# Challenges for machine learning using medical data

Henning Müller

AIDA lecture

# Henning Müller

- **Medical informatics** studies in Heidelberg, Germany

  Exchange with Daimler Benz research, USA

- PhD in **CBIR**, computer vision, Geneva, Switzerland (1998-2002)

  - Exchange with Monash University, Melbourne, AUS

- Professor in radiology and medical informatics at the University of Geneva (2014-)

- Professor in Computer Science at the HES-SO, Sierre, Switzerland (2007-)
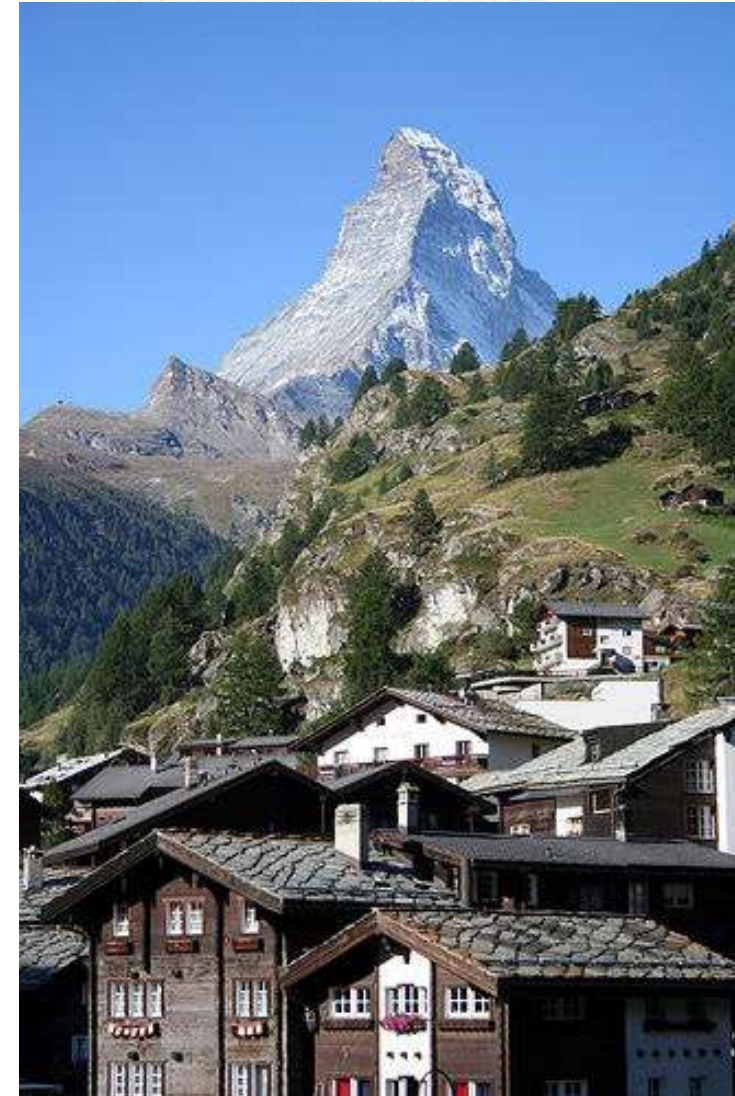
  - Visiting faculty at Martinos Center (2015-2016)

