

FGAI4H-O-047

Berlin, 31 May – 2 June 2022

Source: Institute for Molecular Medicine Finland – FIMM

Title: Workshop: TG-POC & TG-Histo - When is AI good enough for implementation in diagnostics?

Purpose: Discussion

Contact: Johan Lundin E-mail: johan.lundin@helsinki.fi
Institute for Molecular Medicine
Finland – FIMM
University of Helsinki, Finland
Karolinska Institutet, Sweden

Abstract: This PPT contains a presentation from the TG-POC & TG-Histo workshop on “Validation of annotations for AI models within the scope of point-of-care diagnostics (POC)”

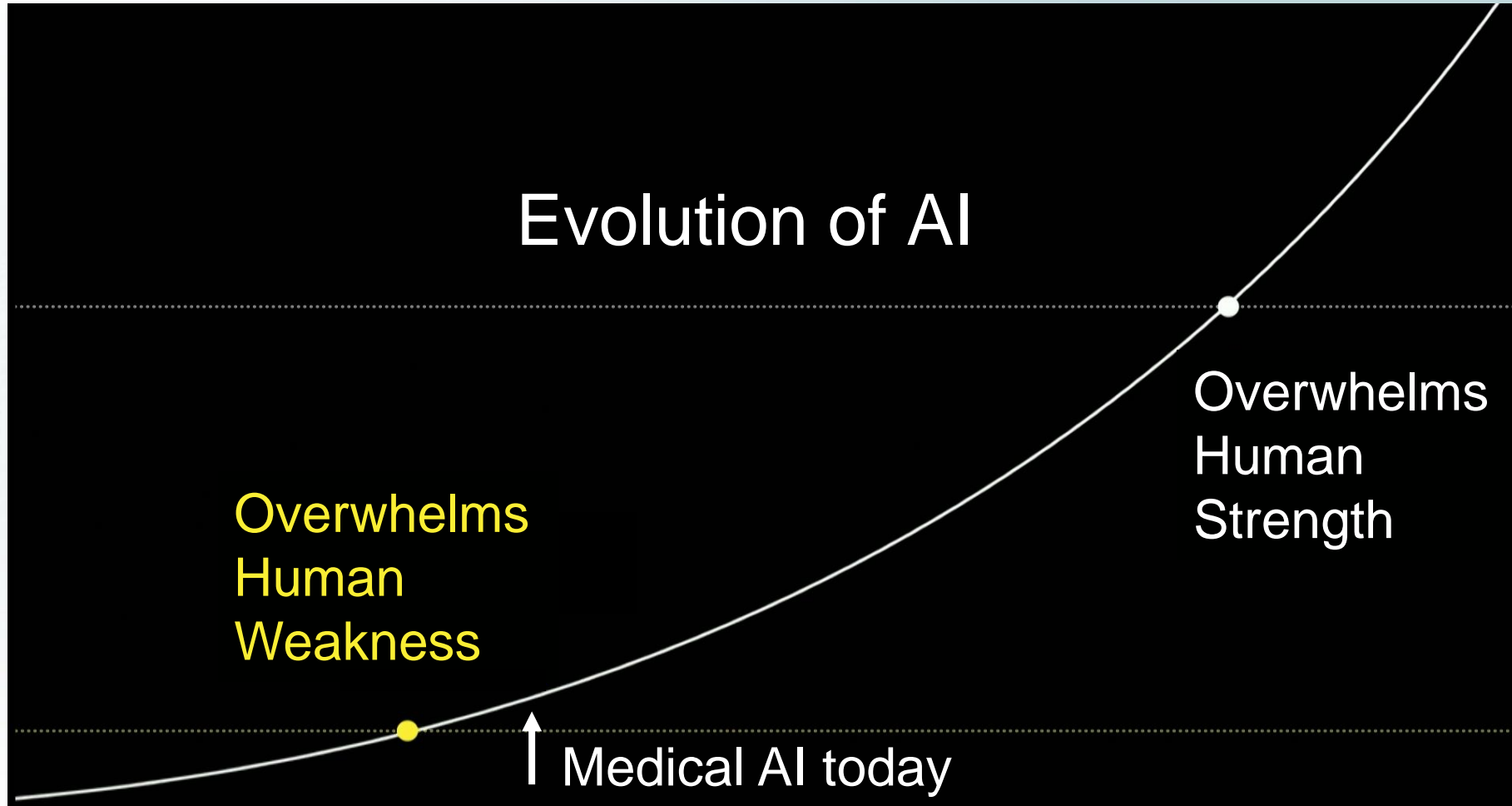
When is AI good enough for implementation in diagnostics?

Johan Lundin, MD, PhD
Research Director
Institute for Molecular Medicine Finland – FIMM
University of Helsinki, Finland
and
Professor of Medical Technology
Karolinska Institutet, Sweden

Disclaimer

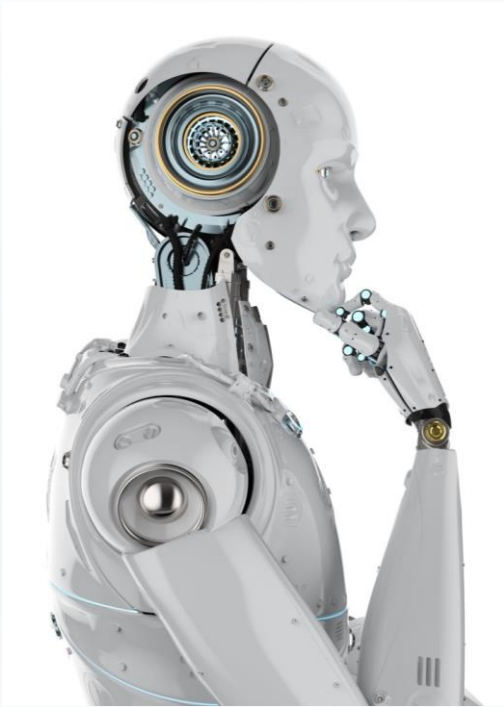
- › Founder, co-owner and board member of Aiforia Technologies



