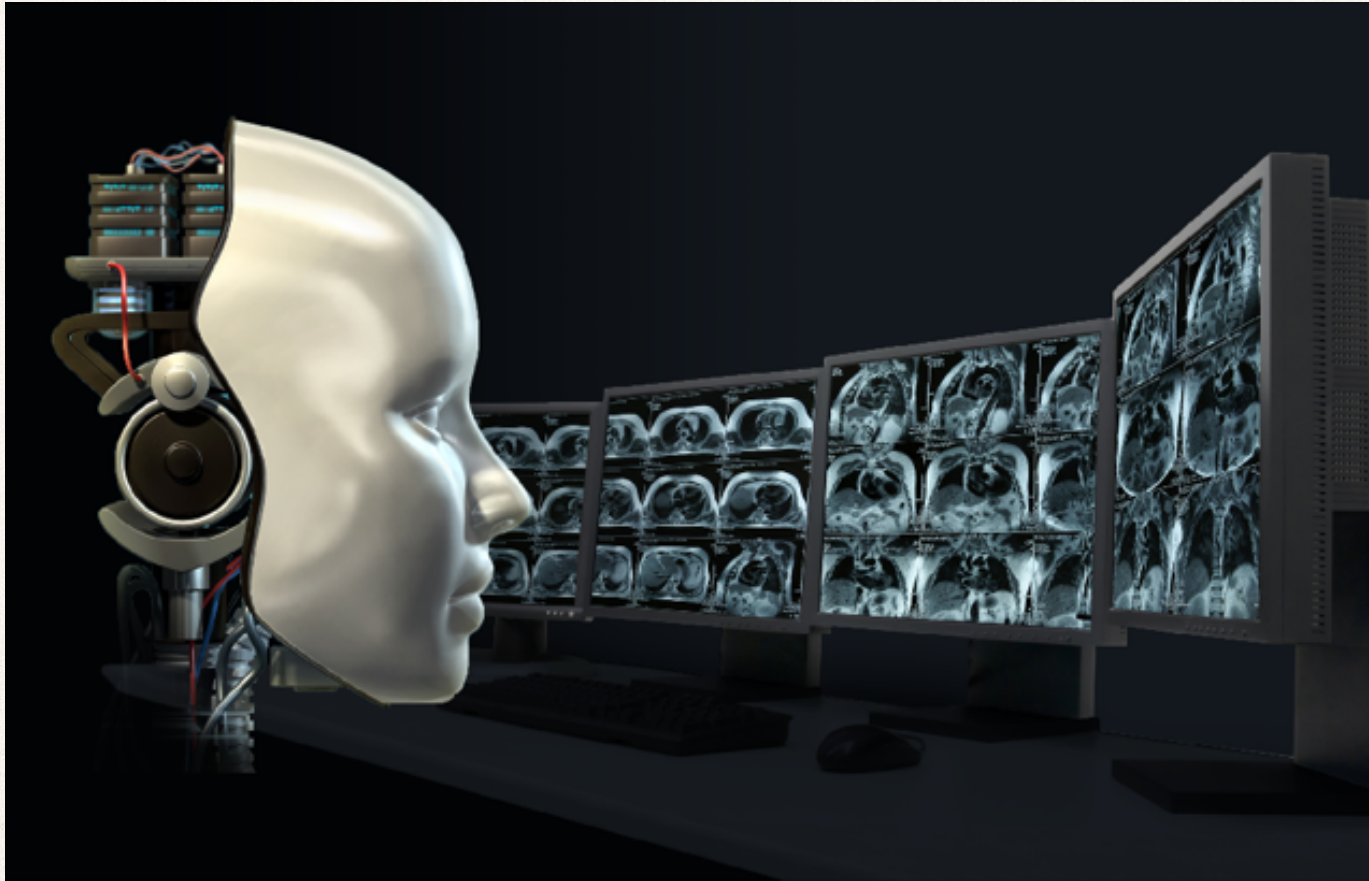# Explainable Deep Learning for High Risk AI Applications