- Member of the Swiss National Research Council

# Where we are

# Overview

- Status of medical AI

  - And its <span style="color:red">challenges</span>

- <span style="color:red">Projects</span> addressing the challenges

  - With a bias towards our work

- Open challenges

- Conclusions and <span style="color:red">discussion</span>
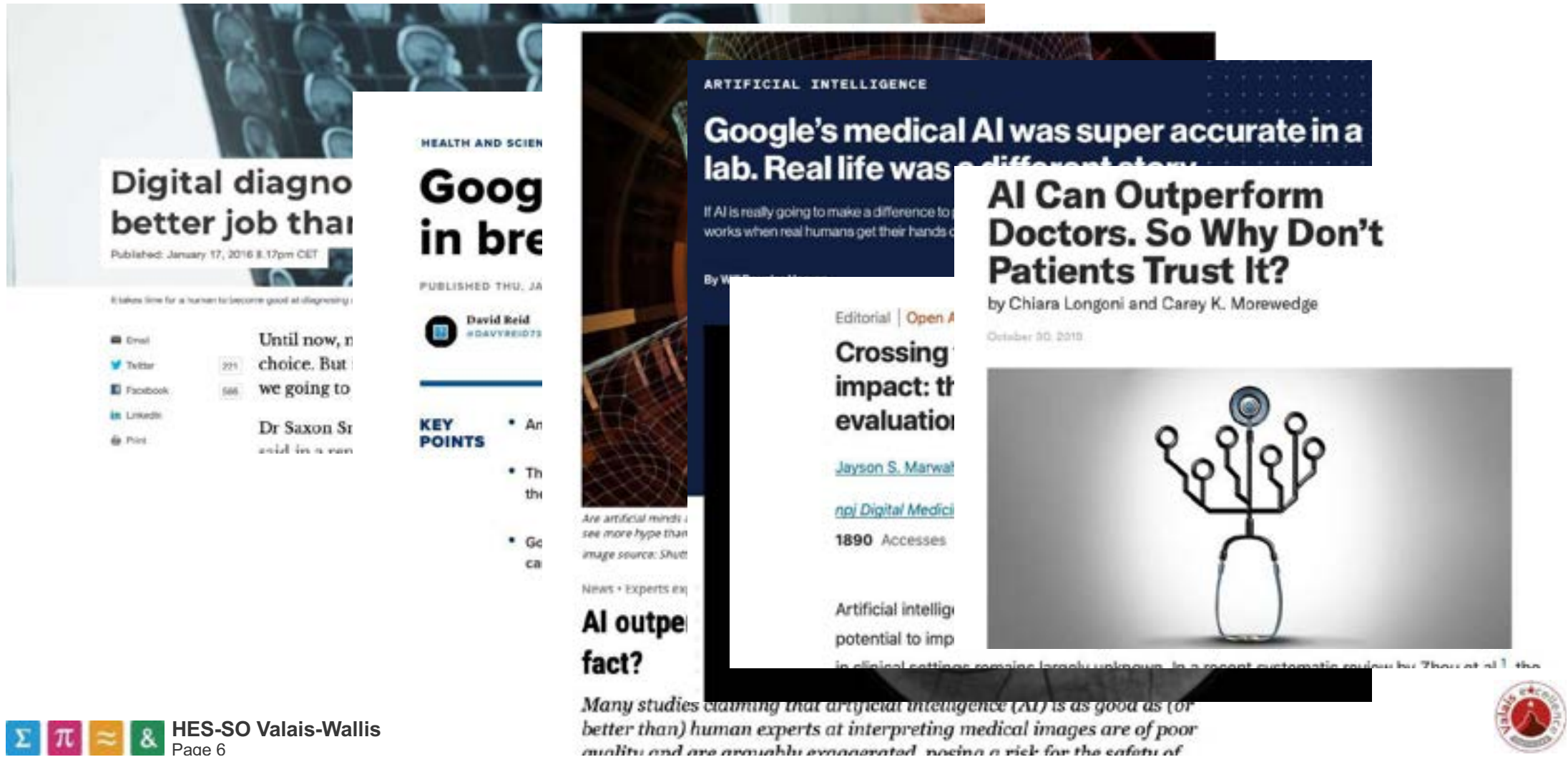
# The promise of medical AI

Geoff Hinton: On Radiology

https://www.youtube.com/watch?v=2HMPRXstSvQ

# Medical AI in the media

# Realistic expectations

- <https://www.aimyths.org/>

# Advantages of medical data

- Images created under <span style="color:red">standardized</span> conditions

- Images are always attached to a case and a <span style="color:red">report</span> describing it, plus a reason for producing the images

  - <span style="color:red">Metadata</span> exist, and other data on the same patient

    - We know the <span style="color:red">context</span> of the images

- Much medical <span style="color:red">knowledge</span> is available

  - Coded and maintained in ontologies

- Much clinical research is done

  - Medical imaging is estimated to occupy 30% of world storage

# Challenges with medical data

- Data privacy and ethics can make sharing data hard

- Medical equipment and procedures vary across hospitals

  - And equipment changes frequently

- Pixel-level annotations are extremely expensive

  - Specialists are needed and often not available (too busy)

- An image is only a very small part in a case with patient history, temporal data, genetics, text, structured data etc.

  - Regions determining a decision are often extremely small

    • Needle in a haystack

# Challenges for medical AI

- Much data are needed, so solutions need to be scalable

  - Diversity is required for generalization

  - Data sets are very unbalanced

- Continuous learning is required due to changing equipment and clinical guidelines (half-life of knowledge)

- Pixel level annotations are not available, as expensive

- Combining multiple sources is needed for proper learning

- Results need to be explainable for workflow integration

  - Deep learning is a priori a black box

# Examode consortium



High Performance Computing (HPC) resources: SURFSARA

**ACADEMIC PARTNERS**
(UNIPD, HESSO)

Scientific experience in extracting knowledge from:
- heterogeneous images
- text

**MEDICAL PARTNERS:**
(AOEC, Radboudumc)

Worldwide unique providers of:
- medical data (imaging and text)
- digital pathology knowledge
- clinical evaluation experience

**INDUSTRIAL PARTNERS:**
(ONTOTEXT, MICROSCOPEIT)

Solid industrial experience and market opportunities in:
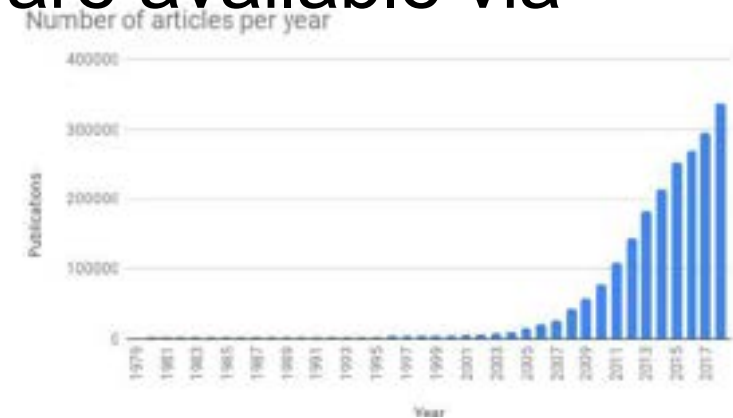- semantics-based services
- advanced AI-based image processing solutions