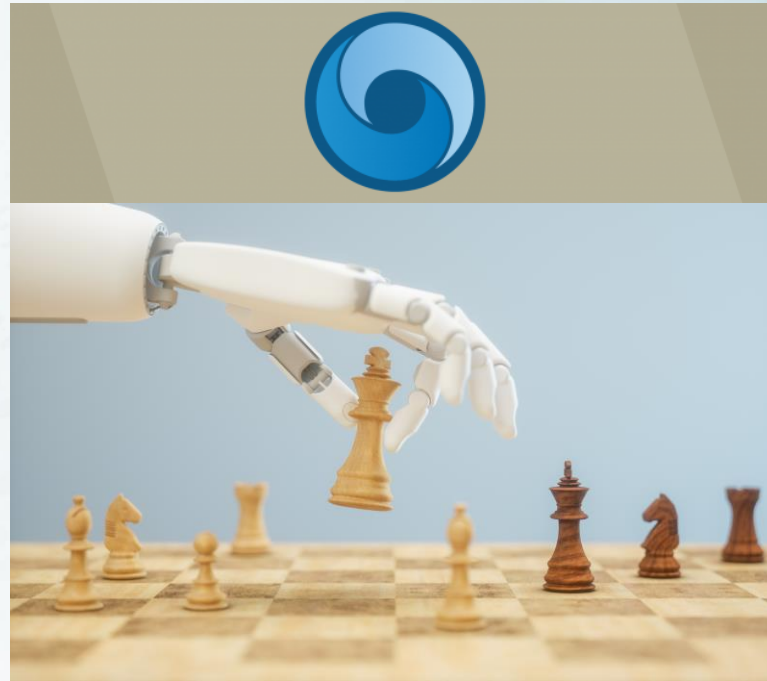


Courtesy of Tristan Harris, Center for Humane Technology

Human chess world champion learns from games played by AI



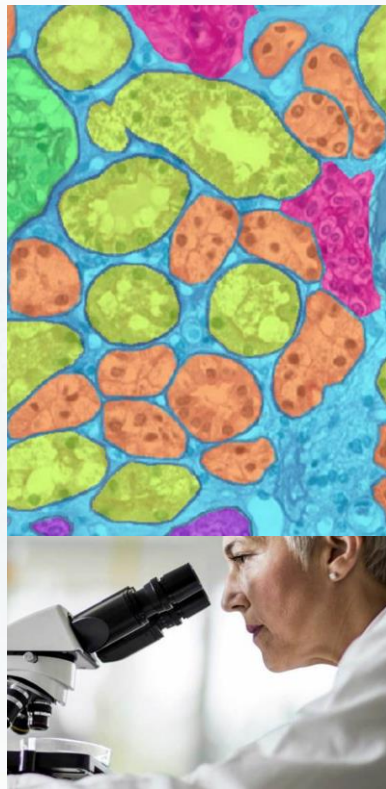
AlphaZero, self-taught AI-based world champion of chess*



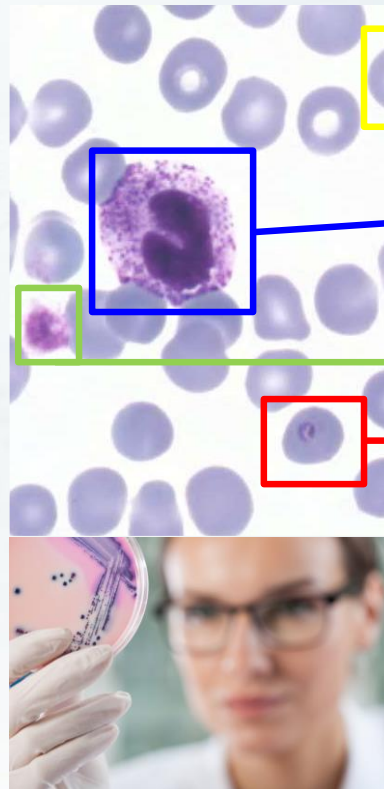
Magnus Carlsen, the current human champion of chess

Medical experts are likely to benefit from AI in a similar way


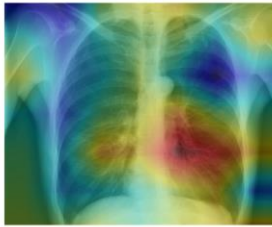
AI will impact all medical fields where an expert makes a visual interpretation



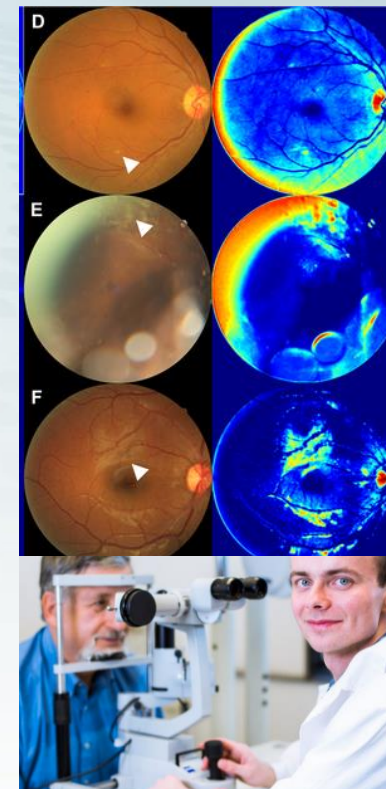
Pathology



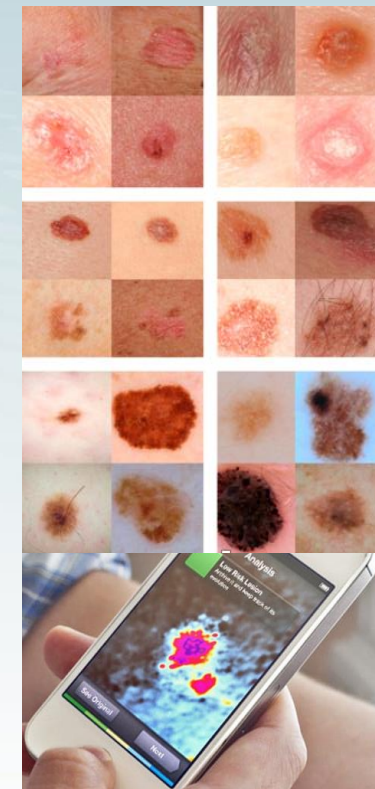
Microbiology


Input Chest X-Ray Image
CheXNet 121-layer CNN
Output Pneumonia Positive (85%)


Radiology



Ophthalmology



Dermatology

Crude estimate: minimum >7-8 billion visual diagnostic assessments globally per year

When is AI good enough for implementation in diagnostics?

- › When AI achieves the same accuracy or exceeds the current gold standard?
- › When AI analyses an order of magnitude more samples than a human expert within a particular time period?
- › When AI complements the human expert and finds a significant number of targets that otherwise would have been missed?
- › When AI is the only alternative in a setting with shortage of experts?

