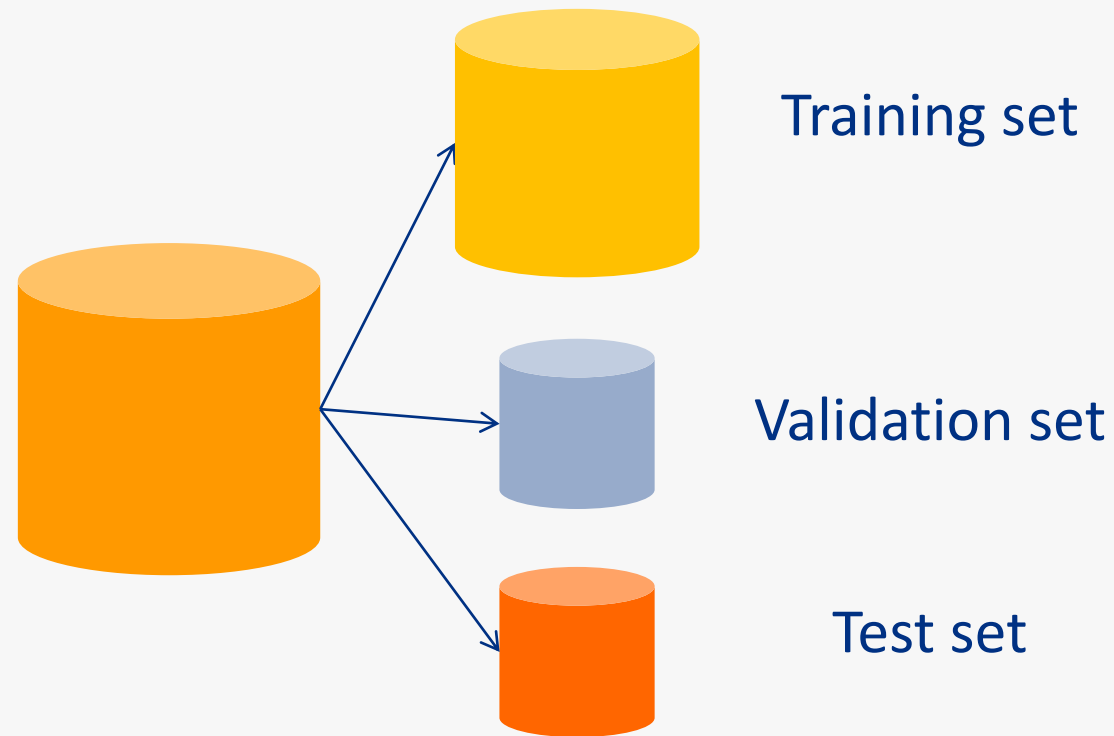
- Publicly available X-ray dataset: ChestXray14¹
- 112.120 Xrays
- 14 labeled diseases
- ~60.000 images with 'No Findings'



Model training



Data split

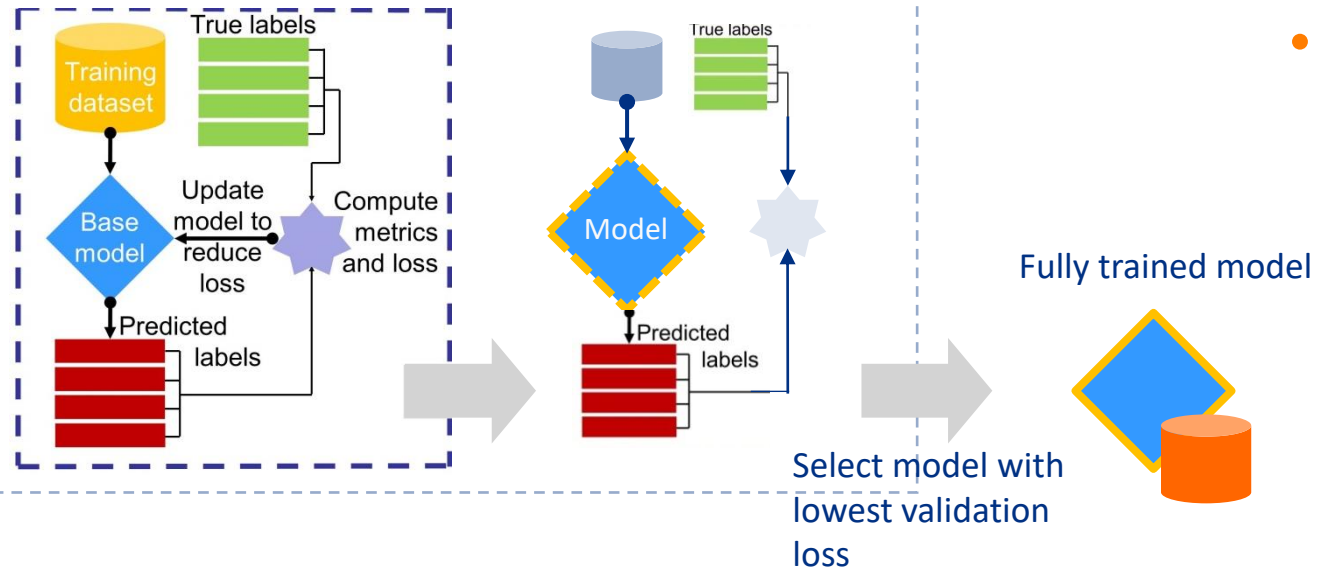


So how does a computer learn?

Repeat for n epochs

Stop training when validation loss does not decrease for x epochs

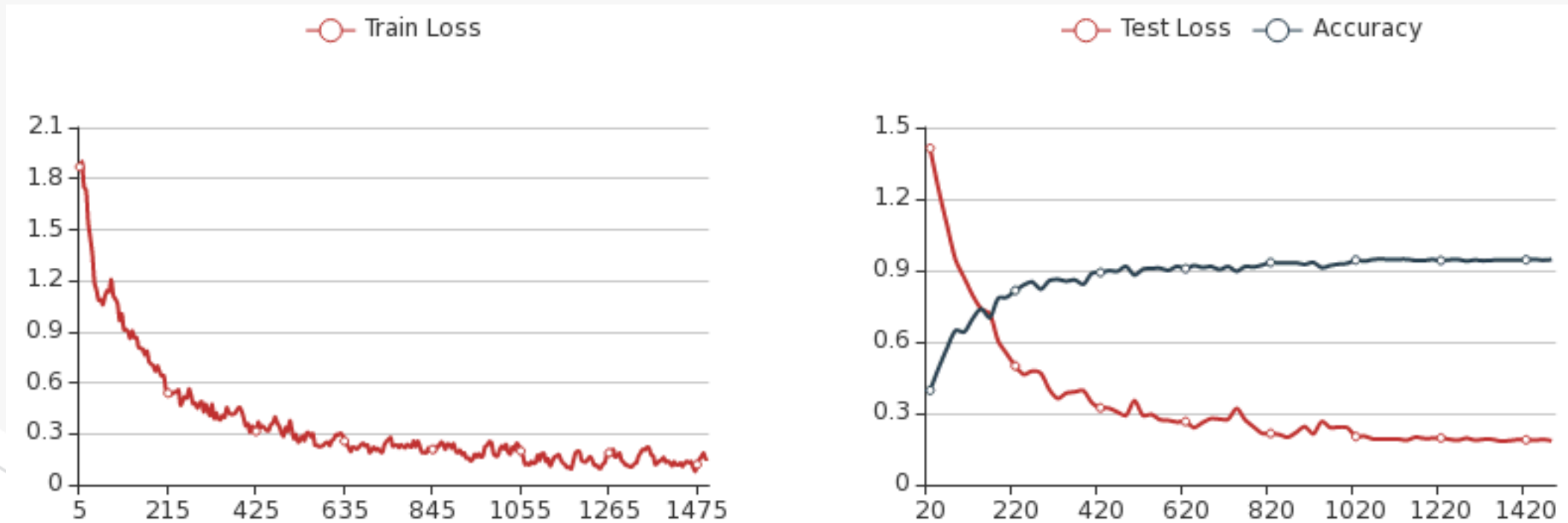
Train for 1 epoch



Epoch: one iteration over the full dataset

- Learns by calculating how wrong it is: The loss
- Training phase:
 - Provide input and the correct output
 - Computer predicts the output
 - Calculates the error from the correct output
 - Reiterates

Loss graph



Model training

InceptionV3 architecture¹

Input: 299x299

Data augmentation:

Rotations $\alpha \in [-25^\circ, +25^\circ]$

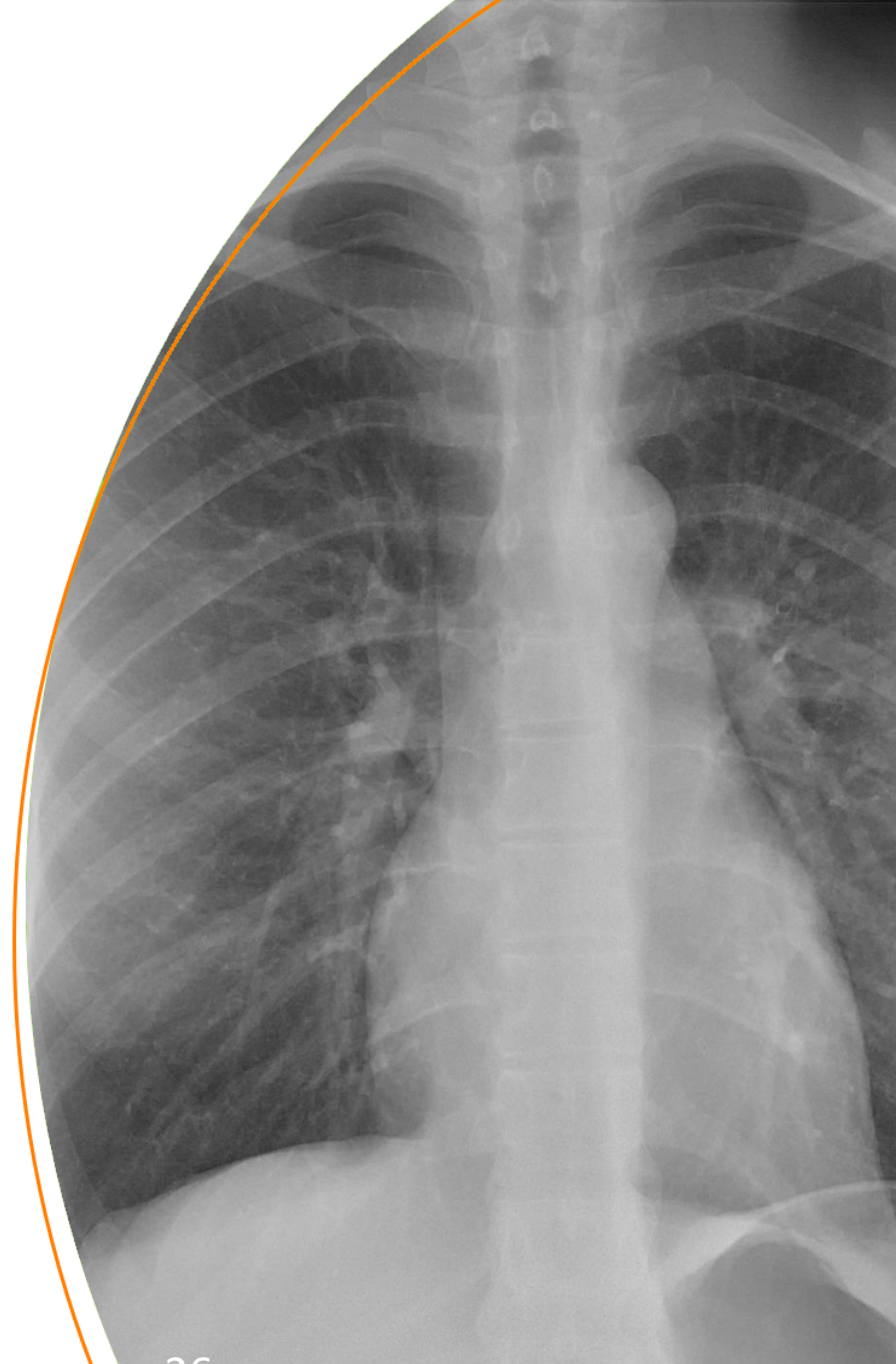
Translations $\Delta x, \Delta y \in [-0.1, +0.1]$