# Image accessibility

- **Open data** policies of funding agencies make large medical data sets available

  - Particularly NIH is pushing towards this

- **TCIA** and **TCGA** are very large repositories

  - There are many scientific challenges

- Images from the Biomedical literature are available via **PubMedCentral**

  - Exponentially increasing
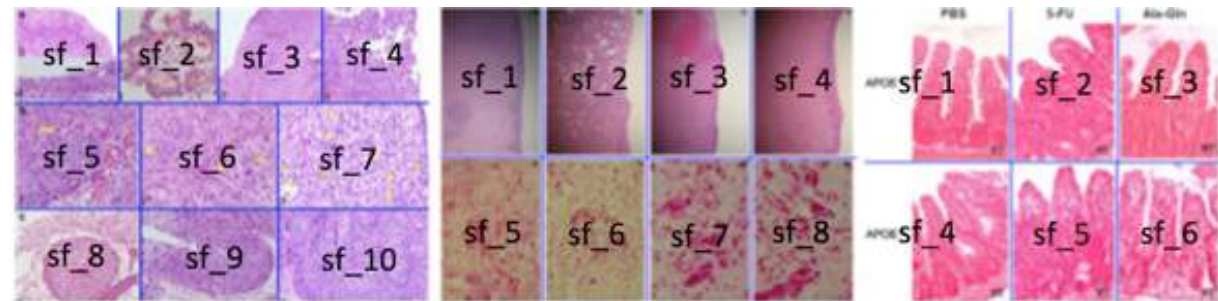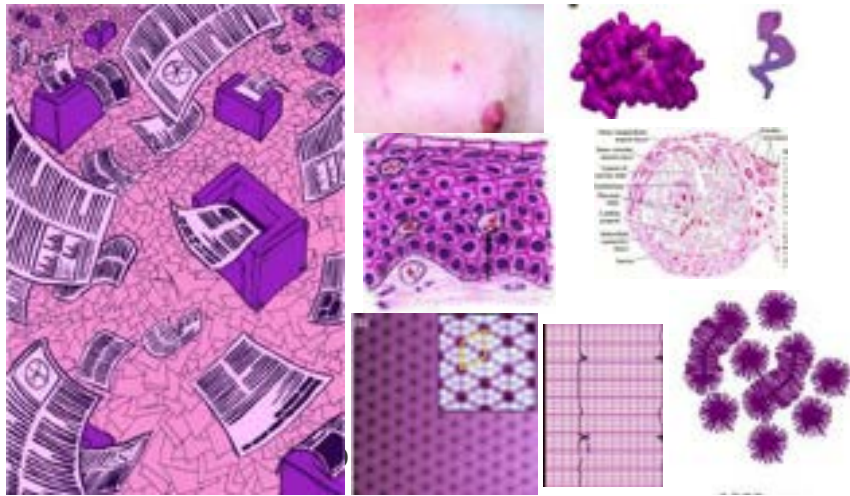
  - Extremely varied, hard to use

# Unbalanced data sets

- Differing <span style="color:red">frequencies</span> of the relevant classes need to be taken into account

    - At cancer screen even high-risk people are ~1% positive

    - <span style="color:red">Sensitivity</span> and <span style="color:red">specificity</span> as measures, not accuracy

        • Weight between false positives and false negatives varies

    - Some cases may occur once/twice per year in large hospitals

- <span style="color:red">Rare cases</span> is what is more commonly described in articles

    - Images from articles can thus help (at least in theory)

    - Variety of imaging parameters and laboratories is very high

# Challenges with PubMed

- >20'000'000 images in 2022, many graphs, charts

- Look-alikes is a problem, and compound figures

  - Very varied and sometimes strange content needs removal

- Compound figures need to be separated

  - Cutting sub figures apart makes content accessible

# Making the images usable

- Removing very small images & strange aspect ratios

- Classify figures into <span style="color:red">figure types</span>

  - Using image data and also text, remove non-relevant images

- Detect and cut <span style="color:red">compound figures</span> into their parts

  - Classify these into figure types again

- Filter <span style="color:red">human</span> vs. animal tissue and specific <span style="color:red">organs</span>

- Check <span style="color:red">diseases</span> or grading/staging images

  - Classes for machine learning
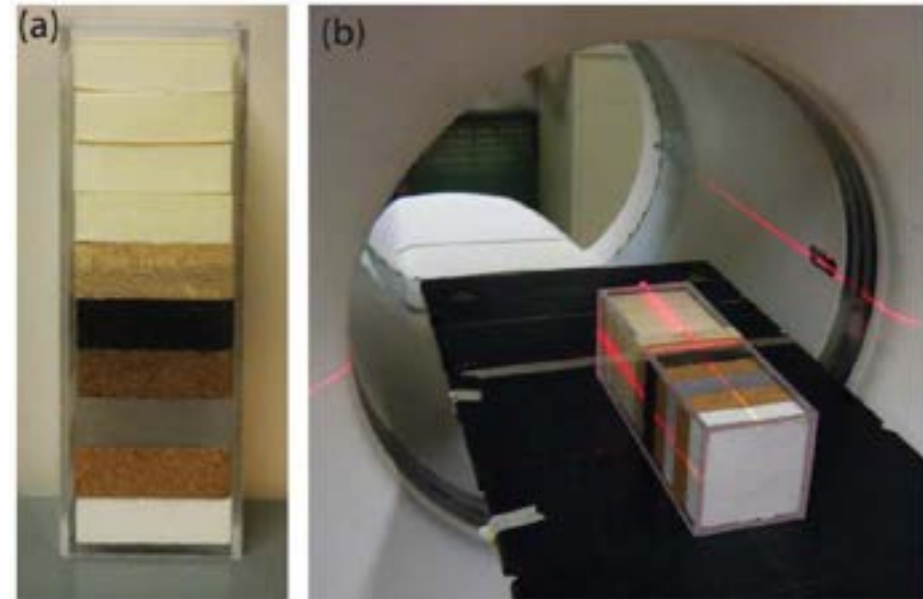
# Advantages of literature images