Ground truth and gold standard tests

- › The term ground truth refers to the underlying absolute state of information
- › The gold standard strives to represent the ground truth as closely as possible.
- › In machine learning and information retrieval, "ground truth" is the preferred term even when classifications may be imperfect
- › The gold standard is assumed to be the ground truth

Definition of a gold standard test

In medicine and statistics, a **gold standard test** is usually the diagnostic **test** or benchmark that is the best available under reasonable conditions. Other times, a **gold standard** is the most accurate **test** possible without restrictions.

[Gold standard \(test\) - Wikipedia](https://en.wikipedia.org/wiki/Gold_standard_(test))

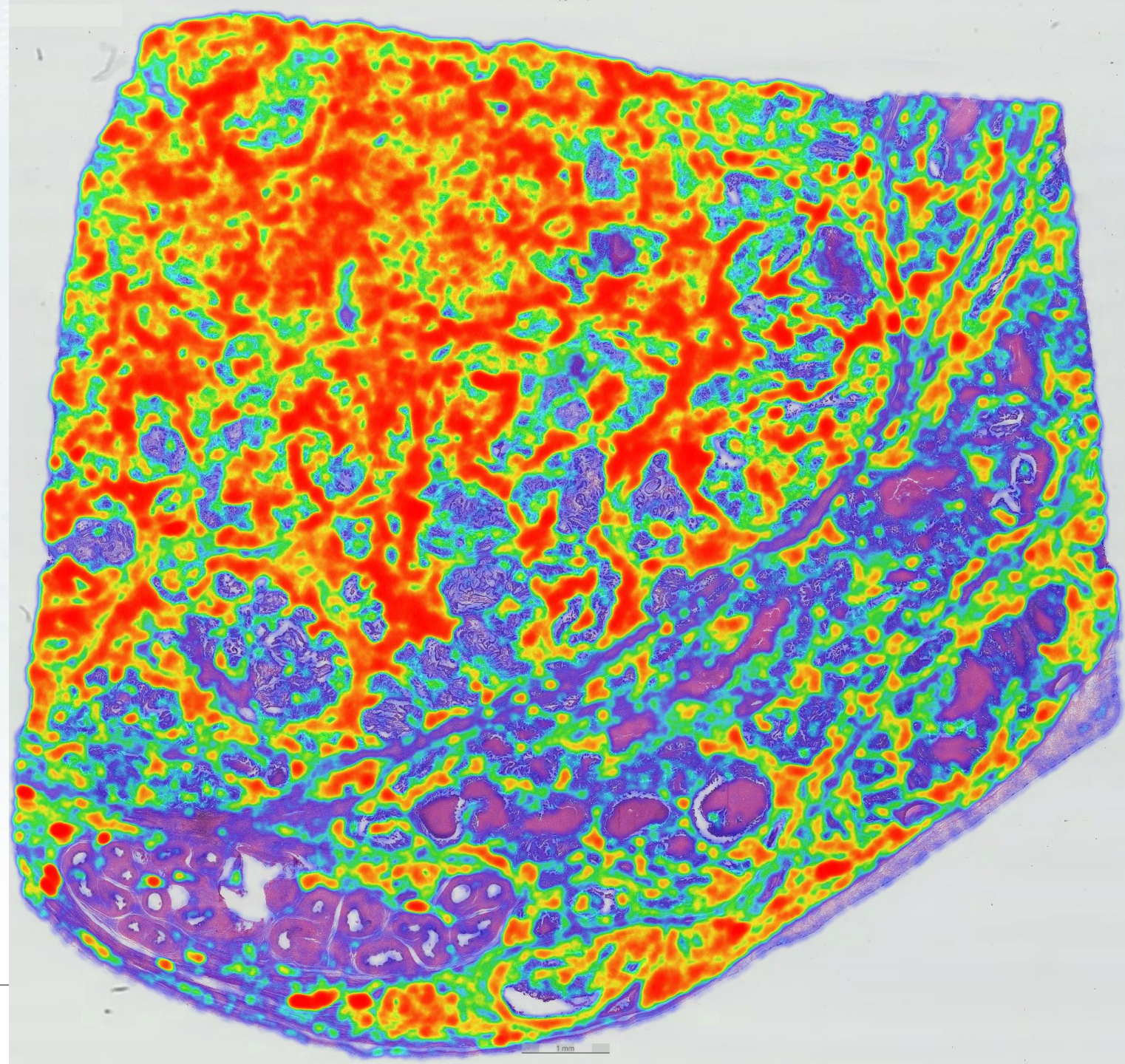
[https://en.wikipedia.org › wiki › Gold_standard_\(test\)](https://en.wikipedia.org/wiki/Gold_standard_(test))

Performance of a gold standard test

- › A hypothetical ideal "gold standard" test has a sensitivity of 100% with respect to the presence of the disease and a specificity of 100%.
- › In practice, there are sometimes no true gold standard tests.
- › Currently, no gold standard tests exist for deep learning applied to pathology or microscopy
- › According to the literature, AI-based algorithms typically reach a good to excellent diagnostic accuracy as compared to the ground truth, but the ground truth is rarely a real gold standard

Challenges in the development of gold standard for machine learning in pathology

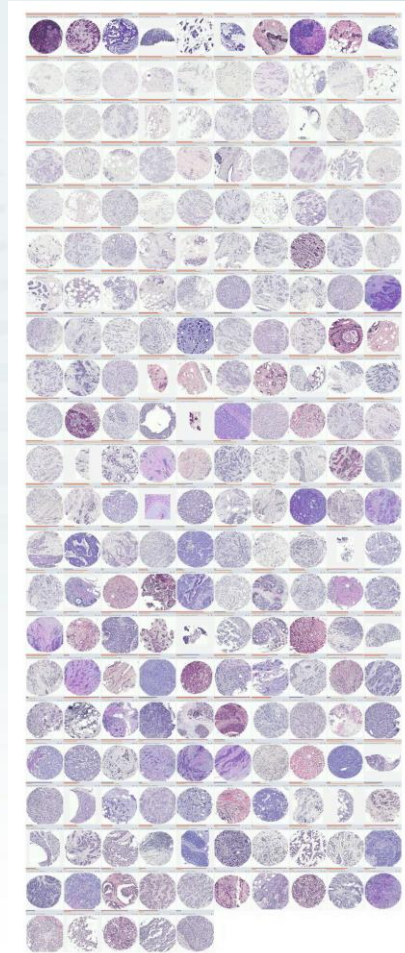
- › Samples and data
- › Annotations
- › Algorithms
- › Other challenges