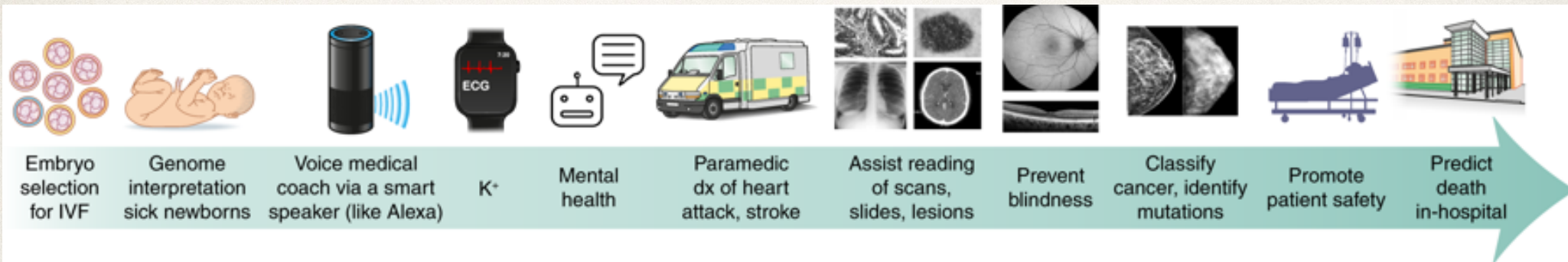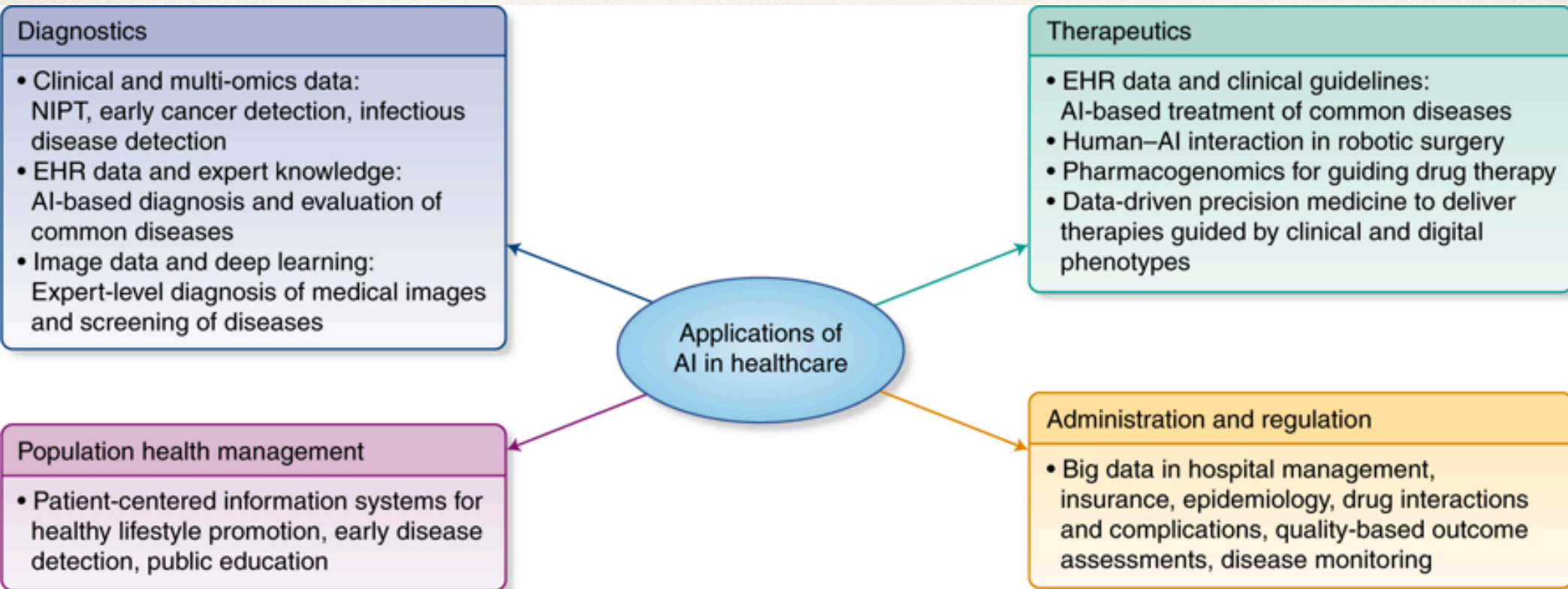**Ulas Bagci, PhD.**, Center for Research in Computer Vision, University of Central Florida.
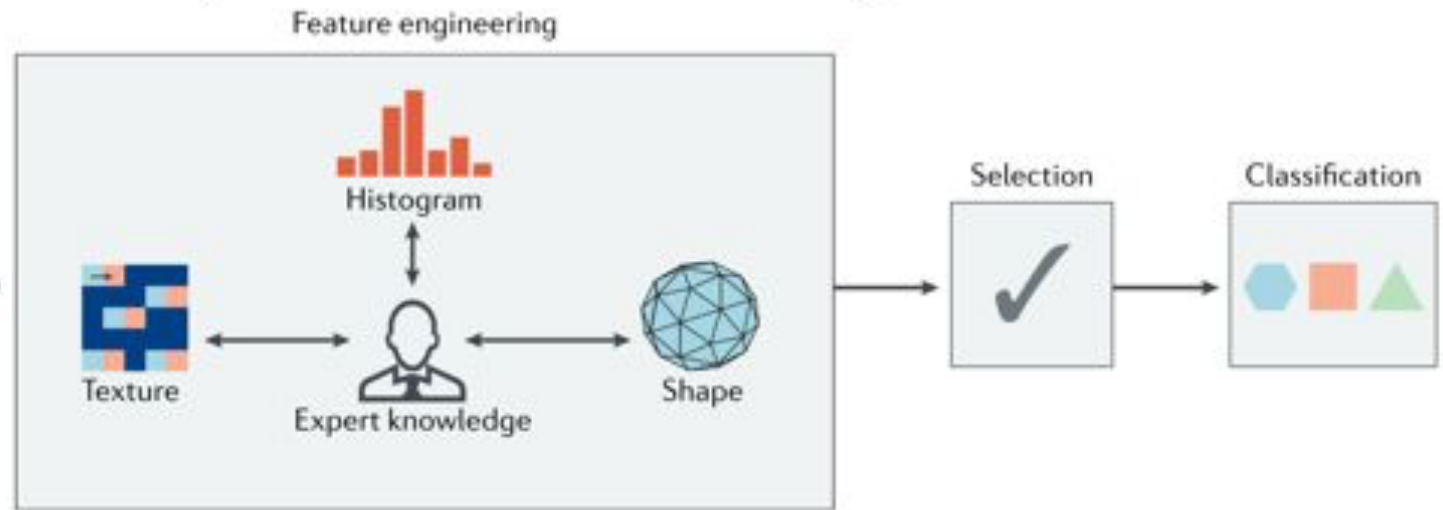
*Picture Credit: "Aidoc"*