- **Rare images** (unusual, untypical) are generally used for articles and case descriptions

  - A few typical cases but mainly extreme cases

  - Creates critical mass for rare diseases

- Images are from **many laboratories** and thus contain many image variations (staining, scanners)

  - Increase generalizability of learned models thanks to this diversity

- Exponentially **increasing** content

# Image harmonization for radiomics

- **Different scanners produce different images**

  - Many protocols, construction kernels, producers, voxel sizes, …
    - Strong influence on features extracted

- **How can we harmonize this?**

  - Deep learning!

- **Phantom study with 17 scanners**

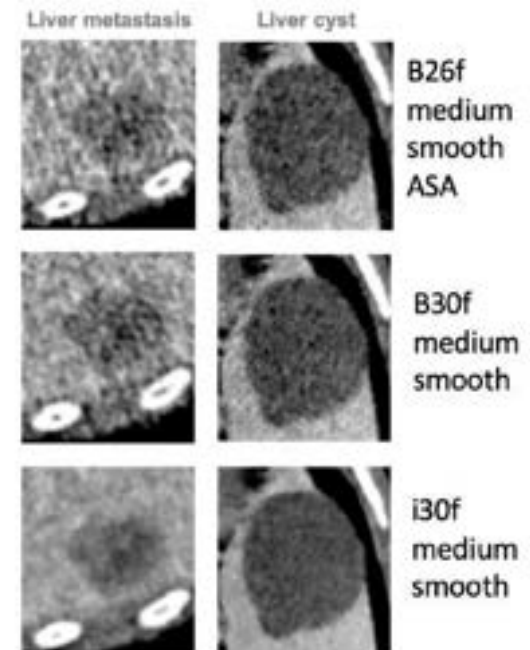  - 10 solid textures

  - Features invariant to scanner

Vincent Andrearczyk, Adrien Depeursinge, and Henning Müller, Neural Network Training for Cross-Protocol Radiomic Feature Standardization in Computed Tomography, Journal of Medical Imaging, 2019.

# Measuring CT variability

- Many CT parameter variations stemming from:

  - Acquisition protocols (radiation dose, …)

  - Image reconstruction parameters

  - Image resolution (slice thickness, overlap, …)

- Variability has a strong influence on the analysis & comparison of radiomics features

- Patient studies evaluating image/feature stability entail ethical concerns with multiple exposures to radiation
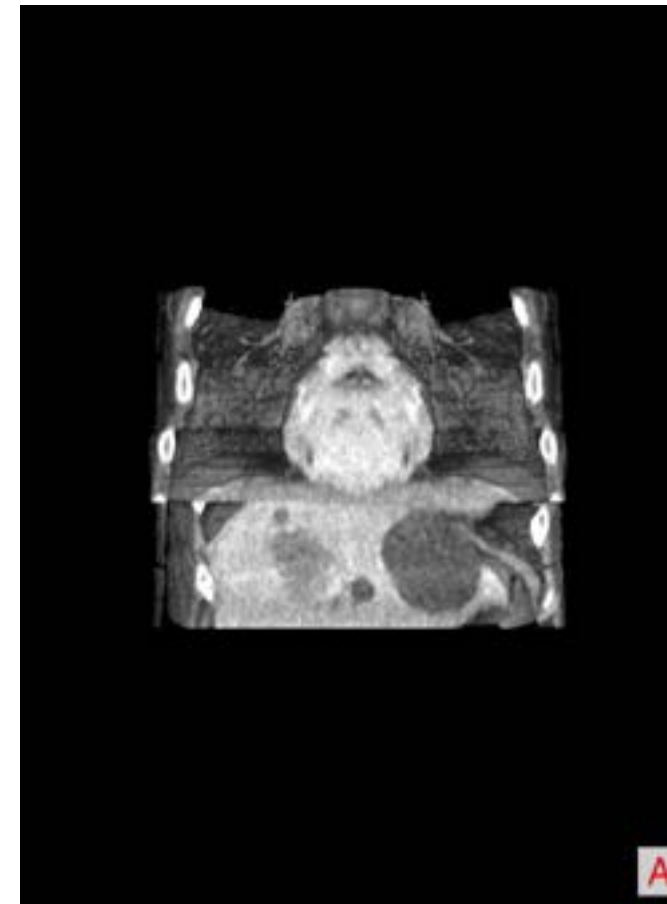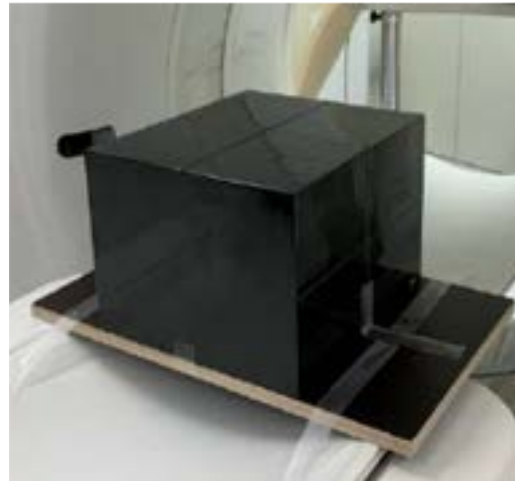
Traverso et al. (2018) *Repeatability and reproducibility of radiomic features: a systematic review.* International Journal of Radiation Oncology Physics **102**.