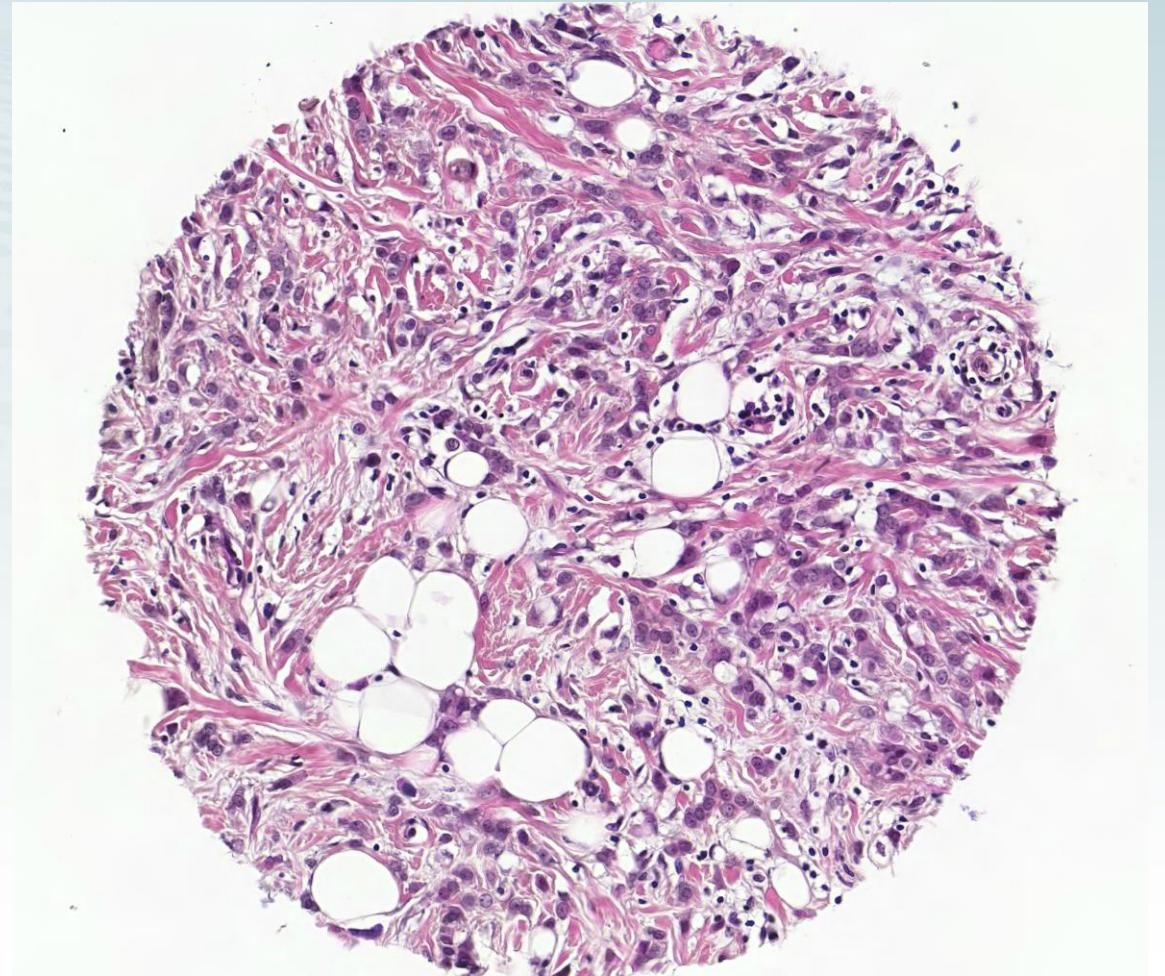
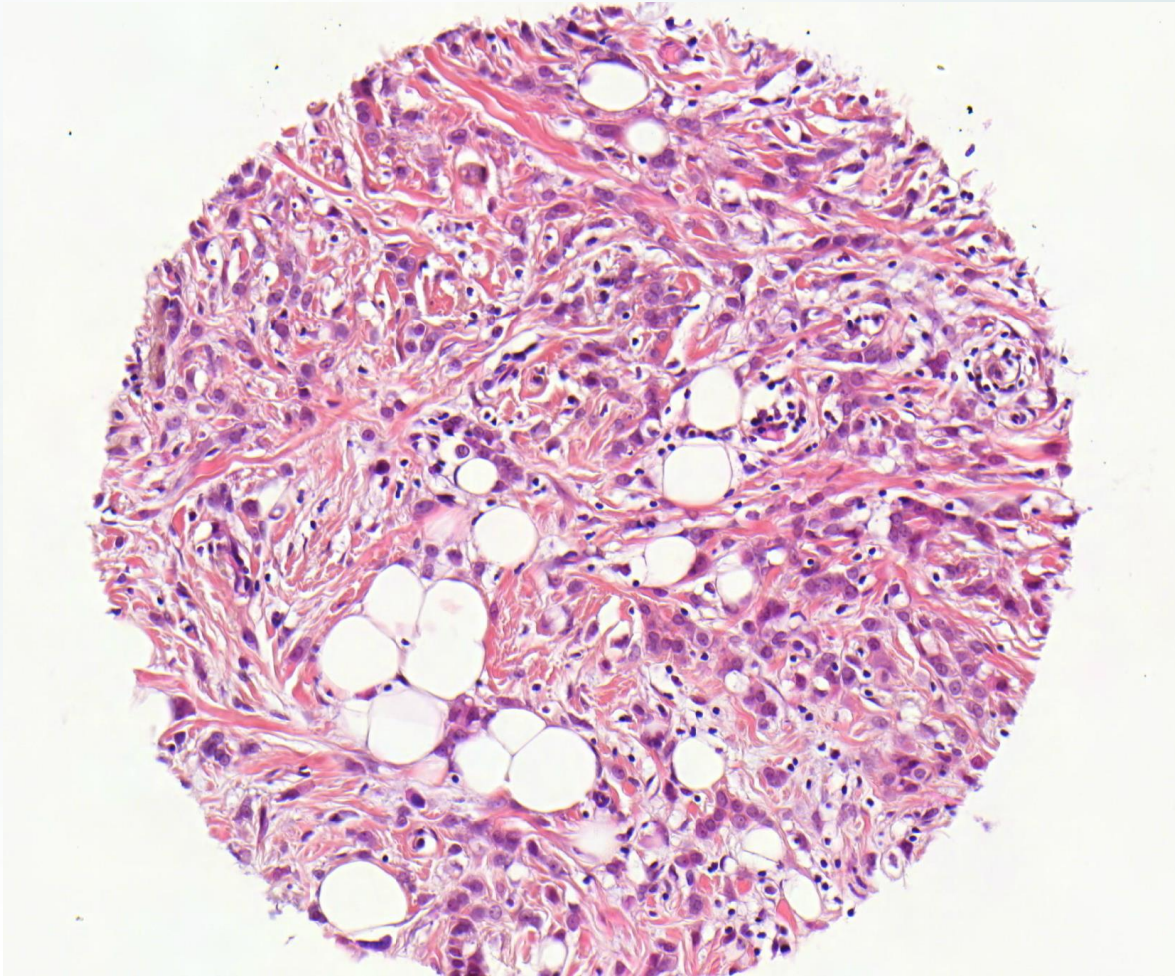
Challenges related to samples, data and annotations