# Outline of today's talk

* Interpretability / Explainability in DL

* Image based diagnosis as a high-risk AI application

* Eye-Tracking for human in the loop DL systems

* Visual attribute learning for building the thrust

# Deep Learning / AI revolutionizes Medicine!
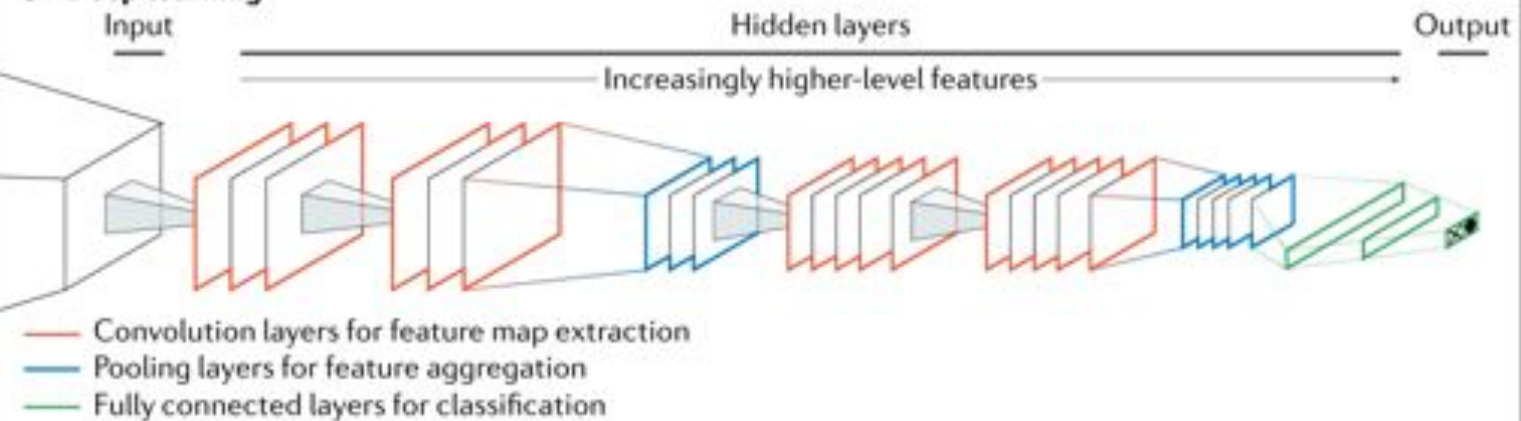


**Diagnostics**
- Clinical and multi-omics data:
  NIPT, early cancer detection, infectious disease detection
- EHR data and expert knowledge:
  AI-based diagnosis and evaluation of common diseases
- Image data and deep learning:
  Expert-level diagnosis of medical images and screening of diseases

**Therapeutics**
- EHR data and clinical guidelines:
  AI-based treatment of common diseases
- Human–AI interaction in robotic surgery
- Pharmacogenomics for guiding drug therapy
- Data-driven precision medicine to deliver therapies guided by clinical and digital phenotypes

**Applications of AI in healthcare**

**Population health management**
- Patient-centered information systems for healthy lifestyle promotion, early disease detection, public education

**Administration and regulation**
- Big data in hospital management, insurance, epidemiology, drug interactions and complications, quality-based outcome assessments, disease monitoring

Embryo selection for IVF | Genome interpretation sick newborns | Voice medical coach via a smart speaker (like Alexa) | K⁺ | Mental health | Paramedic dx of heart attack, stroke | Assist reading of scans, slides, lesions | Prevent blindness | Classify cancer, identify mutations | Promote patient safety | Predict death in-hospital

*Credit: "E.Topol"*

# Less Artificial - More Intelligent!



**a  Predefined engineered features + traditional machine learning**

Feature engineering

Histogram

Texture

Expert knowledge

Shape

Selection

Classification

**b  Deep learning**

Input

Hidden layers

Increasingly higher-level features

Output

— Convolution layers for feature map extraction
— Pooling layers for feature aggregation
— Fully connected layers for classification

# Mostly Black-Box Nature of DL

- Imagine a physician using a DNN to diagnose a patient.