*Solomon et al. (2014) Quantum noise properties of CT images with anatomical textured backgrounds across reconstruction algorithms: FBP and SAFIRE. Medical Physics 41, 091908.*
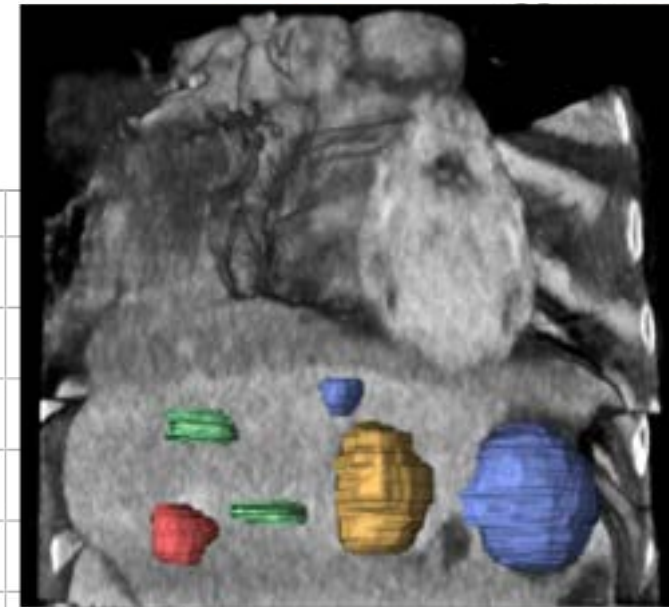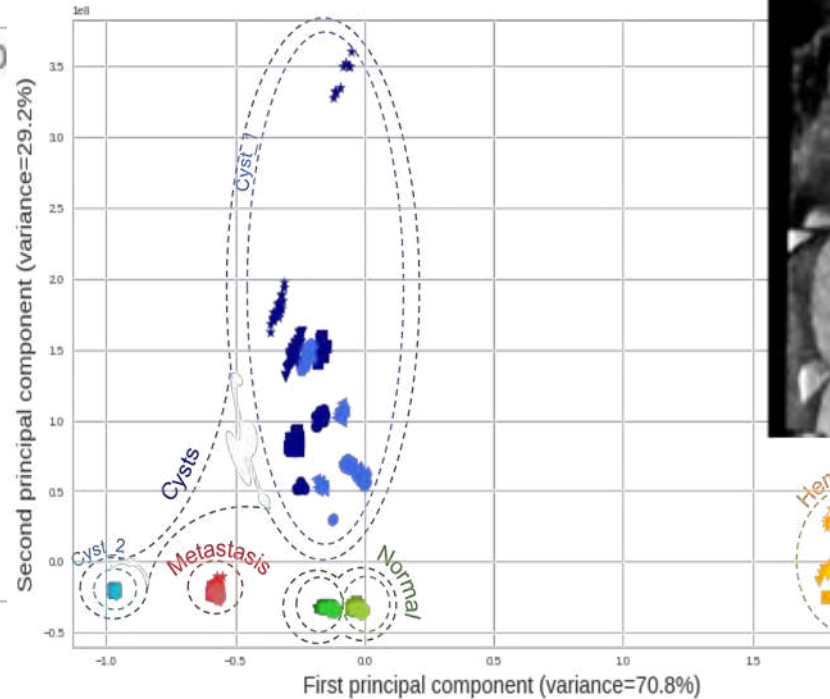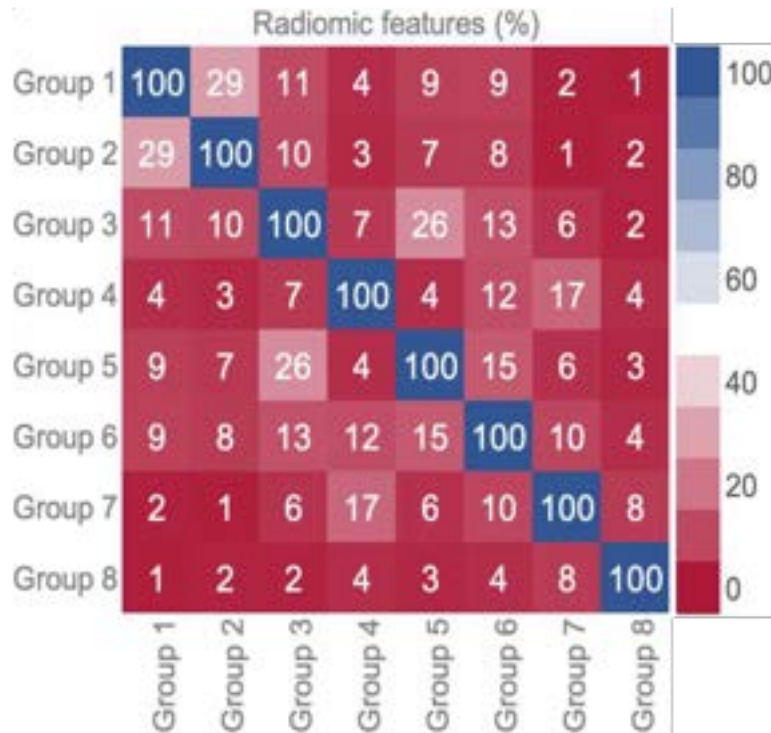
# 3D printed phantom