- › Limited access to sample images with associated clinical data
 - Share data, form joint projects, use federated or swarm learning
- › Lack of annotated images
 - Create public libraries and common repositories of annotated images
- › Biased data due to incompleteness or lack of diversity
 - Strive for completeness of data, collect from many centers
- › Variable quality, artifacts and heterogeneity of samples
 - Perform quality control (with AI?), re-cut, re-stain, re-scan, color calibrate
 - ..or include all types of artifacts and variabilities in the training set

Sample variability in a breast cancer tissue microarray



Variability due to the scanner and camera



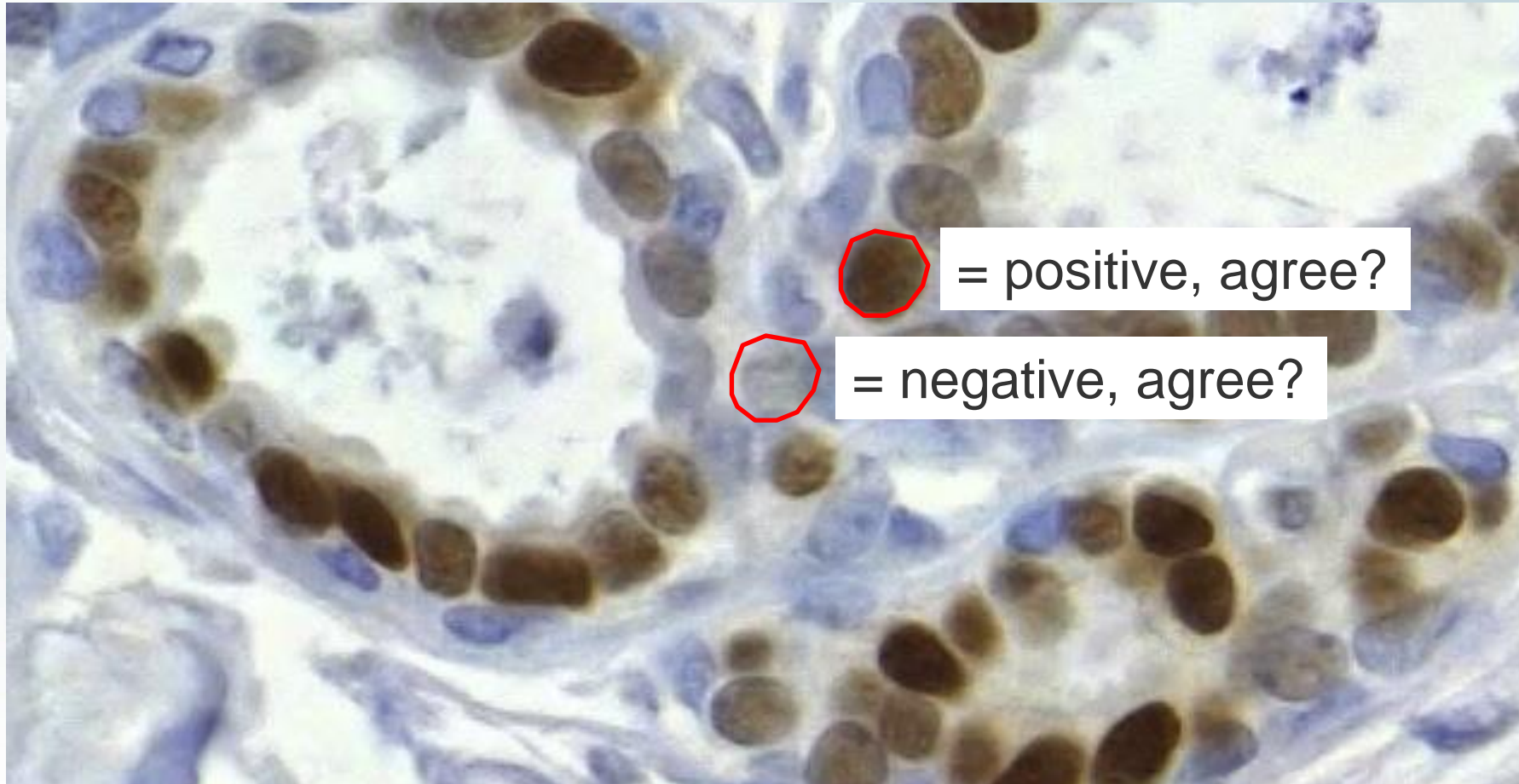
Variable quality of annotations

The screenshot displays the aiforia_create software interface. On the left is a dark sidebar with the 'aiforia_create' logo and a back arrow. The main interface is divided into several sections:

- Anti-CD8**: The current analysis type.
- Annotations**: A list of annotations, currently showing two 'CD8 positive region' entries.
- LAYERS**: Includes '+ ADD LAYER' and 'EDIT LAYER' buttons. A layer named 'CD8 positive region' is selected.
- IMAGES**: A dropdown menu showing '94_32_ACD_ICD_089_05_A_SLS_CD8' and a navigation indicator '4 / 44'.
- TRAINING REGIONS / ANALYSIS REGIONS**: A section with icons for creating, editing, and deleting regions. It lists 'Training region 4' and 'Training region 3'.
- ANNOTATIONS (4)**: A list of annotations with edit and delete icons.