Brightness $\Delta B \in [-0.3, +0.3]$

Dataset split:

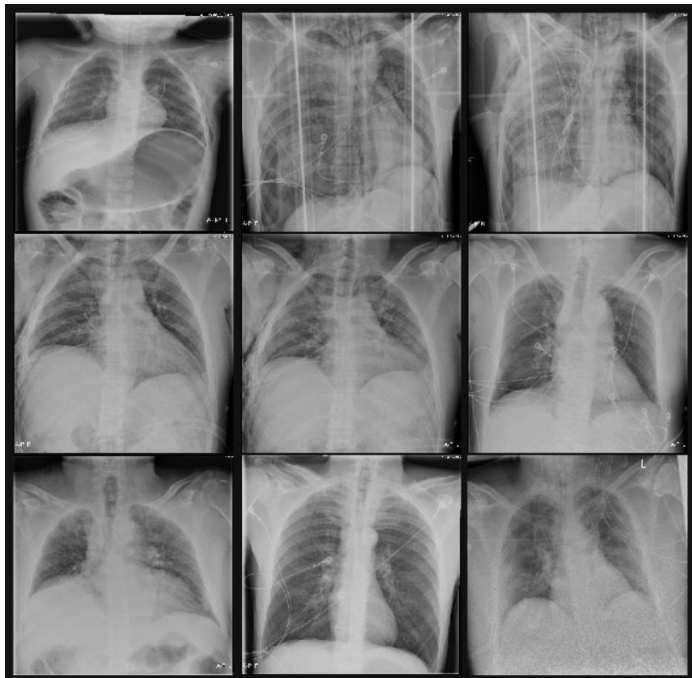
Train/validation/test: 80/10/10

Trained for 100 epochs

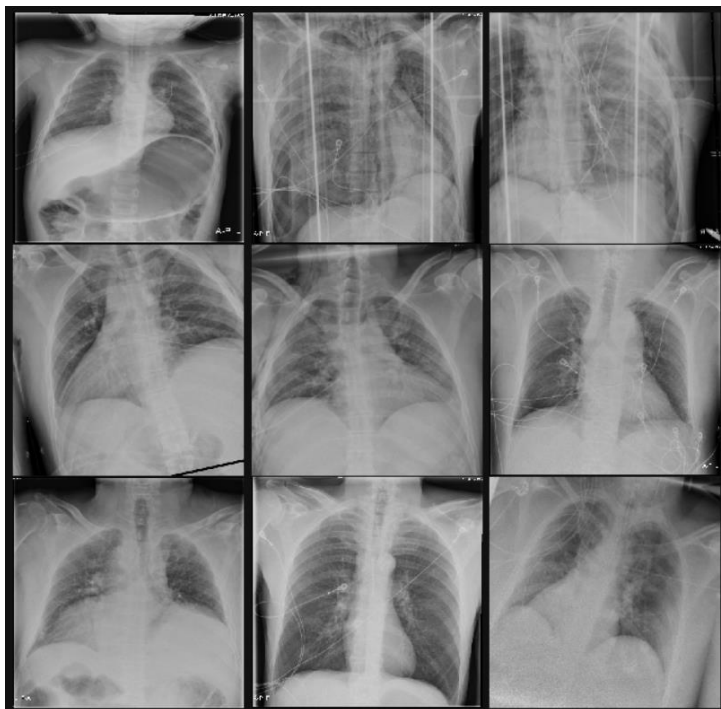


Data augmentation

No augmentations



Random horizontal flip
Random rotation/scaling/translation
Random brightness/contrast adjustment



Model Validation and Testing



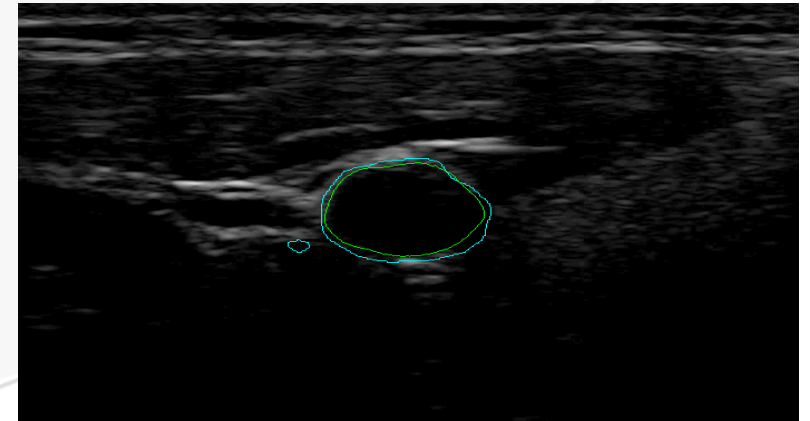
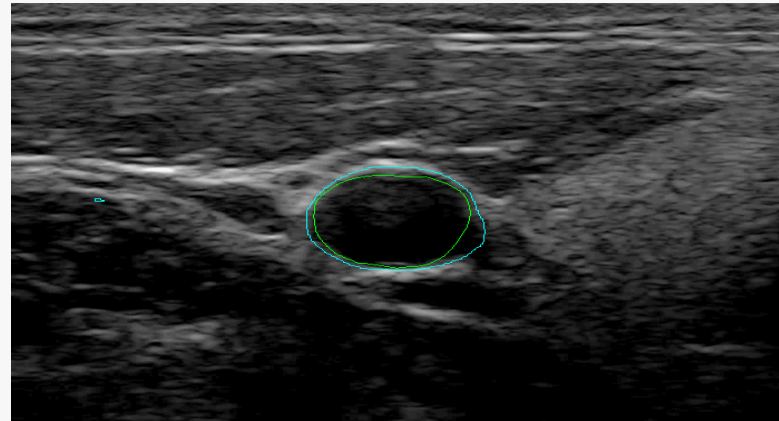
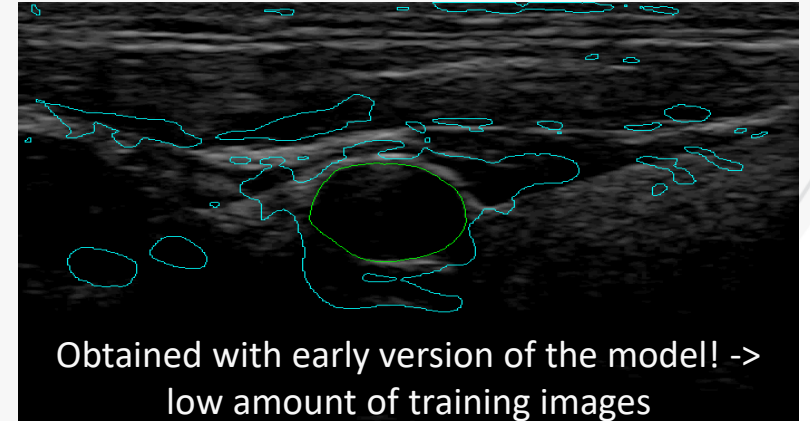
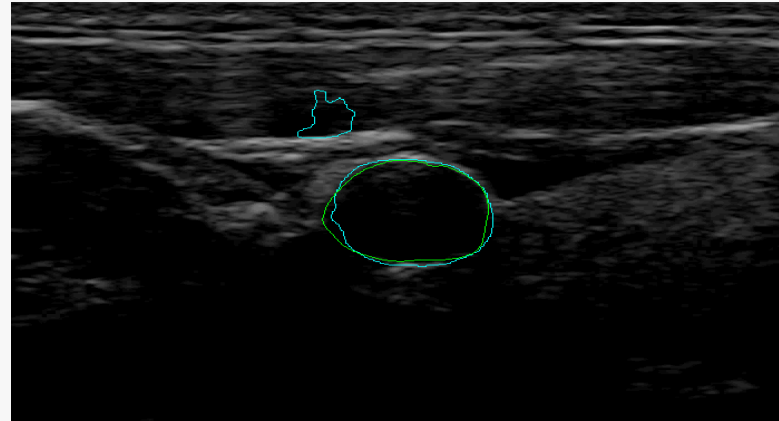
Example Carotid artery ultrasound

AI training Summary	
Task	Segmentation of carotid artery lumen
Number of images	1060 labeled, 2500 total
Annotation procedure	Manual annotation
Performance measure	Pixelwise dice coefficient $dc = \frac{2TP}{2TP + FP + FN}$
Current model performance on test set	$dc = 0.93$
Train/validation/test split	848/106/106
Network architecture	Unet