- Phantoms allow repeated radiation exposure

- Highly controlled acquisitions

  - No patient movement

  - No breathing

  - Precise positioning

- Limitations

  - In density (-100 HU to1000 HU)

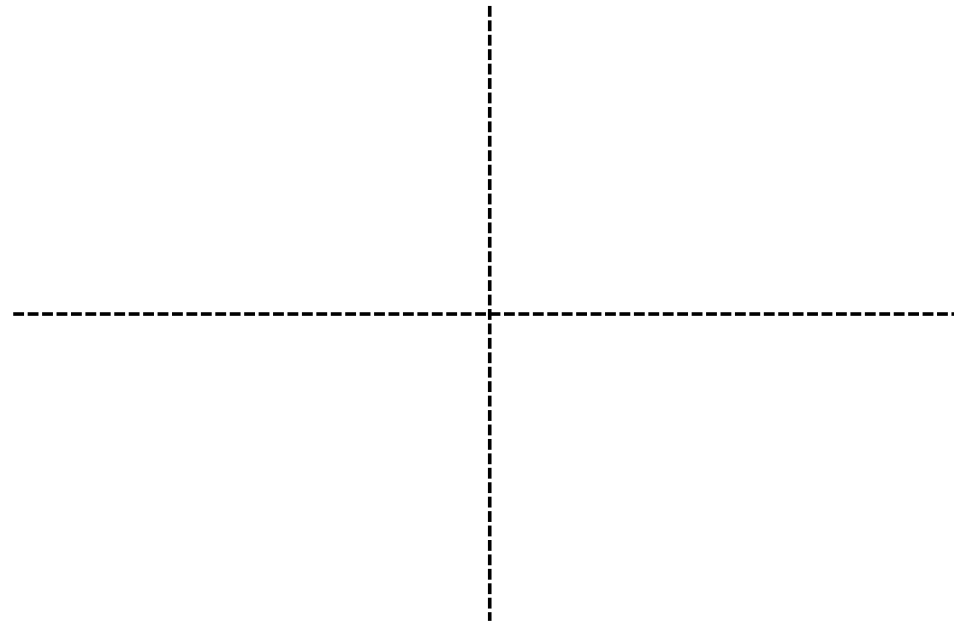  - Small blocks are glued (artifacts)

# First results



Oscar Jimenez-del-Toro, Christoph Aberle, Michael Bach, Roger Schaer, Markus Obmann, Kyriakos, Ender Konukoglu, Bram Stieltjes, Henning Müller, Adrien Depeursinge, The discriminative power and reproducibility of radiomics features with CT variations: Task-based analysis in a realistic CT liver phantom, Investigative Radiology, 2021.

# Stability vs. discriminative power

- **Image Biomarker Standardization Initiative**

  - Define all visual features used in radiomics

    - And compare implementations on the same data

      Digital phantoms

    - When there are differences then check the implementations

  - **Installment 1** is finished,

    - Simple statistical (texture) features

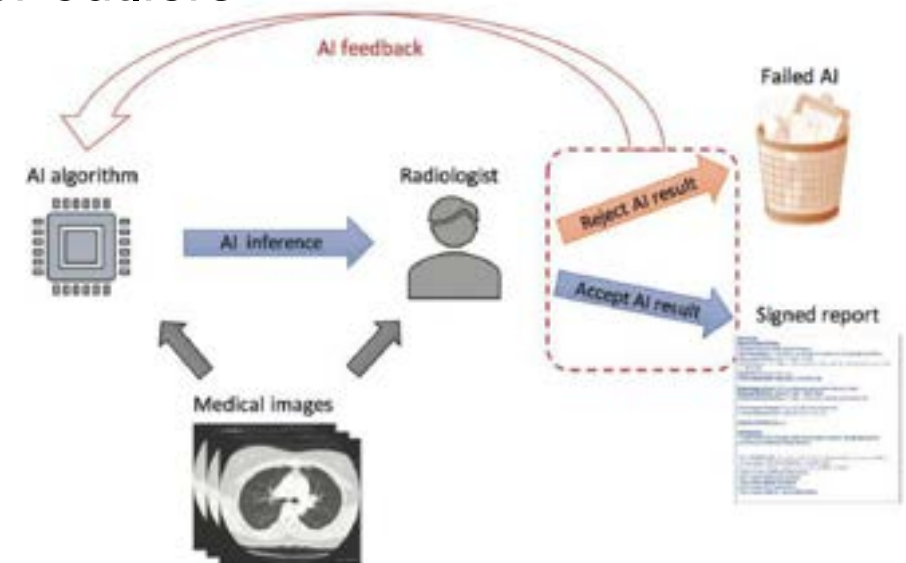  - **Installment 2** is under way

    - Filter banks (Wavelets, Gabor, …)

- All information is available at: https://theibsi.github.io

Zwanenburg, A., et al. (2020). The image biomarker standardization initiative: standardized quantitative radiomics for high-throughput image-based phenotyping. *Radiology*, 295(2), 328-338.