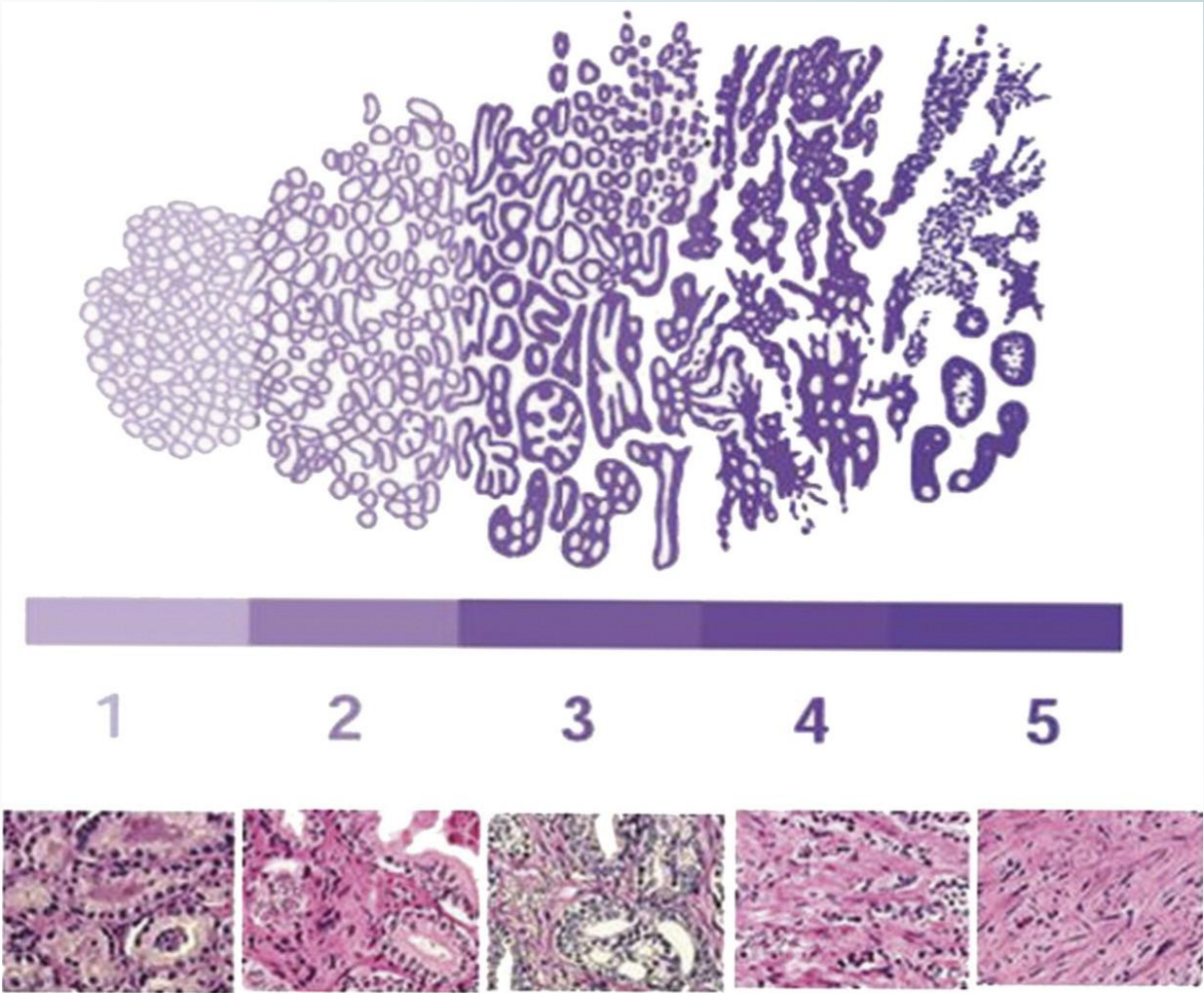
The central area shows a histology image of a tissue section. A large purple-stained area is visible, with a smaller region at the bottom outlined in black and containing several red dots. A scale bar at the bottom indicates '1 mm'. A zoomed-in inset of the red-dotted region is shown in the bottom right corner.

Consistency and representativeness of ground truth -deep learning algorithms just as good as their teacher?

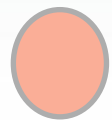
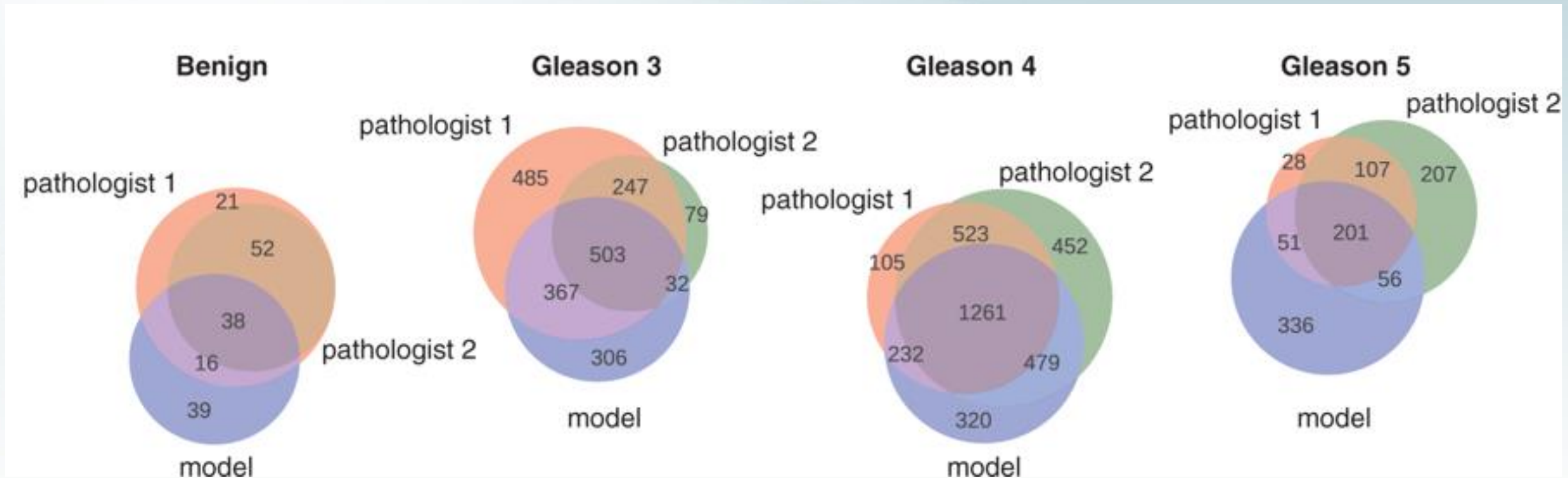