Examples of low dice results

Ground Truth

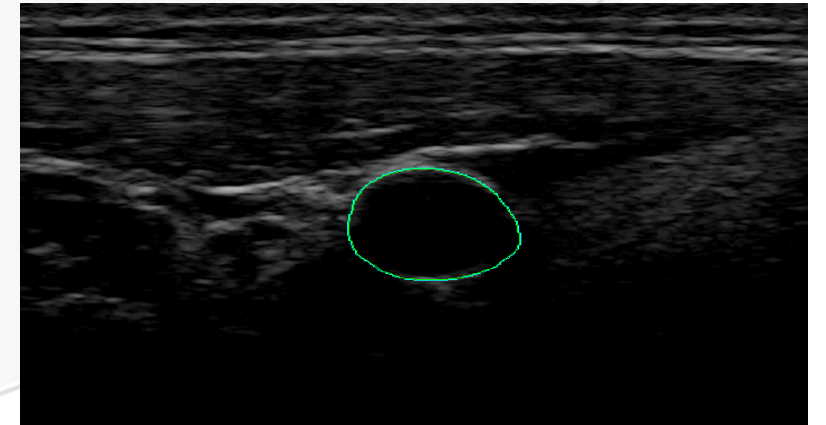
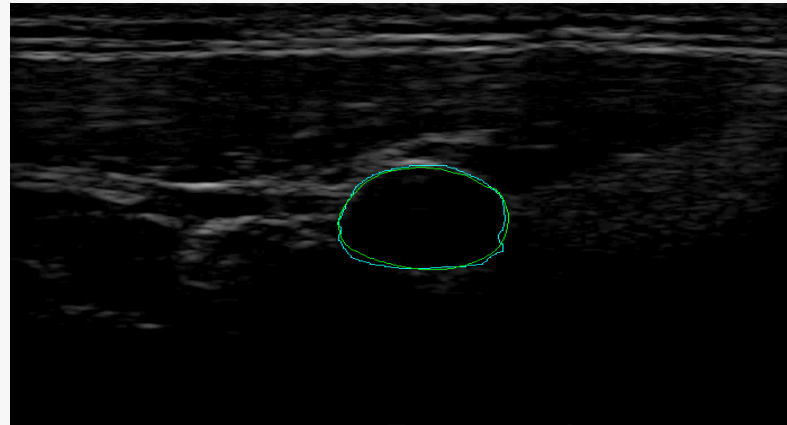
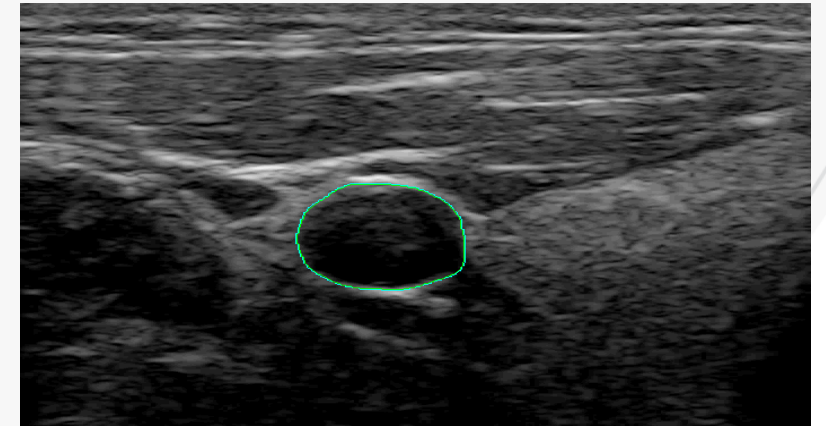
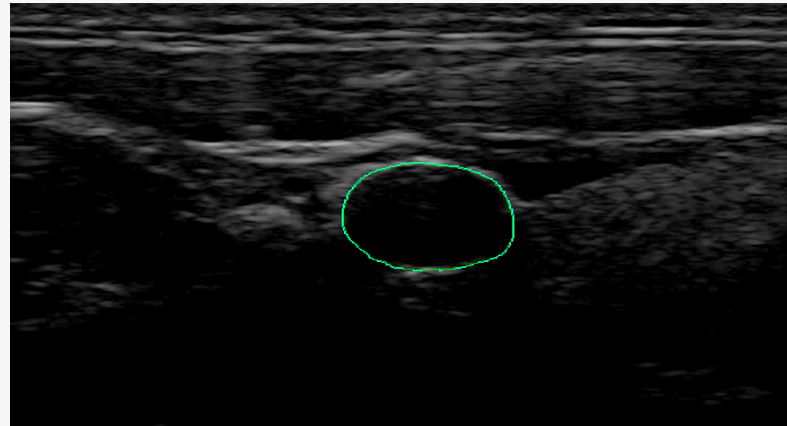
AI output



Examples of high dice results

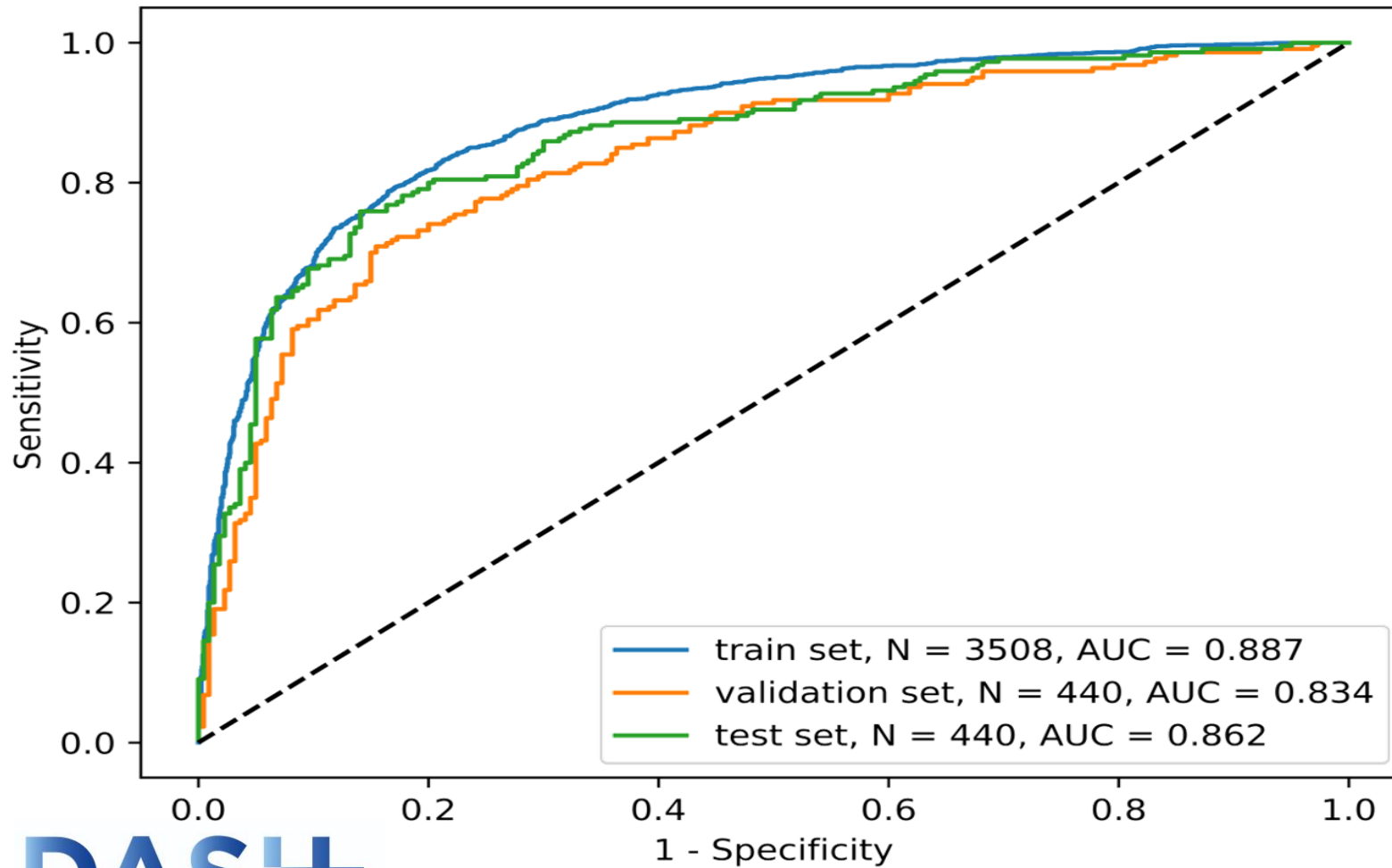
Ground Truth

AI output



Model Validation and Testing

ROC curves Pneumothorax detection (Inception v3)



Explainability

Understanding “why” the network arrives at its prediction

1.Failure mode analysis

Inspect misclassified samples, see if there is a pattern

2.Do the classification pixel-wise (i.e. make it a segmentation problem)

Requires contours as labels (often not available)