- S/he will most likely **not trust** an automated diagnosis unless s/he **understands the reason** behind a certain prediction (e.g. highlighted regions in the brain that differ from normal subjects) allowing him/her to verify the diagnosis and reason about it.

**World**

**Data**

⇧ capture

**World**

**Black Box Model**

⬆ learn

**Data**

⬆ capture

**World**

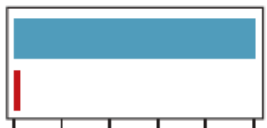# Why should I trust AI?



**Original image**

Dermatoscopic image of a benign melanocytic nevus, along with the diagnostic probability computed by a deep neural network.
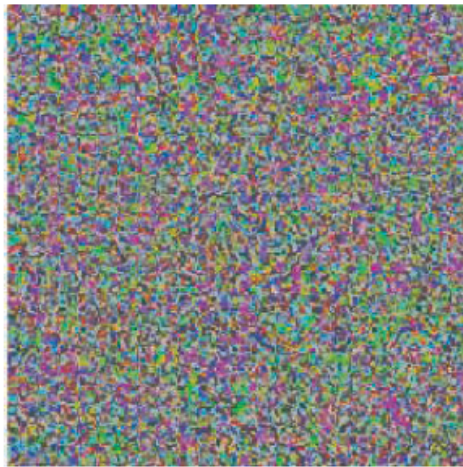
Benign
Malignant
Model confidence

**Diagnosis: Benign**

+ 0.04 ×

**Adversarial noise**

Perturbation computed by a common adversarial attack technique.
See (7) for details.
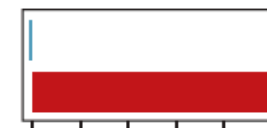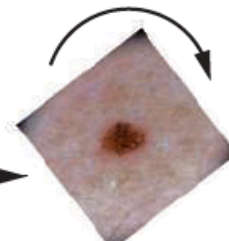
=

**Adversarial example**

Combined image of nevus and attack perturbation and the diagnostic probabilities from the same deep neural network.

Benign
Malignant
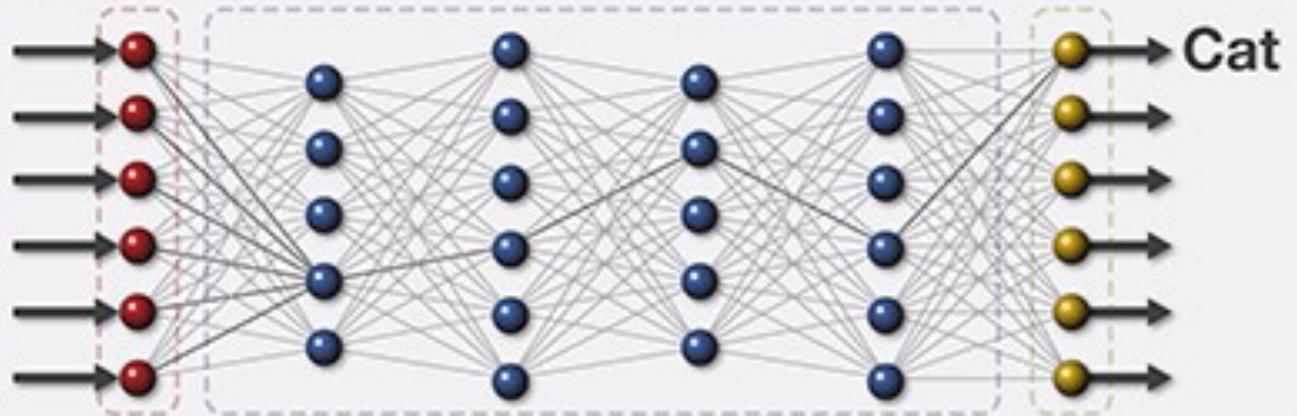Model confidence

**Adversarial rotation** (8)

**Diagnosis: Malignant**

# AI fails badly too!

- Uber self-driving car kills a pedestrian

- Amazon AI recruiting tool is gender biased

- Google Allo suggested man in turban emoji as response to a gun emoji

- Google Translate shows gender bias in Turkish-English translations (doctors, hard-working —> he, nurses, lazy —-> she)