# Continuous learning

- Add new annotated samples regularly to update algorithms (for example with new machines)

    - Avoid <span style="color:red">catastrophic forgetting</span>

        • By adapting to a few specific cases or outliers

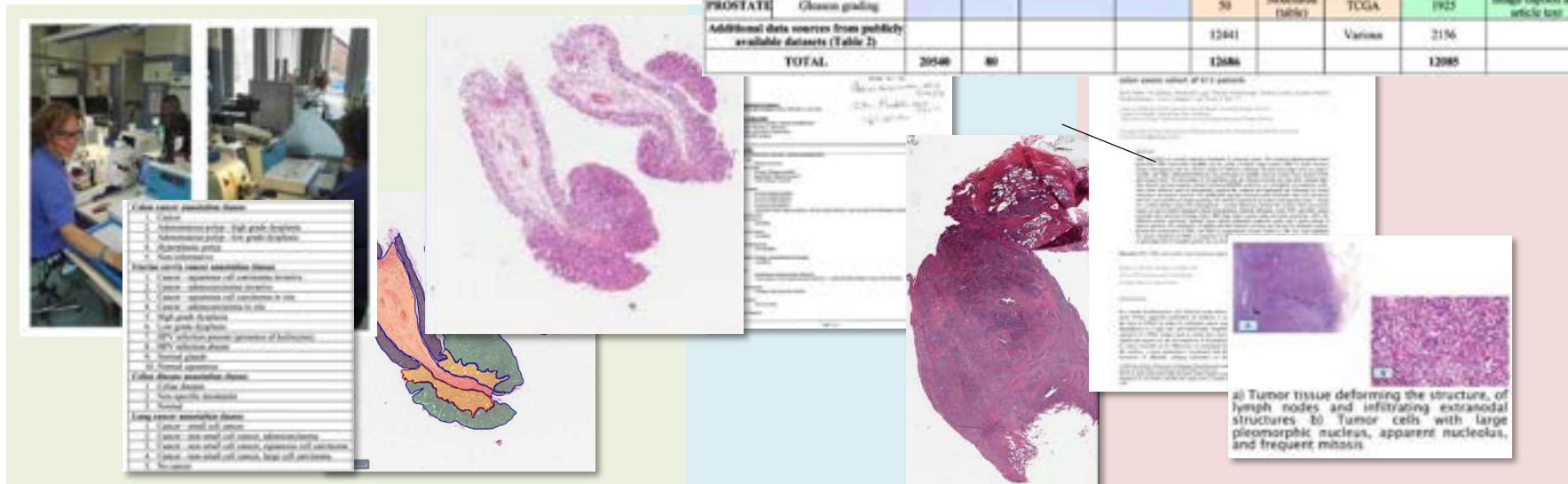- Regular <span style="color:red">feedback</span> loop

    - With clinicians using AI

Pianykh, O.S., Langs, G., Dewey, M., Enzmann, D.R., Herold, C.J., Schoenberg, S.O. and Brink, J.A., 2020. Continuous learning AI in radiology: implementation principles and early applications. *Radiology*, *297*(1), pp.6-14.

# Data used in ExaMode

# Weakly supervised learning

# Project status

# Weakly supervised learning from reports

# First results on multimodal data

# Clinical workflow and AI

- A clinician orders an image

- A radiologist/pathologists produces and views the image and writes a report based on the question and anamnesis

  - Much data on the patient (environment, prior diseases, genetics, blood tests, development of a condition, …)

  - Differential diagnosis, under much time pressure

- Any AI needs to be integrated into the workflow and tools

  - Adding evidence, identifying bias, uncertainty, …

    - Explaining the decisions and their context

# Interpretability of Deep Learning

- Make decisions <span style="color:red">understandable</span> & remove black box image

- Make sure that decisions are sound

- Explain why things may not be working



- In medicine it is particularly important to make sure that results can be explained & reproduced

  - High <span style="color:red">impact of wrong decisions</span>

- There are many approaches interpretability

  - 2D projections, PCA, TSNE

  - Class activation maps, saliency, …

# A taxonomy for explainability

- Many terms have been used in slightly different ways for AI: interpretability, explainability, transparency, accountability, fairness, (opacity) …
  - Bias, reliability, robustness, uncertainty, confidence
- A workshop was held in the summer of 2021 on this with views from several domains: legal, technical, philosophical, social, cognitive, ethical, …
  - https://taxonomyinterpretableai.wordpress.com
- EU is preparing the way

  M Graziani, L Dutkiewicz, D Calvaresi, J Pereira Amorim, K Yordanova, M Vered, R Nair, P Henriques Abreu, T Blanke, V Pulignano, JO. Prior, L Lauwaert, W Reijers, A Depeursinge, V Andrearczyk, H Müller, A Global Taxonomy of Interpretable AI: Unifying the Terminology for the Technical and the Social Sciences, Artificial Intelligence Reviews, 2022.