Requires different model architecture

3.“Saliency map” or “heatmap” of pixel relevance

Indicates how much each pixel in the input contributes to model output

Layer-wise relevance propagation

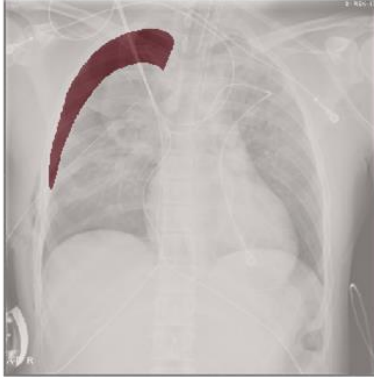


Explainability – Grad-CAM

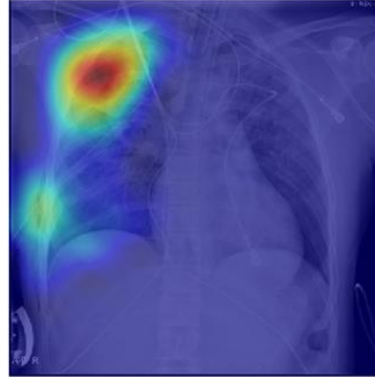
Input Image



Mask



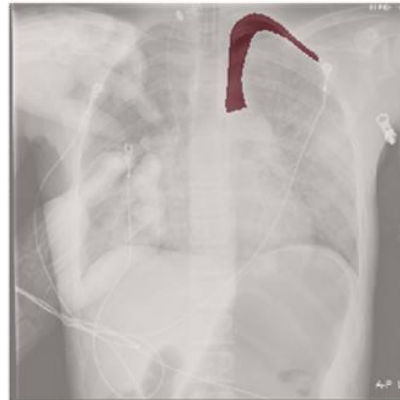
Grad-CAM



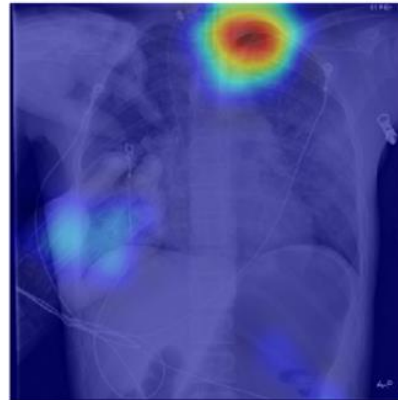
Input Image



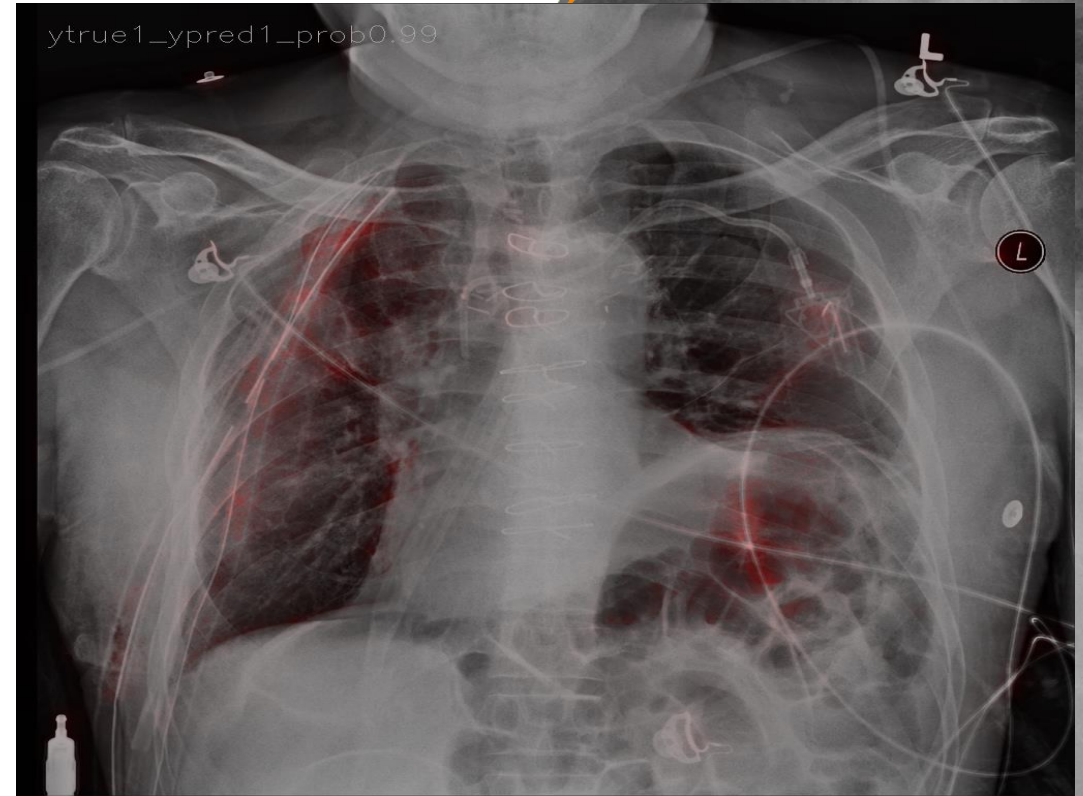
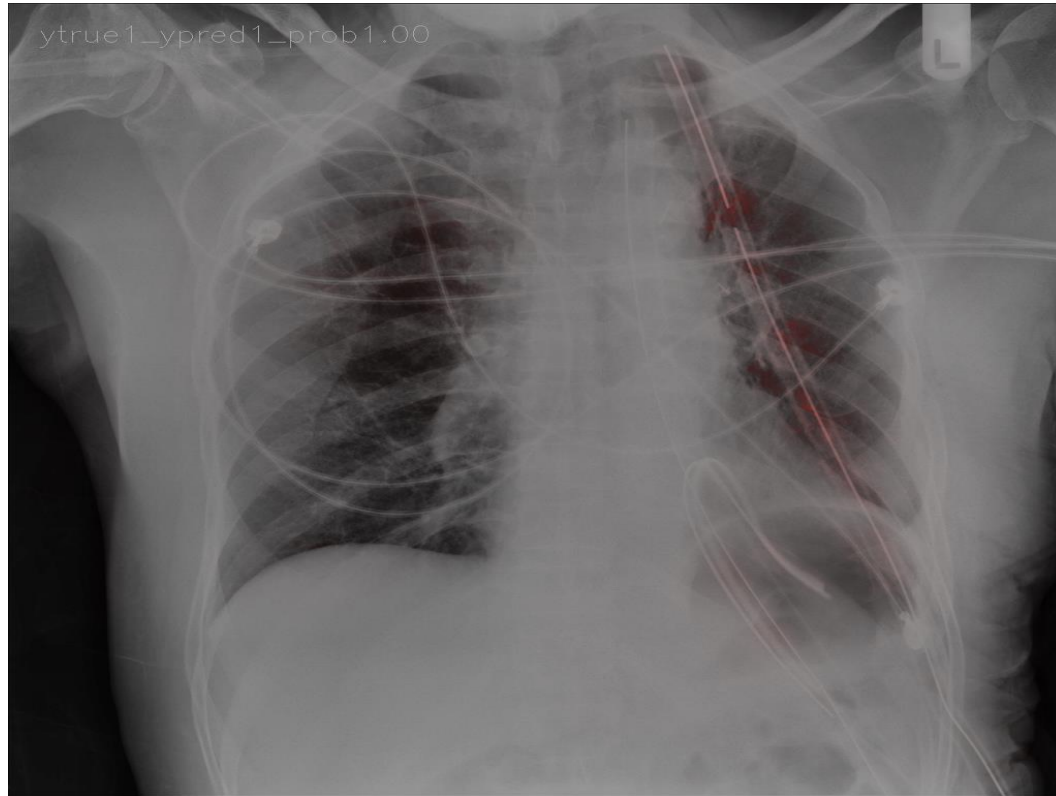
Mask



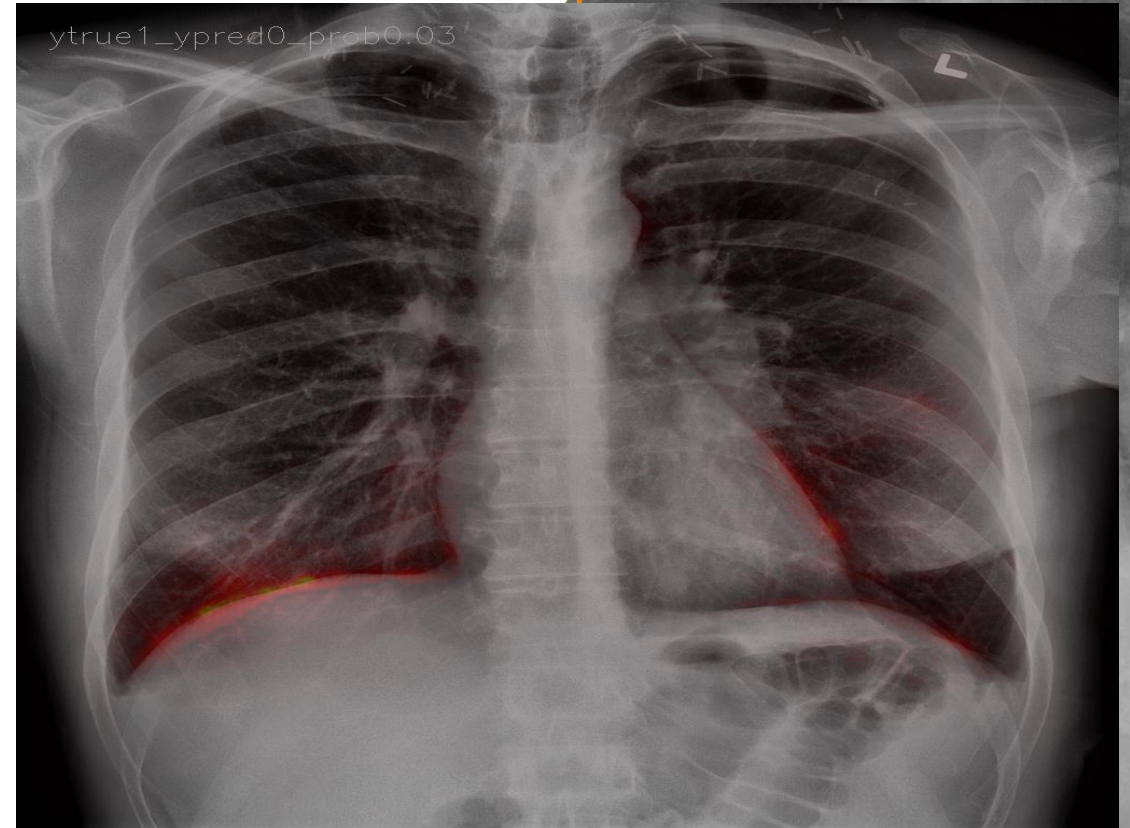
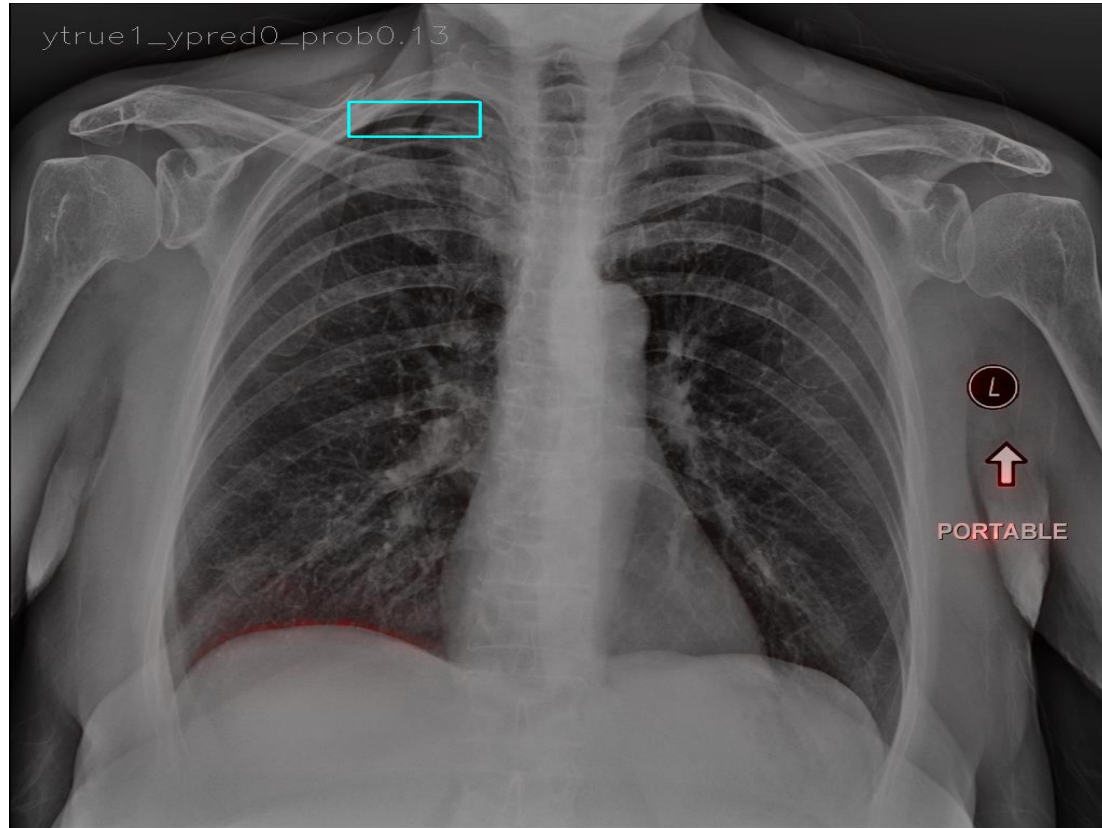
Grad-CAM