- Facebook chatbots shut down after developing their own language

- AI misses the mark with Kentucky Derby predictions

- Google Home Minis spied on their owners

- …

Source: https://www.darpa.mil/program/explainable-artificial-intelligence
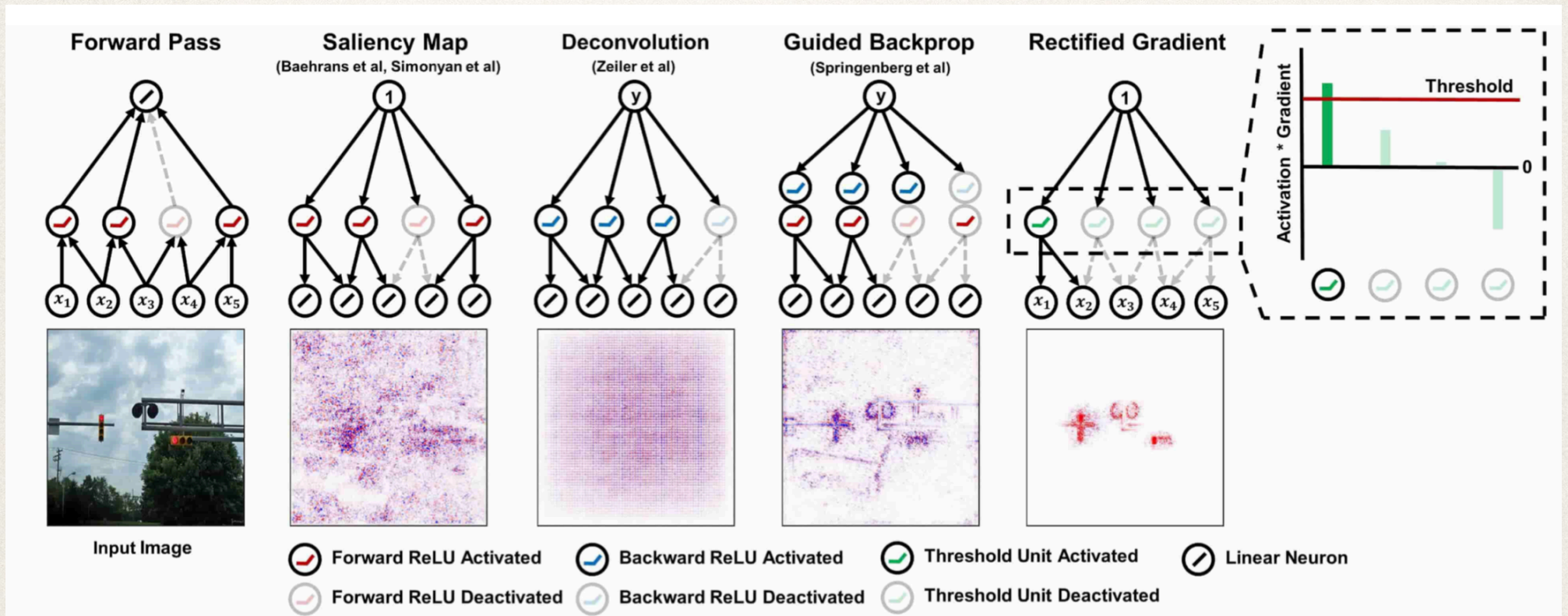
# Current XAI Approaches

✤ One qualitative approach is to **highlight areas** that provide evidence in favor of, and against choosing a certain class.

- Filtering/Visualization

- Perturbation based methods (saliency maps, CAM, etc)