Ground truth digital samples crucial in development of medical AI

More complicated patterns – e.g. Gleason grade in prostate cancer



There is a risk that AI will be just another subjective "expert" - example of automated grading of Gleason



Pathologist 1



Pathologist 2



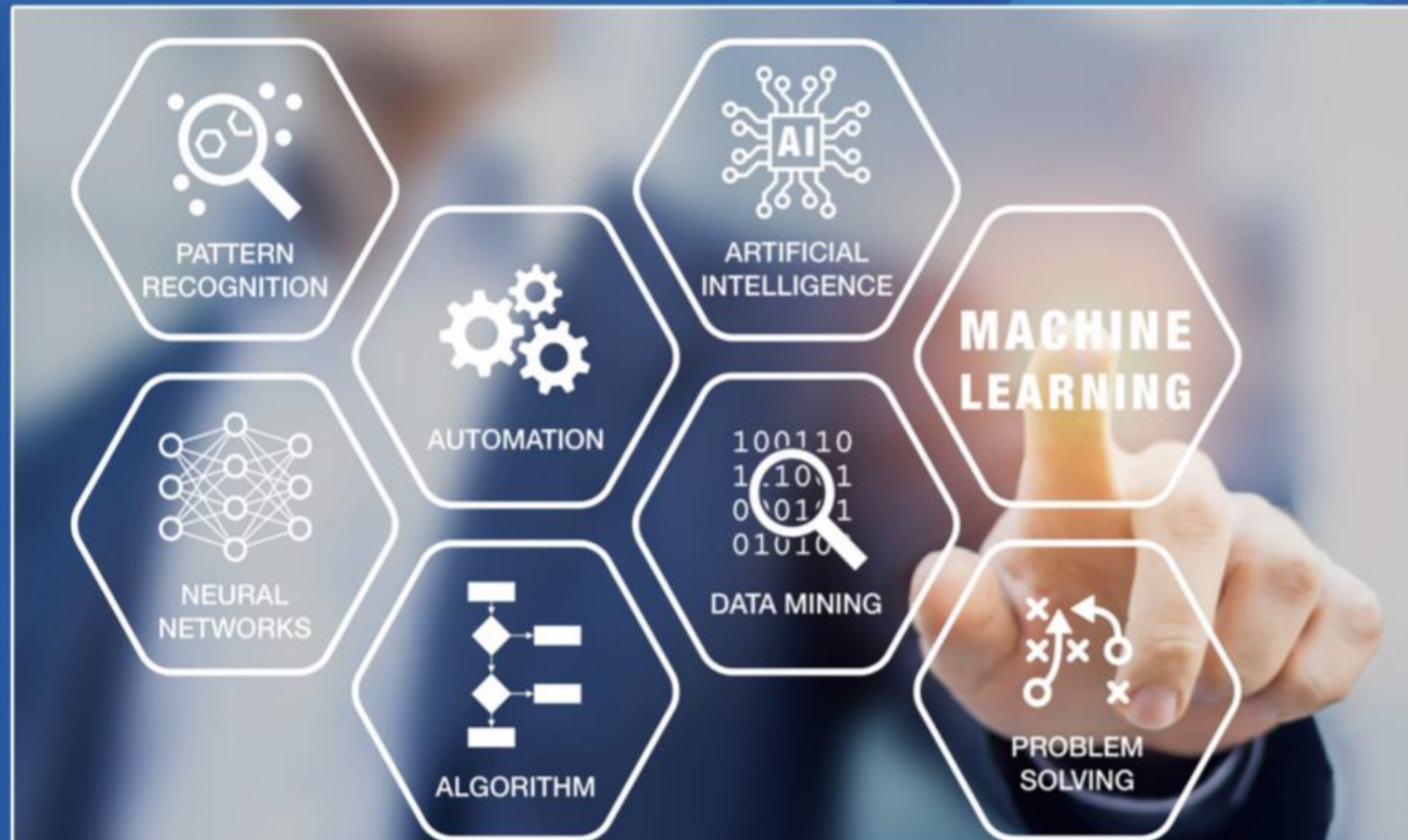
Model = AI

Other challenges related to AI for diagnostics

- › Ownership and access to the images
 - Create public libraries of annotated images for developers
- › What to do with rapidly improving and updated algorithms?
 - Allow algorithm performance to be a moving target
 - FDA white paper
- › How to handle tens or hundreds of AI:s for the same purpose?
 - Create consensus algorithms? Use swarm learning to adjust parameters?

Proposed Regulatory Framework for Modifications to Artificial Intelligence/Machine Learning (AI/ML)-Based Software as a Medical Device (SaMD)

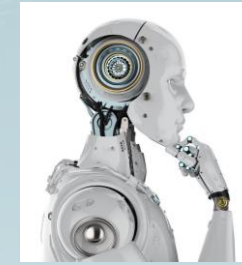
Discussion Paper and Request for Feedback



Human vs machine



Human



Machine

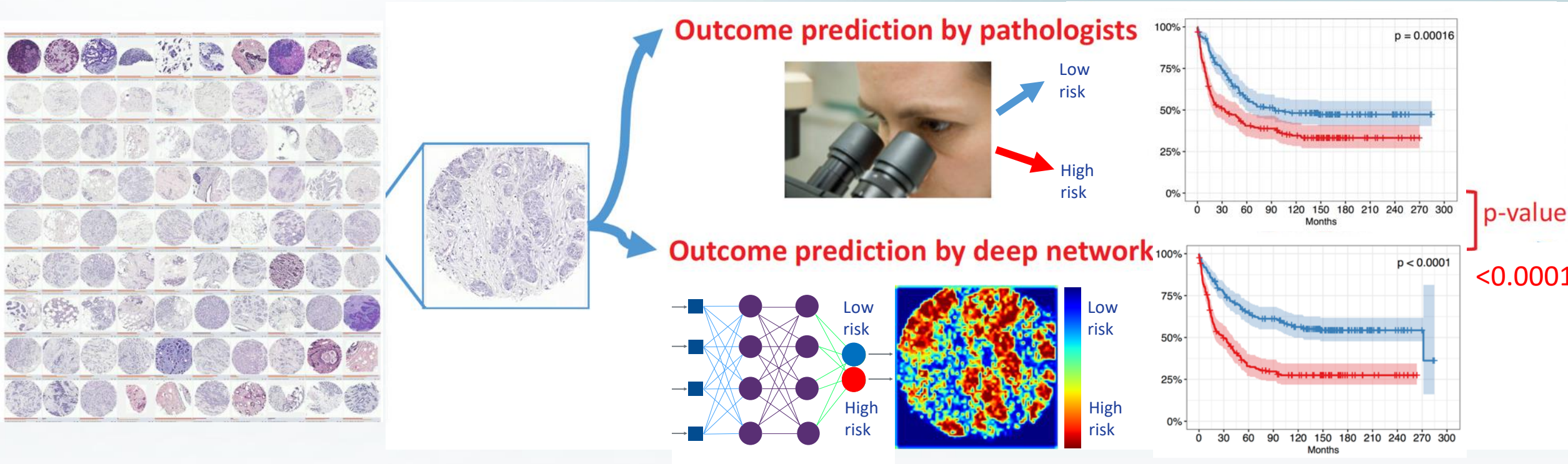
- › Understand context
- › Reproduce assessment
- › Handle outliers
- › Find rare events
- › Generalize
- › Achieve high throughput



When is an AI algorithm good enough?