Explainability – Deep Taylor Decomposition

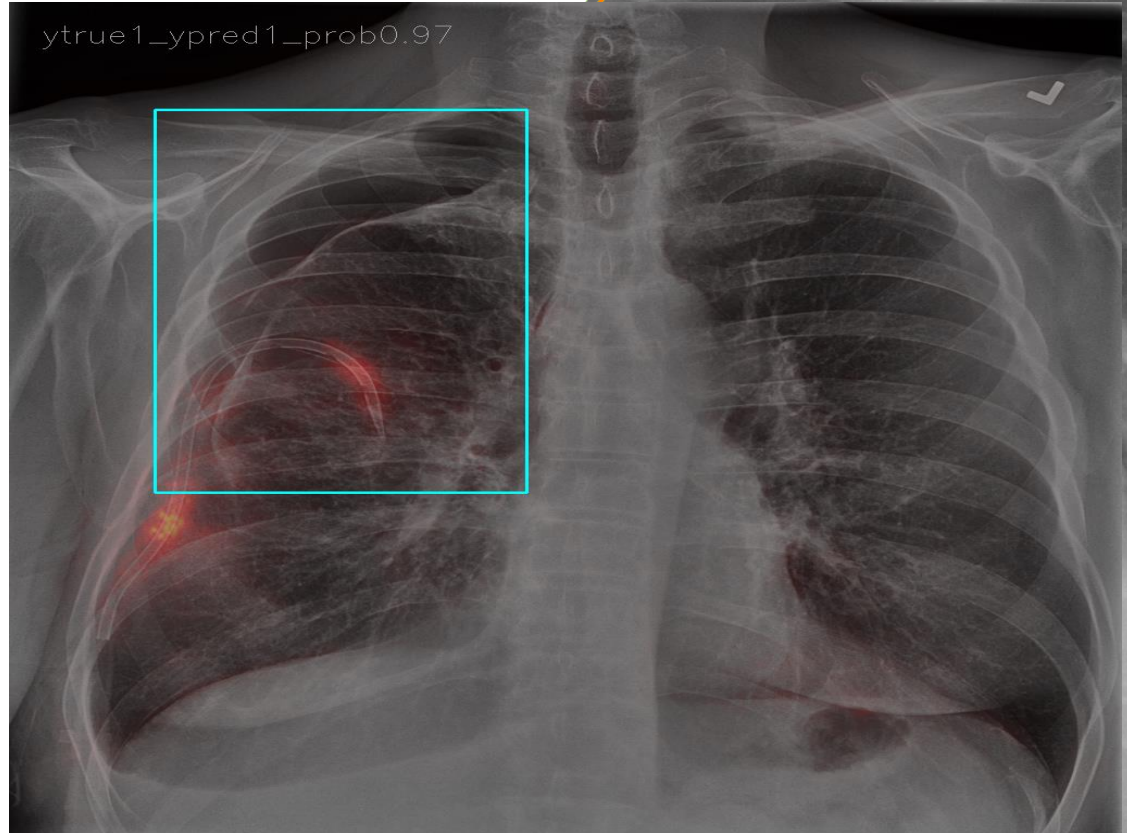
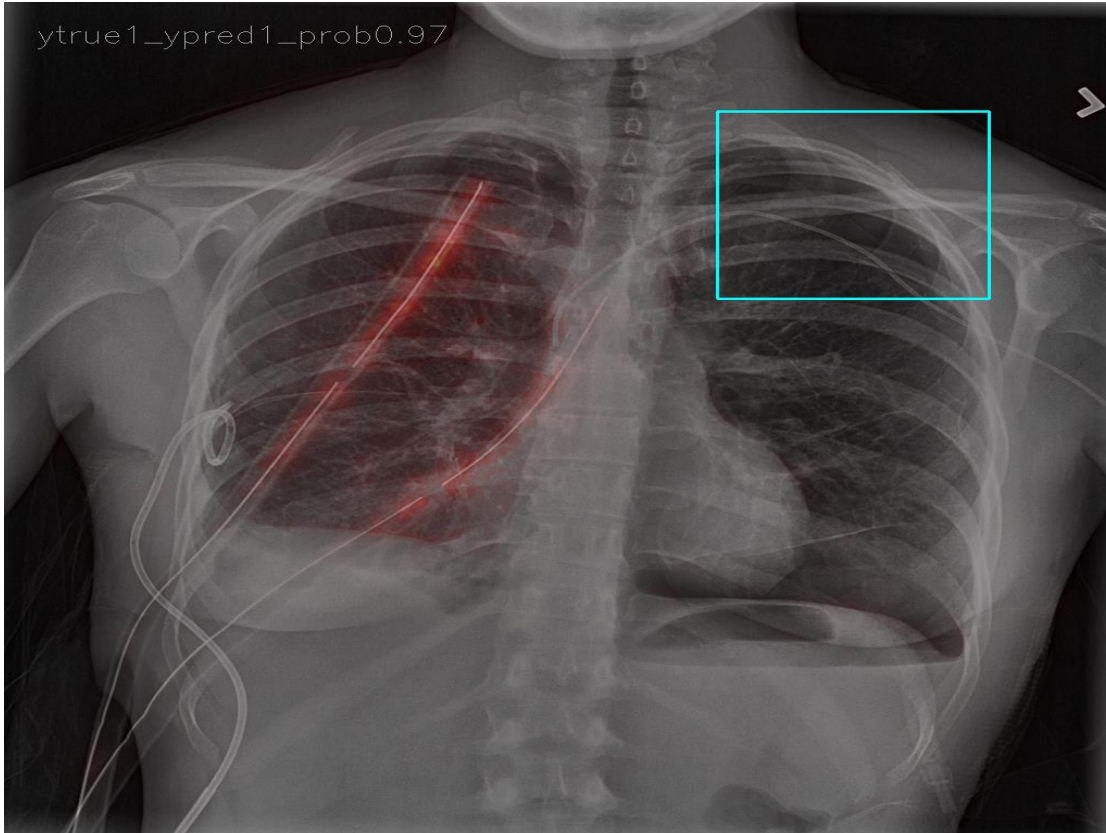


Explainability – Deep Taylor Decomposition



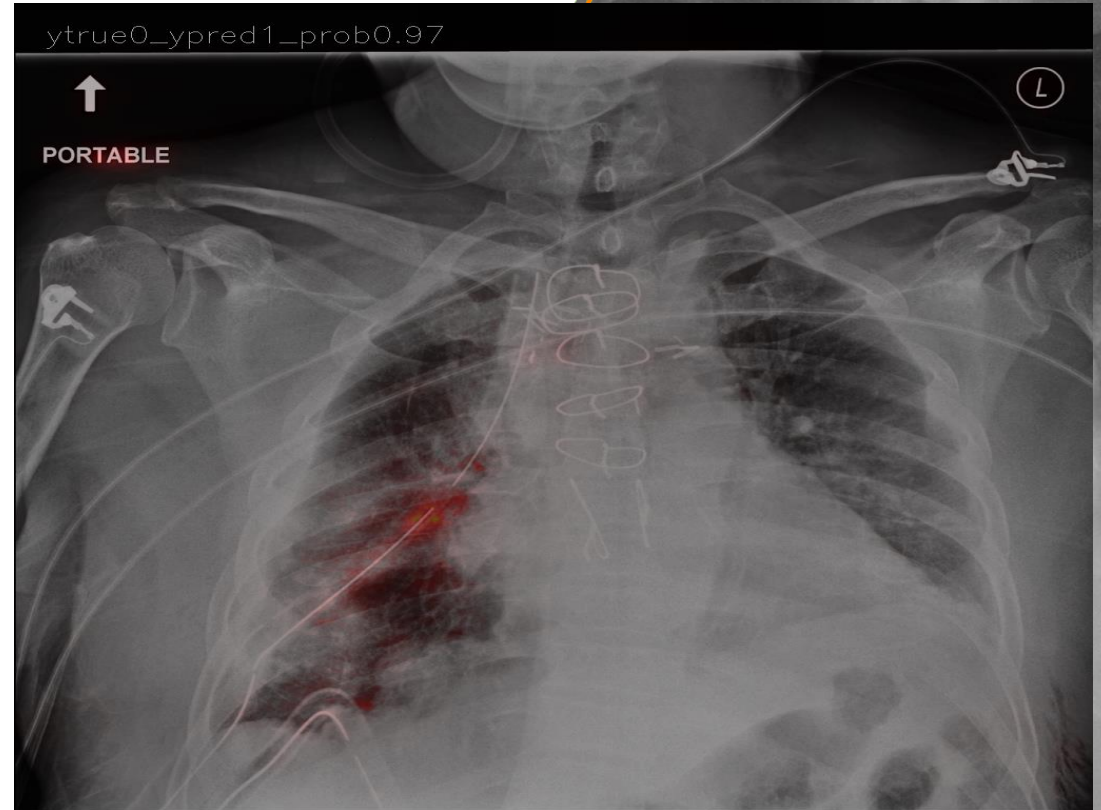
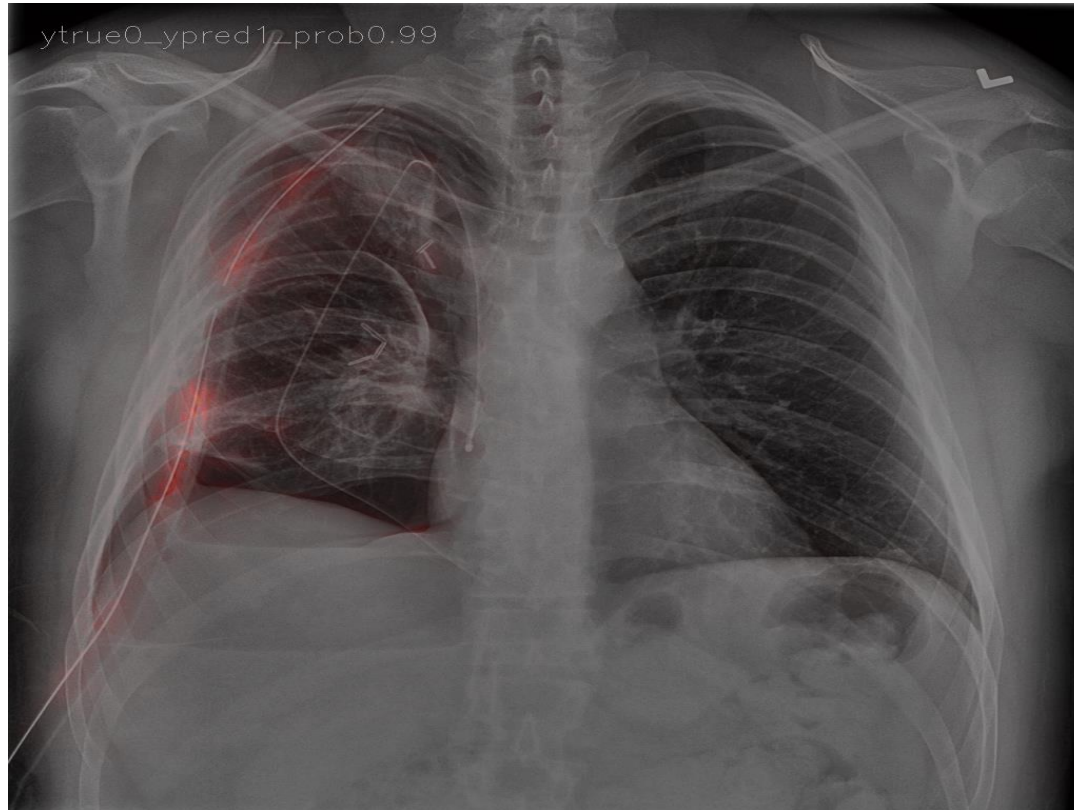
FALSE NEGATIVES