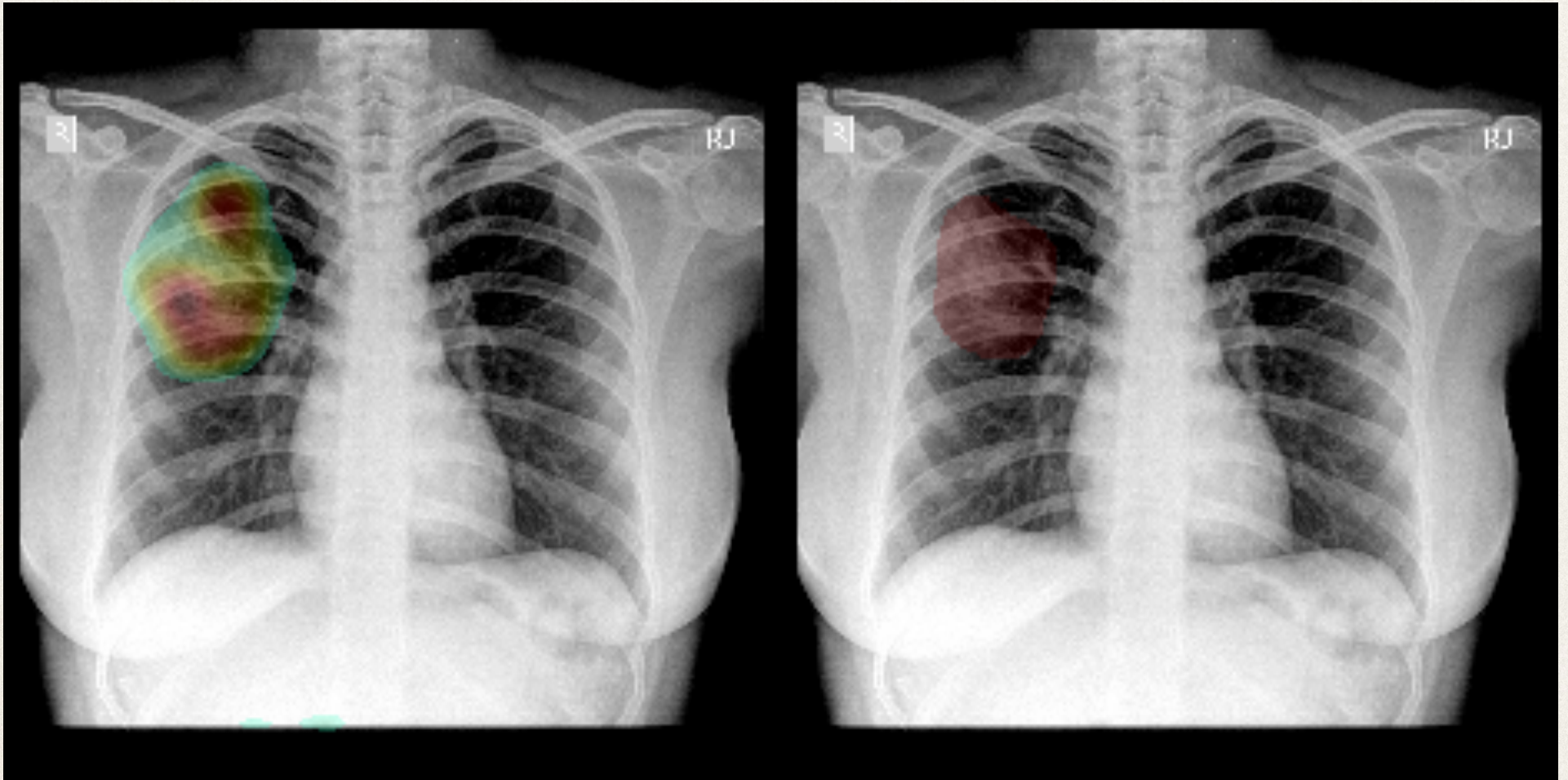# Current XAI Approaches



* The pixels which contribute maximally to the prediction, once altered, would drop the probability by the maximum amount.
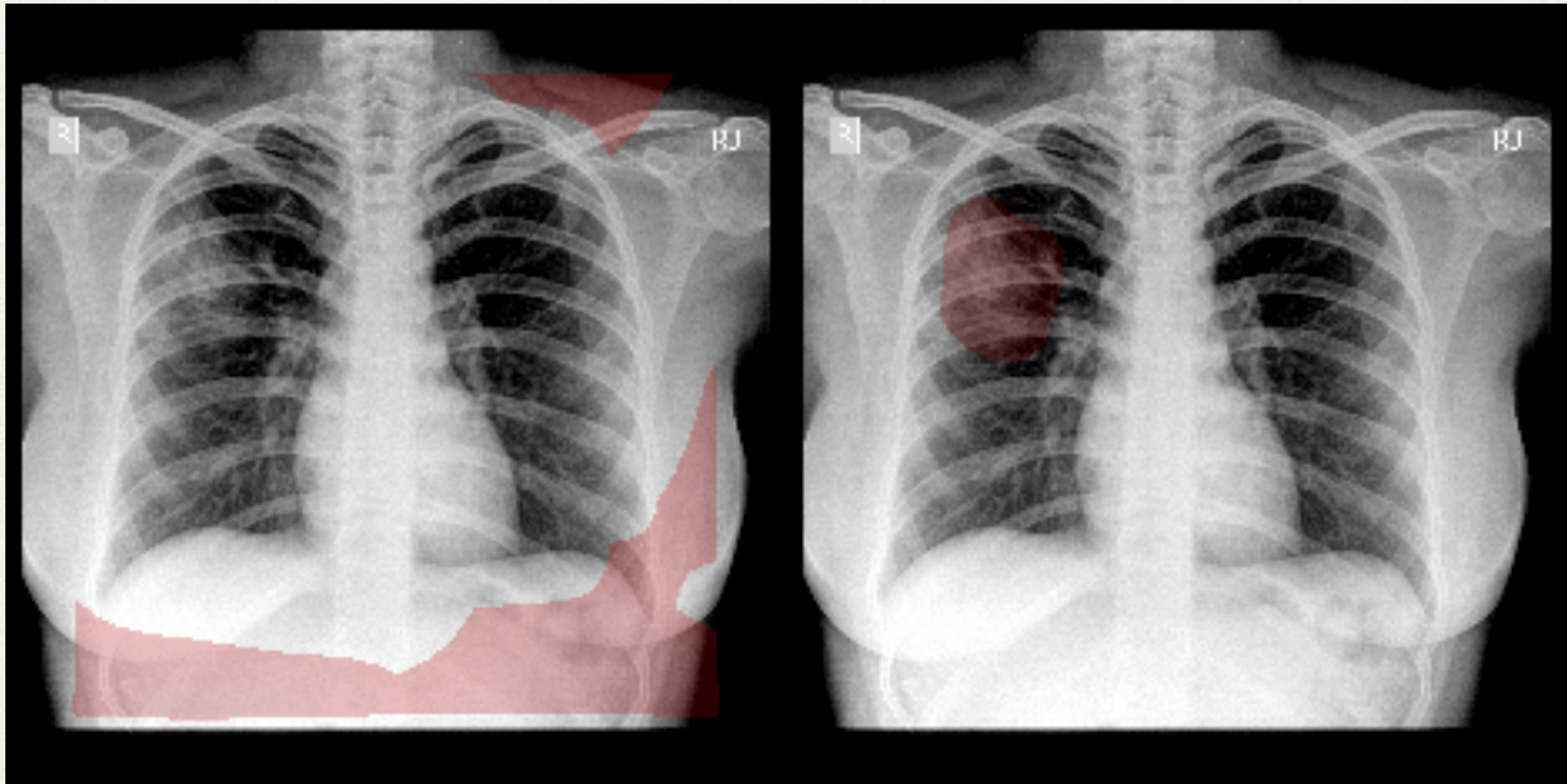
# Current XAI Approaches