- › 🤖?
- › Needs to outperform or supplement human experts in at least one of the following: Sensitivity, speed, reproducibility
- › Generalizability needs to be established
 - robustness to artifacts, outliers and local variations in protocols
- › Superiority shown in prediction of clinical endpoints rather the replication of annotations?
 - Outcome and biomarker supervised learning

Outcome supervised learning in colorectal and breast cancer



Digitized tissue samples from cancer patients with known outcome of cancer i.e. survivor or non-survivor

Comparison of human expert-based and AI-based outcome prediction

Deep learning outperformed experienced pathologists in outcome prediction ¹⁻³

¹Bychkov et al, Scientific Reports 2018;8:3395 ²Turkki et al, Breast Can Res Tr 2019;177:41-52 ³Bychkov et al, J Pathol Informatics 2022;13:9

Human and machine combination: Sensitive AI algorithm – specific human observer

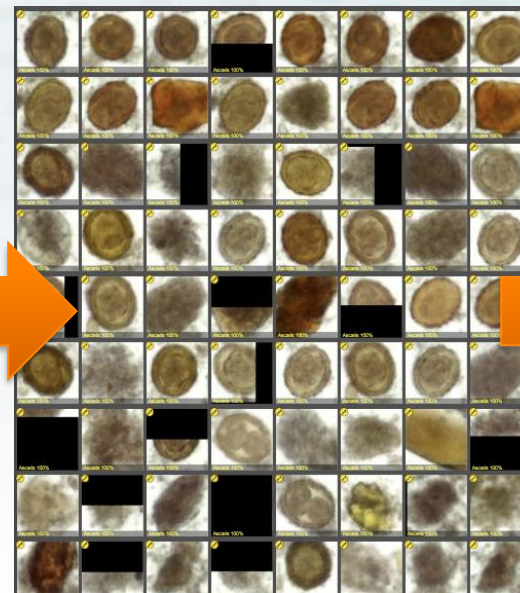
- › Example of application to neglected tropical diseases for better access to diagnostics
- › Assisted detection of helminth eggs in stool samples and verification by human expert



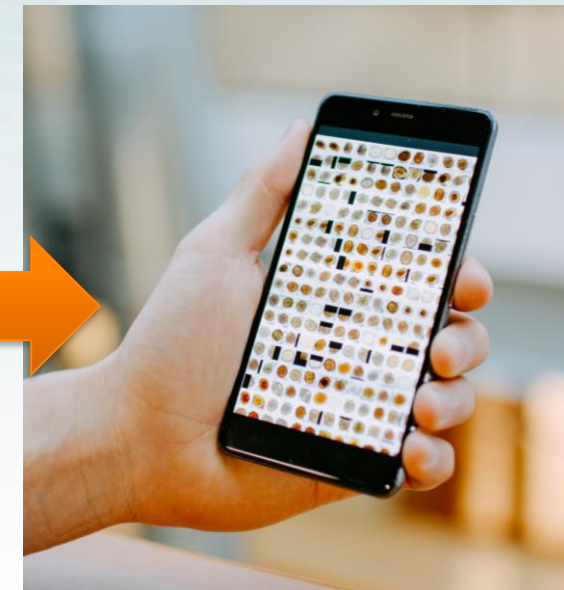
Stool sample



Scan with mobile microscope



Parasite egg candidates

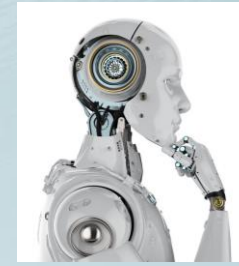


Human expert reviews results

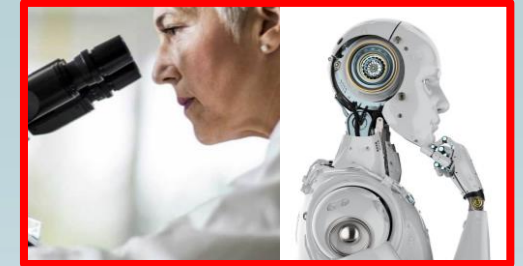
Human **AND** machine?



Human



Machine



Human&Machine

- › Understand context
- › Reproduce assessment
- › Handle outliers
- › Find rare events
- › Generalize
- › Achieve high throughput



Some arguments and thoughts as a conclusion

- › The quality, selection and annotation of training data for AI is crucial
- › How could we get more reliable ground truth?
- › There is a risk that we end up with hundreds of AIs for the same purpose but with unclear accuracy
- › Ground truth based on human observation and annotation will always be subjective
- › Alternative ground truths, such as sample level diagnosis, patient outcome and response to treatment should be explored

nature reviews cancer

[Explore content](#) ▾

[About the journal](#) ▾

[Publish with us](#) ▾

[nature](#) > [nature reviews cancer](#) > [viewpoint](#) > [article](#)

Viewpoint | [Published: 17 September 2021](#)

Artificial intelligence in cancer research, diagnosis and therapy

[Olivier Elemento](#) ✉, [Christina Leslie](#) ✉, [Johan Lundin](#) ✉ & [Georgia Tourassi](#) ✉

[Nature Reviews Cancer](#) (2021) | [Cite this article](#)

2977 Accesses | **110** Altmetric | [Metrics](#)

Lundin Group, FIMM



Klaus Kessel

Sebastian Stenman

Oscar Holmström

Johan Lundin

Hakan Kucukel

Mikael Lundin

Nina Linder

Dmitrii Bychkov

Antti Suutala





Institute for Molecular Medicine Finland
Nordic EMBL Partnership for Molecular Medicine

Building a bridge from discovery to medicine