Explainability – Deep Taylor Decomposition



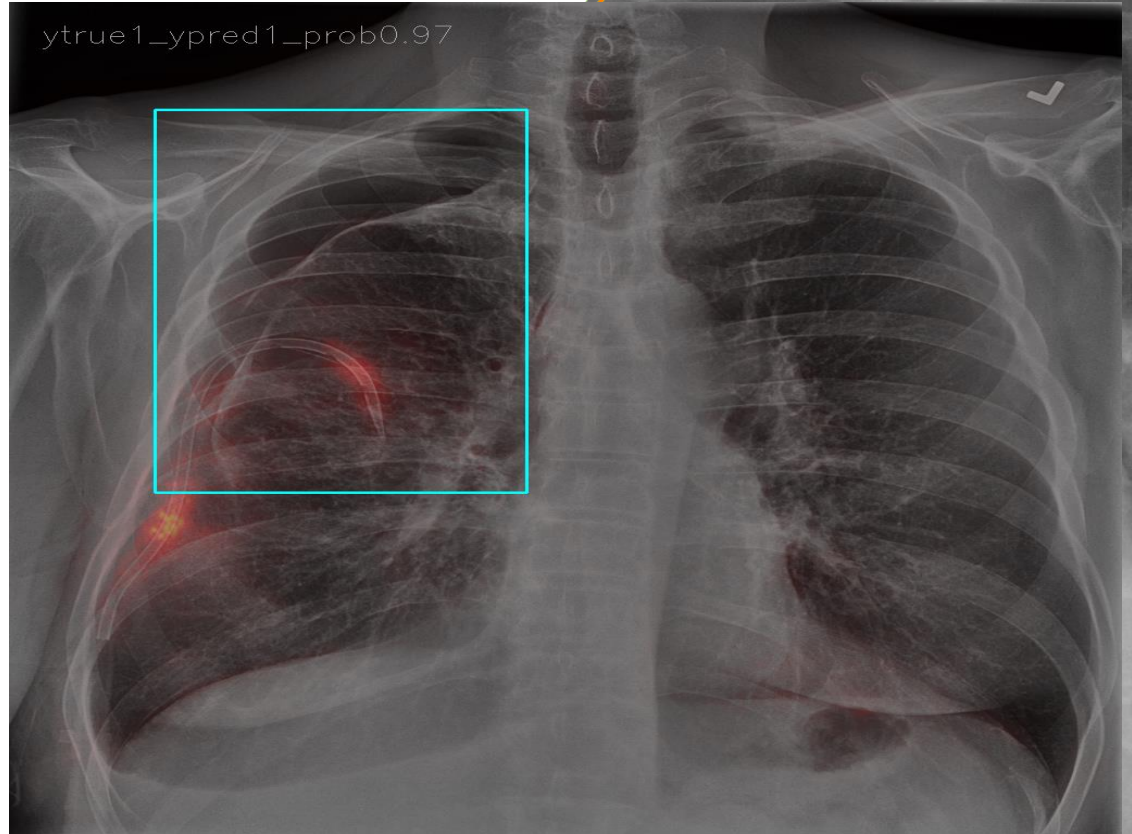
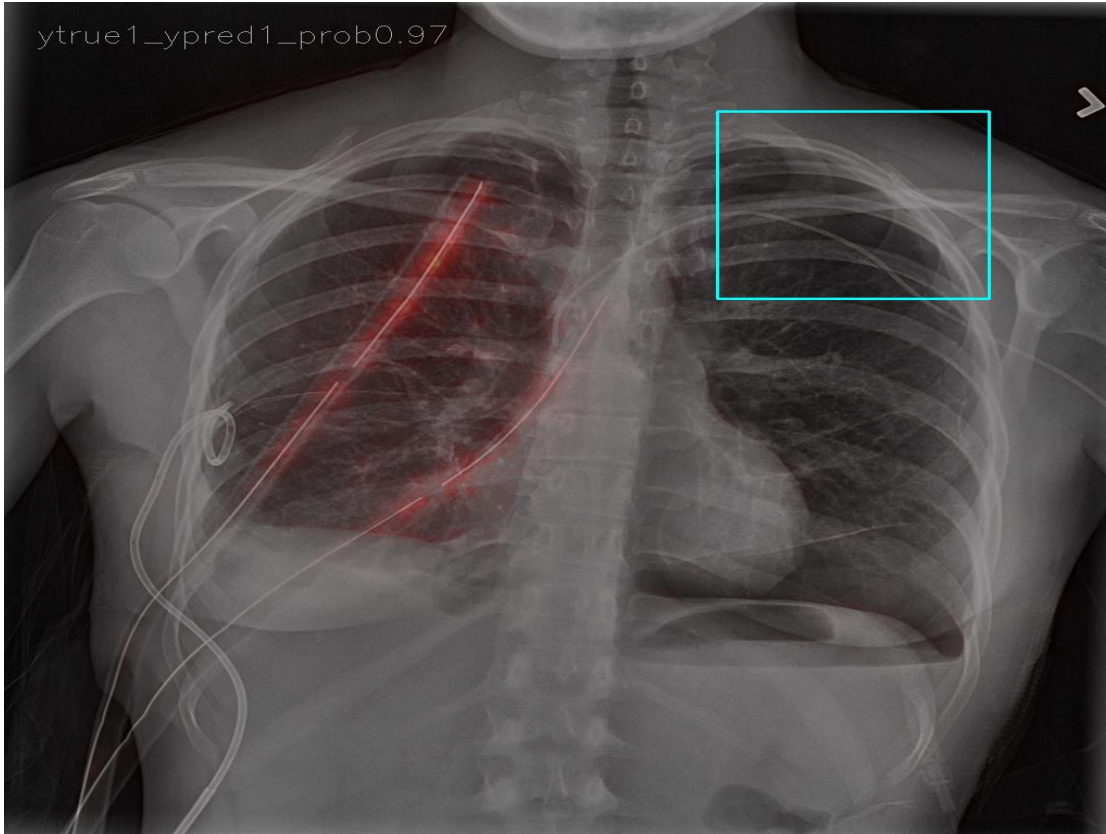
TRUE POSITIVES

Explainability – Deep Taylor Decomposition



FALSE POSITIVES

Explainability – Deep Taylor Decomposition



TRUE POSITIVES

Dataset Bias!

- CHX 14: PTX images are biased
 - Large fraction of post-treatment images in dataset
- Causes a bias towards X-rays with drain
- Correct classification, but not suitable for diagnostics
- Impossible to tell from the metrics alone
 - Heatmap/relevance map useful tool to understand model decision making



Model Deployment and Quality Assurance



AI Market situation

- <https://grand-challenge.org/aiforradiology/>

European Radiology (2021) 31:3797–3804
<https://doi.org/10.1007/s00330-021-07892-z>




AI for Radiology
an implementation guide

IMAGING INFORMATICS AND ARTIFICIAL INTELLIGENCE



Artificial intelligence in radiology: 100 commercially available products and their scientific evidence

Kicky G. van Leeuwen¹  • Steven Schalekamp¹ • Matthieu J. C. M. Rutten^{1,2} • Bram van Ginneken¹ • Maarten de Rooij¹



AI Market situation

- 100 CE Marked AI products from 54 vendors
- 64/100 have no peer review evidence
- 36/100 have evidence from 237 papers
- 116/237 papers were independent and not (co-)funded or (co-)authored by the vendor
- Only 18/100 AI products have demonstrated (potential) clinical impact!!



To buy or not to buy—evaluating commercial AI solutions in radiology (the ECLAIR guidelines)

Patrick Omoumi¹  • Alexis Ducarouge² • Antoine Tournier² • Hugh Harvey³ • Charles E. Kahn Jr⁴ • Fanny Louvet-de Verchère⁵ • Daniel Pinto Dos Santos⁶ • Tobias Kober⁷ • Jonas Richiardi¹



Checklist AI assessment

1. Relevance
2. Performance and Validation
3. Usability and Integration
4. Regulatory and Legal Aspects
5. Financial and Support services considerations



Table 2 Top 10 questions to consider

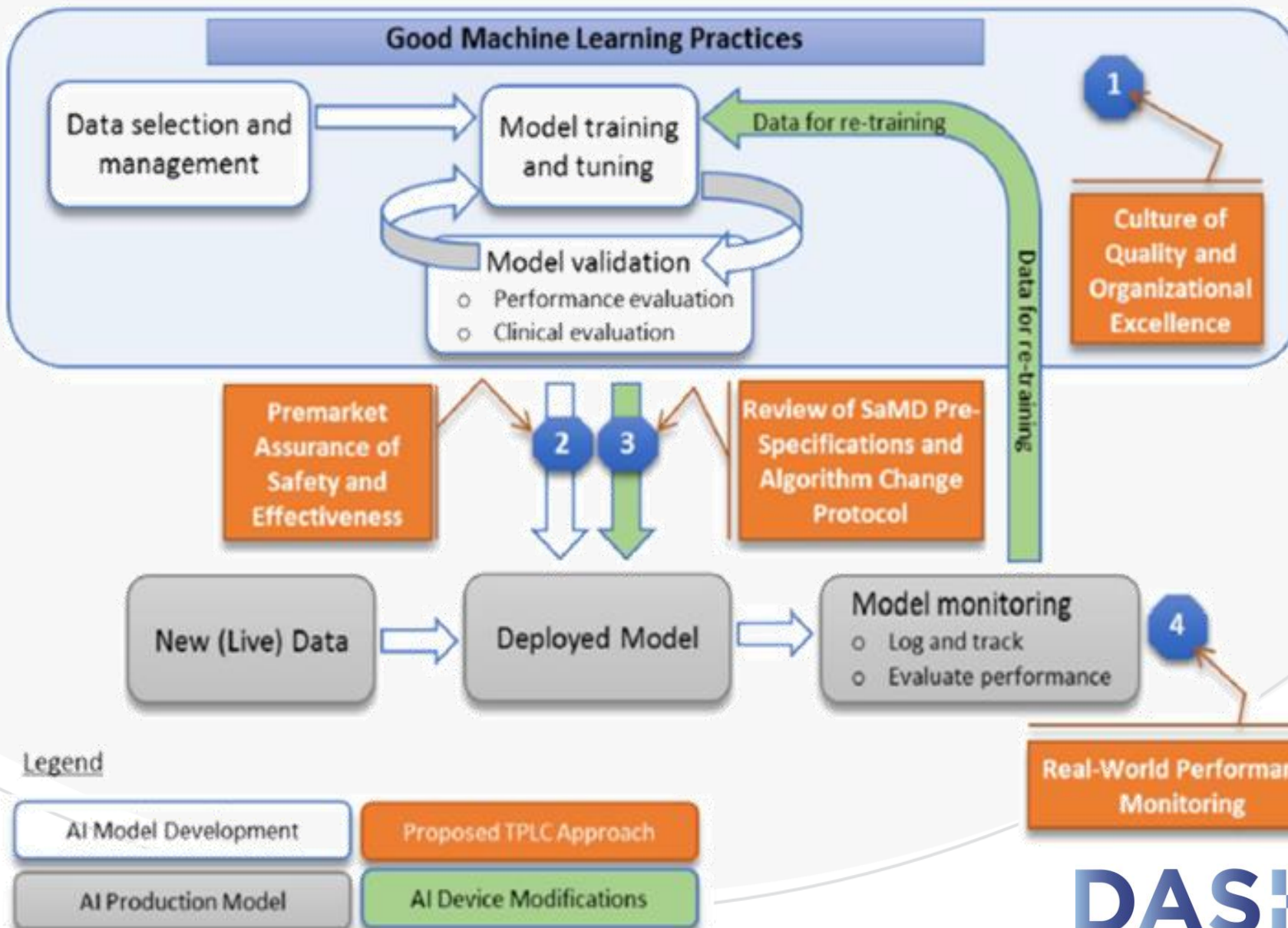
1. What problem is the application intended to solve, and who is the application designed for?
 2. What are the potential benefits and risks, and for whom?
 3. Has the algorithm been rigorously and independently validated?
 4. How can the application be integrated into your clinical workflow and is the solution interoperable with your existing software?
 5. What are the IT infrastructure requirements?
 6. Does the application conform to the medical device and the personal data protection regulations of the target country, and what class of regulation does it conform to?
 7. Have return on investment (RoI) analyses been performed?
 8. How is the maintenance of the product ensured?
 9. How are user training and follow-up handled?
 10. How will potential malfunctions or erroneous results be handled?
-



Challenges in QA

Performance of AI is influenced by:

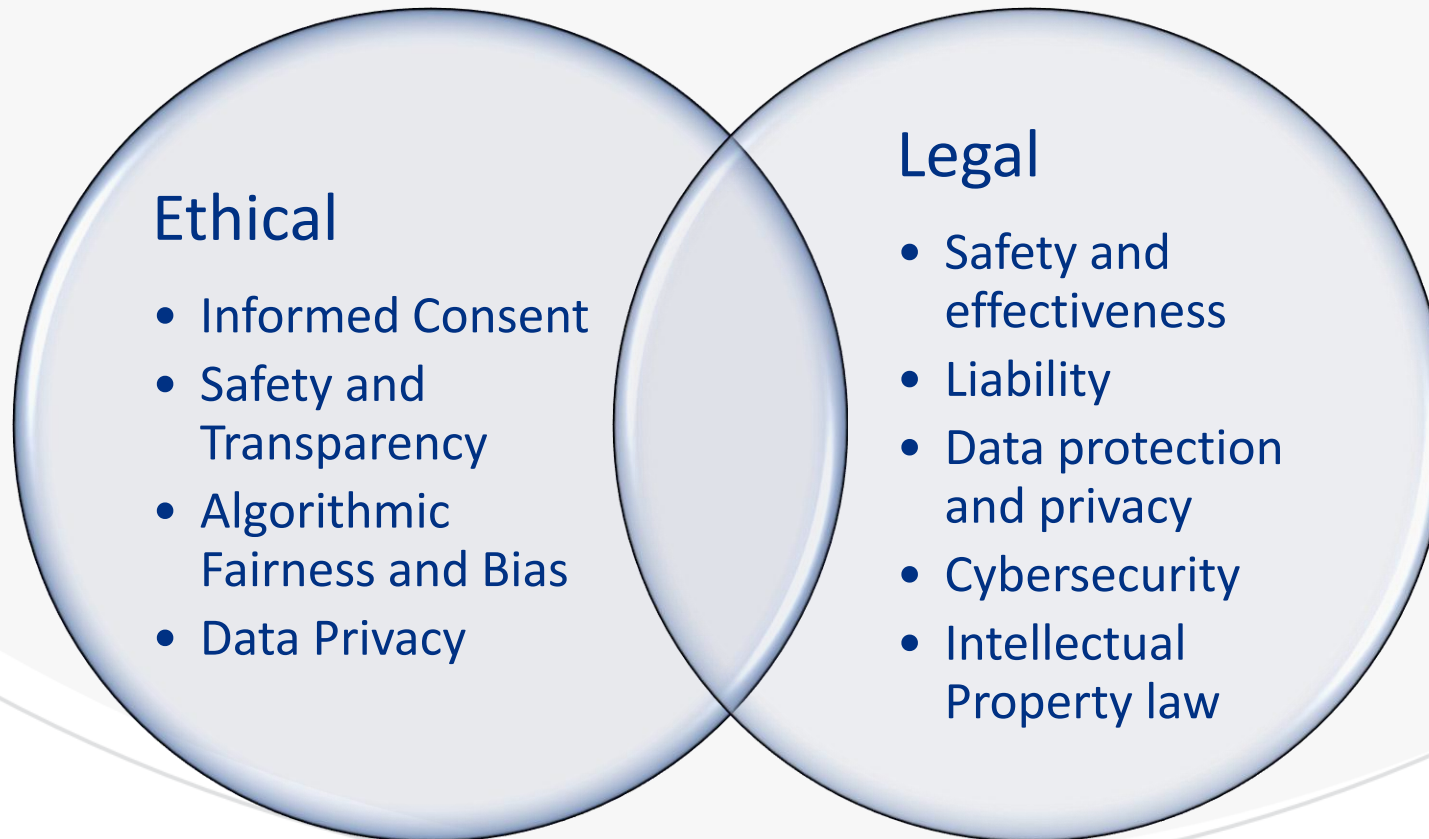
- Data Drift
 - New scanner coming in
- Concept Drift
 - Ideas about diagnosis/conclusion changes
- Population Drift
 - Patient groups change



ELSA Issues



ELSA Issues



Regulatory Concepts EU

- General Data Protection Regulation (GDPR)
- Medical Device Regulation (MDR)
 - Software as a Medical Device (SaMD)
- Ethical Guidelines for Trustworthy AI



Regulatory Concepts - GDPR

“Right for an individual to obtain a meaningful explanation when automated (algorithmic) decision-making is involved.”



Regulatory Concepts - MDR

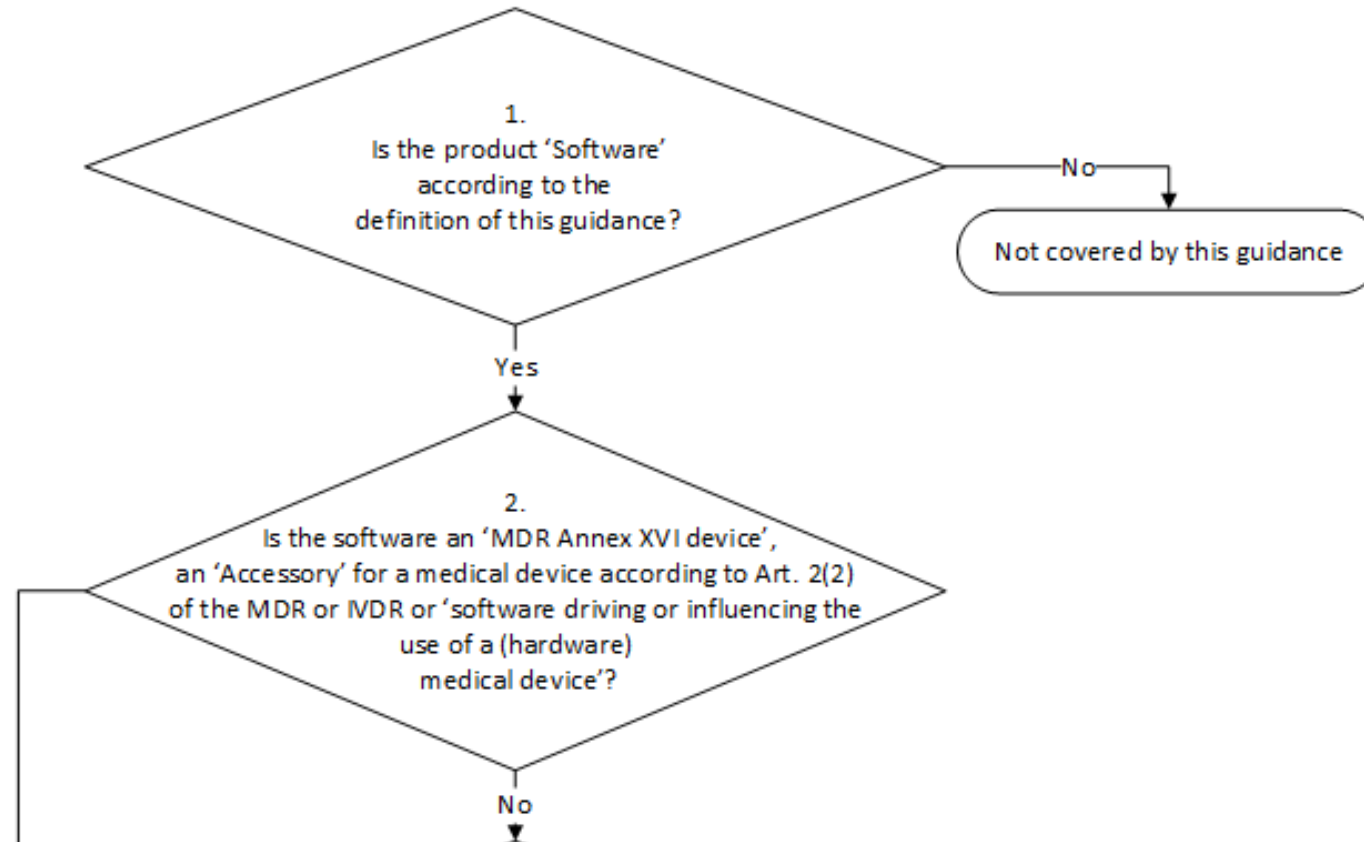


Regulatory Concepts - MDR

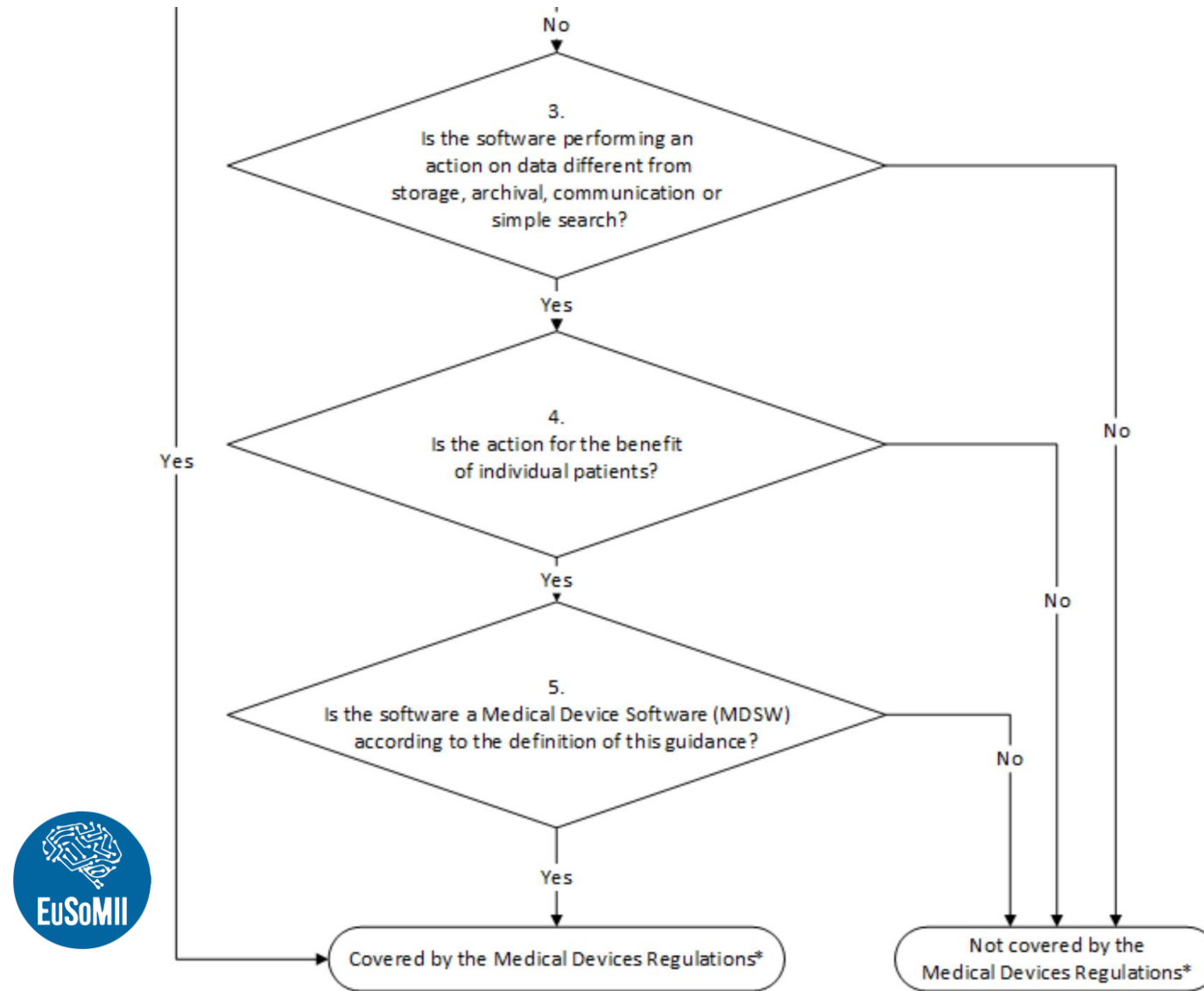
Significance of Information provided by the MDSW to a healthcare situation related to diagnosis /therapy				
State of Healthcare Situation of patient condition		High Treat or diagnose ~IMDRF 5.1.1	Medium Drives Clinical Management ~IMDRF 5.1.2	Low Informs Clinical Management (everything else)
	Critical situation or patient condition ~IMDRF 5.2.1	Class III Category IV.i	Class IIb Category III.i	Class IIa Category II.i
	Serious situation or patient condition ~IMDRF 5.2.2	Class IIb Category III.ii	Class IIa Category II.ii	Class IIa Category I.ii
	Non-Serious situation or patient condition (everything else)	Class IIa Category II.iii	Class IIa Category I.iii	Class IIa Category I.i



Regulatory Concepts - MDR



Regulatory Concepts - MDR



Regulatory Concepts - MDR

- Also for in-house developed software!
- CE marking is not required when software is developed, used, and maintained only within a health institution, under certain conditions:
 - software cannot be transferred to another legal entity
 - health institution justifies that the target patient group's specific needs cannot be met or cannot be met at the appropriate level of performance by an equivalent device available on the market.
 - The health institution needs to prepare documentation containing e.g. design and performance information of the device AND to review experience gained from clinical use of the software and take necessary corrective actions.