**Qure.AI:** Heatmap by GuidedBackprop against original annotation.

# Drawbacks of Current XAI Algorithms

**Qure.AI:** Heatmap by GuidedBackprop against original annotation.



Not completely "true" explanation/reasoning
Artifact generation in visual maps
Limitation to specific architectures

….

# What we propose

✤ Building-in-thrust (human in the loop)

✤

# What we propose

* Building-in-thrust (human in the loop)

* Explainable / interpretable DL system

# What we propose

* Building-in-thrust (human in the loop)

* Explainable / interpretable DL system

**Interpretable to Whom? A Role-based Model for Analyzing Interpretable Machine Learning Systems**

Richard Tomsett[1]   Dave Braines[1,2]   Dan Harborne[2]   Alun Preece[2]   Supriyo Chakraborty[3]

**Defn.** Interpretability is a domain-specific notion, so there cannot be all purpose definition.

# Radiologist Centered AI Protocols

- Human + AI   >      AI
  - (human in the loop ML)

- get a computer system to learn some intelligence behavior by training it on <u>large</u> amount of data.
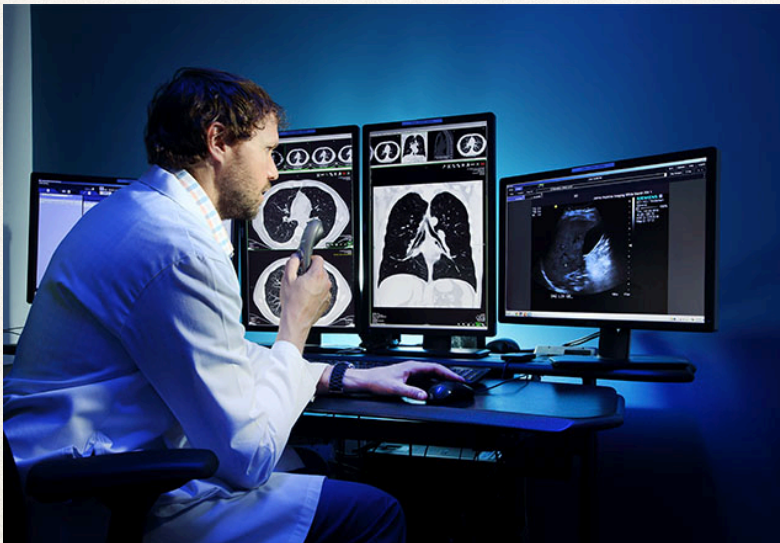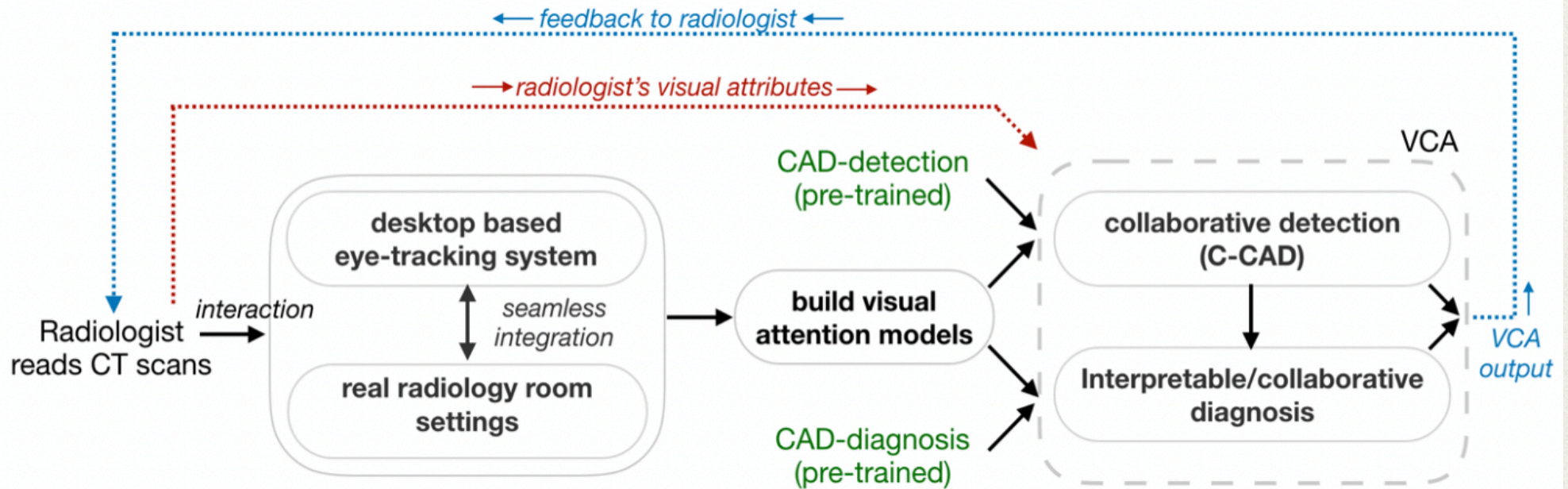
# Radiologist Centered AI Protocols