  - GDPR on data protection and AI policy
    - Limit the strong risks of AI and its use and abuse

# Taxonomy

# Regression concept vectors

- Identify <span style="color:red">existing clinical features</span> and check how the decision layers correlate to these features
  - i.e.: nuclei size, internal heterogeneity, borders, …
  - How much can a decision be explained with these?



M Graziani, V Andrearczik, H Müller, Concept attribution: Explaining CNN decisions to physicians, *Computers in Medicine and Biology*, 2020.

# Improve with interpretability

- Pre-trained models often include <span style="color:red">scale invariance</span>

- In medical applications this can be problematic, as scale carries information



M. Graziani, T. Lompech, H. Müller, A. Depeursinge, V. Andrearczyk, On the Scale Invariance in State of the Art CNNs Trained on ImageNet, *MDPI Make*, 2021.

# Visualizations

- Improve visualizations of regions that are relevant for the decision of a DNN

  - LIME is commonly used to highlight regions, but interpretations can be difficult



M. Graziani, I. Palatnik de Sousa, M.B.R. Vellasco, E. Costa da Silva, H. Müller, V. Andrearczyk, Sharpening Local Interpretable Model-agnostic Explanations for Histopathology: Improved Understandability and Reliability, MICCAI conference proceedings 2021, Springer LNCS, Strasbourg, France, 2021.

# The importance of user tests!

- Most systems are scripts run under laboratory conditions

  - Does not give many indications of routine use

- Impact of the system is hard to measure

  - Better decisions, more confidence, faster, satisfaction?

- What is the influence on the patient?

  - Better treatment? Longer survival? Quality of life?

- User tests are more complex to set up but can really help

- AI and users are usually best together

# Scientific challenges

- Cooperation, <span style="color:red">Coopetition</span>, Competition
- Many data sets are now available

  - Also many medical data sets

- <span style="color:red">Strong baselines</span> help to judge quality

  - Not only the results count!

- Challenges can be run without sharing confidential data

  - Provide VMs or Docker containers



Jimenez-del-Toro, Oscar, et al. "Cloud-based evaluation of anatomical structure segmentation and landmark detection algorithms: VISCERAL anatomy benchmarks." *IEEE transactions on medical imaging* 35.11 (2016): 2459-2475.

# Some more best practices

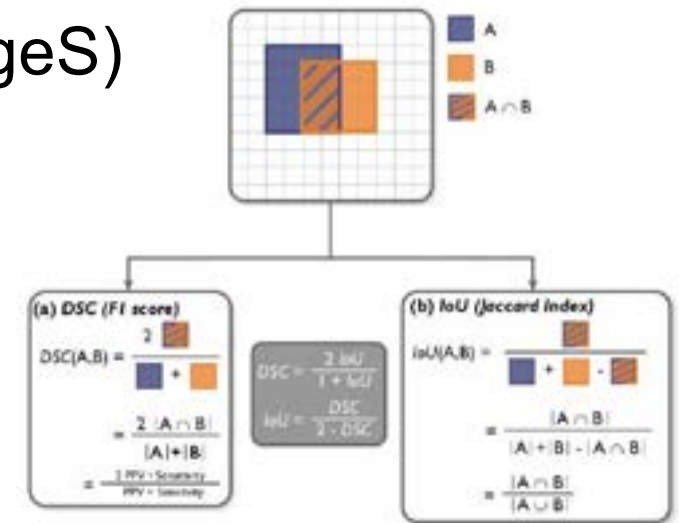- Reporting for scientific challenges in medical imaging

  - BIAS (Biomedical Image Analysis challengeS)

  - Avoid bias, use the right measures

  - Use meaningful data sets and scenarios

  - How to chose the best evaluation metrics

  - …

Maier-Hein, L., Eisenmann, M., Reinke, A., Onogur, S., Stankovic, M., Scholz, P., Arbel, T., Bogunovic, H., Bradley, A.P., Carass, A. and Feldmann, C., 2018. Why rankings of biomedical image analysis competitions should be interpreted with care. Nature communications, 9(1), pp.1-13.
Maier-Hein, L., Reinke, A., Kozubek, M., Martel, A.L., Arbel, T., Eisenmann, M., Hanbury, A., Jannin, P., Müller, H., Onogur, S. and Saez-Rodriguez, J., 2020. BIAS: Transparent reporting of biomedical image analysis challenges. *Medical image analysis*, *66*, p.101796.
Reinke, Annika, et al. "Common limitations of image processing metrics: A picture story." *arXiv preprint arXiv:2104.05642* (2021).

# Tasks and measures

# Certification of medical SW

- Any use of AI in medicine needs to be <span style="color:red">certified</span> (CE, FDA)

    - Software is a "medical device"

    - Unless only for a research study

    - Avoid risks for the patient, tedious process

- <span style="color:red">In-vitro diagnostics</span> is more complex since 2022

    - Transition period for already certified tools

- Expensive to do, so not usable for research tools

# Conclusions

- Medical AI is an extremely <span style="color:red">interesting</span> domain

    - With <span style="color:red">high impact</span> on people's lives!

- AI in medical imaging has many challenges remaining!

    - Some can be addressed relatively easily

    - Many will require much more research

- Consequences of (wrong) decisions are important

- Run user tests (also on prospective data)

# Contact

- More information can be found at

  - http://medgift.hevs.ch/

  - http://publications.hevs.ch/

- Contact: Henning.mueller@hevs.ch