Regulatory Concepts - Ethical Guidelines

Trustworthy AI

- Ethical Guidelines for trustworthy AI
- Independent High-Level experts group on AI
- Started by the European Commission in June 2018



Regulatory Concepts - Ethical Guidelines

Trustworthy AI

7 Requirements:

1. Human Agency and Oversight;
2. Technical Robustness and Safety;
3. Privacy and Data Governance;
4. Transparency;
5. Diversity, Non-discrimination and Fairness;
6. Societal and Environmental Well-being;
7. Accountability.





ALTAI for Test

[Notes](#)

Sections of the ALTAI

- 📁 Human Agency and Oversight
- 📁 Technical Robustness and Safety
- 📁 Privacy and Data Governance
- 📁 Transparency
- 📁 Diversity, Non-Discrimination and Fairness
- 📁 Societal and Environmental Well-being
- 📁 Accountability

Legend of progression symbols

Human Agency and Oversight

AI systems should support human autonomy and decision-making, as prescribed by the principle of respect for human autonomy. This requires that AI systems should both act as enablers to a democratic, flourishing and equitable society by supporting the user's agency and upholding fundamental rights, which should be underpinned by human oversight. In this section, we are asking you to assess the AI system in terms of the respect for human agency, as well as human oversight.

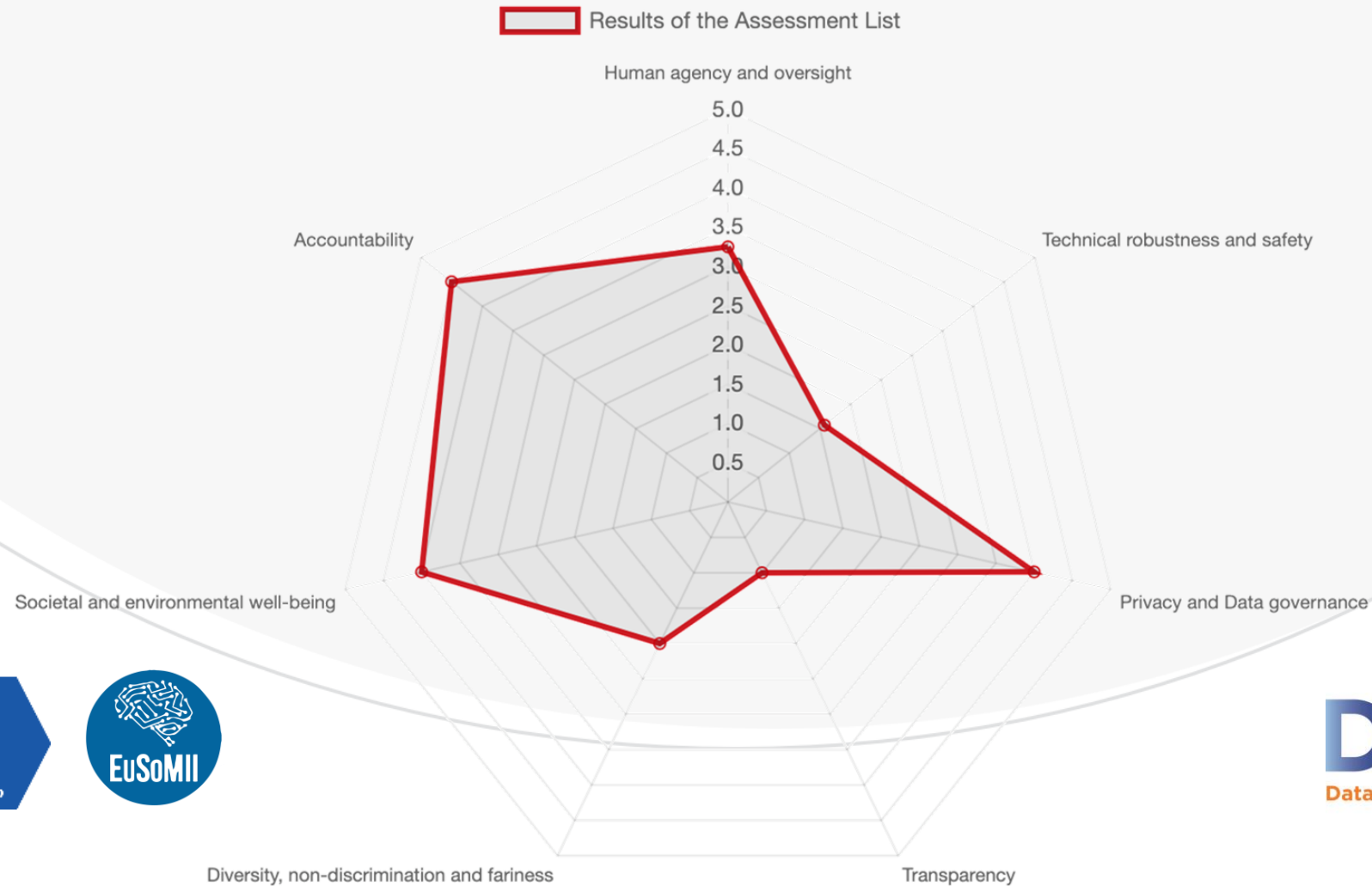
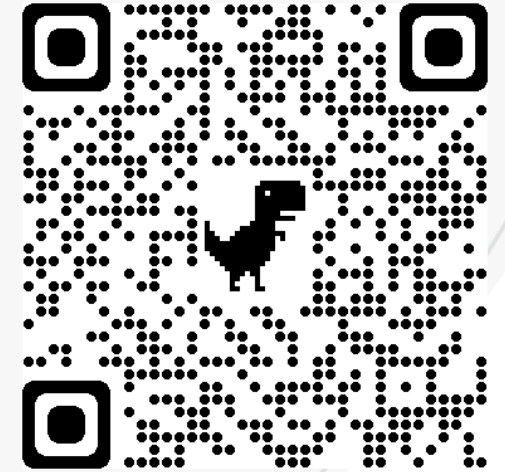
Human Autonomy

<https://altai.insight-centre.org/>

This subsection deals with the effect AI systems can have on human behaviour in the broadest sense. It deals with the effect of AI systems that are aimed at guiding, influencing or supporting humans in decision making processes, for example, algorithmic decision support systems, risk analysis/prediction systems (recommender systems, predictive policing, financial risk analysis, etc.). It also deals with the effect on human perception and expectation when confronted with AI systems that 'act' like humans. Finally, it deals with the effect of AI systems on human affection, trust and (in)dependence.

Is the **AI system** designed to interact, guide or take decisions by human **end-users** that affect humans ('**subjects**') or society? **?** *

Self assessment results and recommendations



Self assessment results and recommendations

Recommendations

Human agency and oversight

Put in place any procedure to avoid that the system inadvertently affects human autonomy.

Deploy a "stop button" or procedure to safely abort an operation when needed.

Technical robustness and safety

No recommendation for this requirement.

Privacy and Data Governance

Consider the privacy and data protection implications of data collected, generated or processed over the course of the AI system's lifecycle.

Whenever possible and relevant, align the AI-system with relevant standards (e.g. ISO, IEEE) or widely adopted protocols for (daily) data management and governance.



Hulpmiddel Handelingsruimte Waardevolle AI voor gezondheid en zorg

- Om onderzoekers en ontwikkelaars in het traject van ontwikkeling tot opschaling van waardevolle artificiële intelligentie (AI) te helpen, biedt dit hulpmiddel aanwijzingen in de handelingsruimte binnen de wet- en regelgeving. Zo kan er vroegtijdig gestart worden met het voorbereiden op gevraagde minimale eisen of standaarden. En reflecteren op acties om tot mensgerichte en betrouwbare AI-toepassingen te komen.



General Lesson

The future in medical imaging is in
Data Science Developments



General Lesson

The future in medical imaging is in Data Science Developments

When applied carefully with the aim to support and facilitate healthcare practice and continuously evaluated and adapted to current possibilities and needs. All with a strong emphasis on technical, clinical, ethical and legal validation.

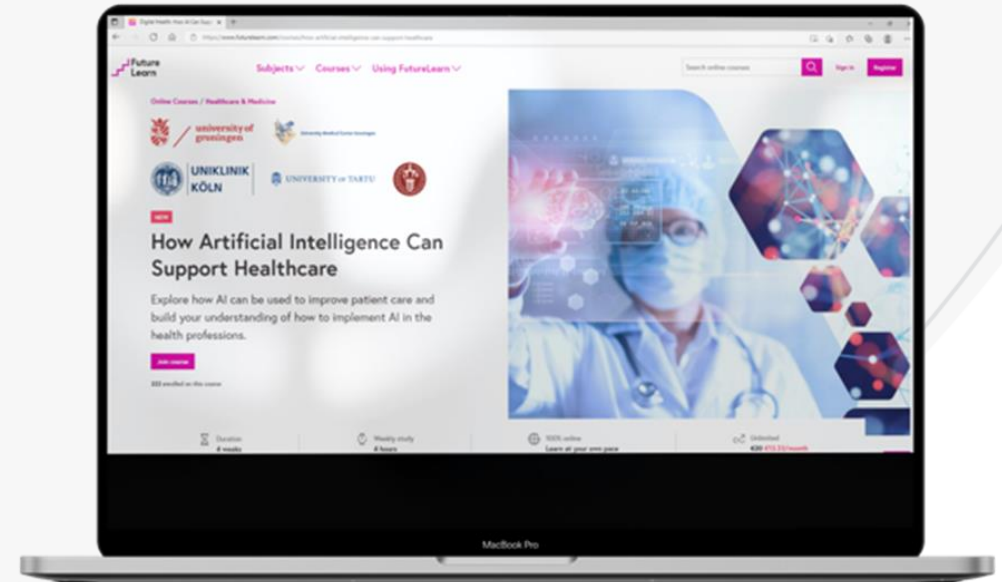
Want to learn more about the use of artificial intelligence (AI) in healthcare?

Follow our free online course

'How Artificial Intelligence can support Healthcare'

Discover how AI can be used to improve patient care and gain a deeper understanding of AI implementation in health professions.

<https://www.futurelearn.com/courses/how-artificial-intelligence-can-support-healthcare>



Contact us



linkedin.com/company/dash-umcg



umcgresearch.org/w/dash



dash@umcg.nl



twitter.com/DASH_umcg