- Human + AI > AI
  - (human in the loop ML)

- get a computer system to learn some intelligence behavior by training it on <u>large</u> amount of data.



## Example High Risk AI Application:
*Detection and Malignancy characterization of lung nodule in CT images*

# RCAI (radiologist centered AI)





- ✤ Dedicated light source
- ✤ darkened environment
- ✤ limited distraction
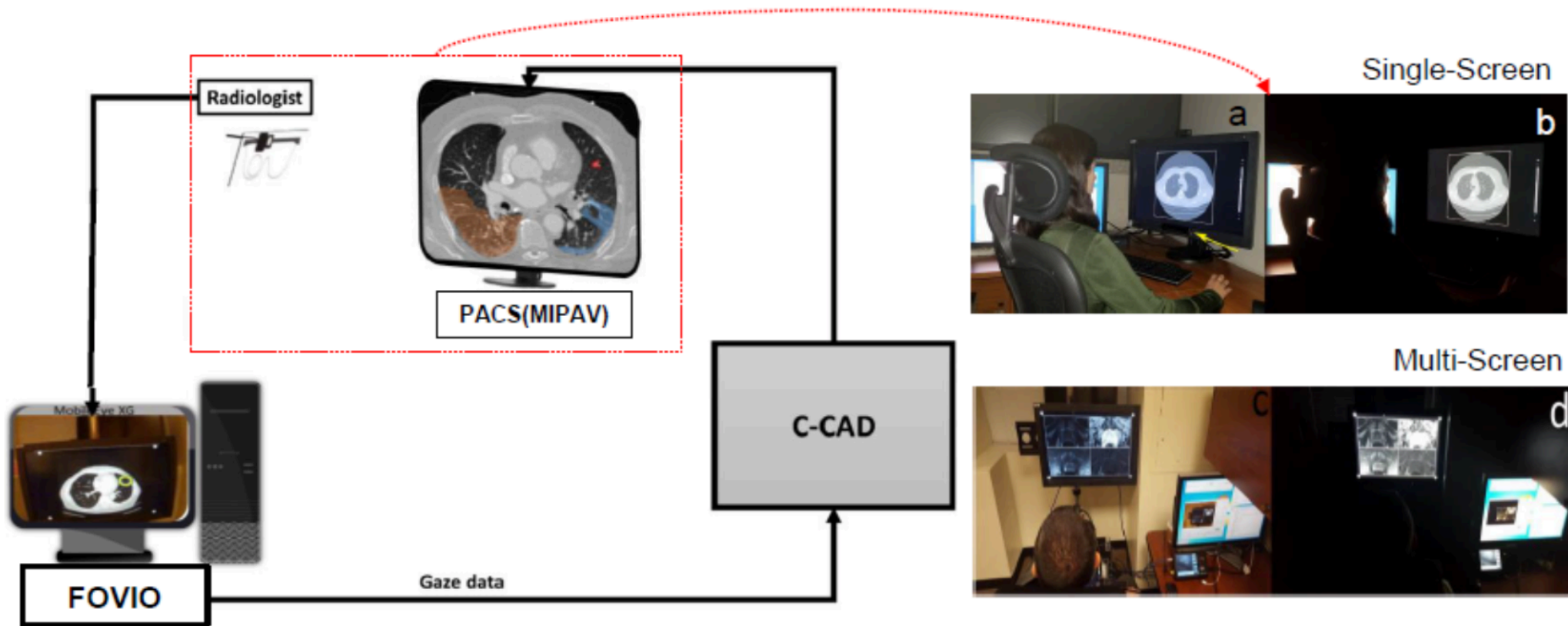
# Detection via Real-Time Eye-Tracking

# Visual Search (Eye-Tracking) + AI Integration

**Soln:** Combine complementary strengths of radiologists and AI



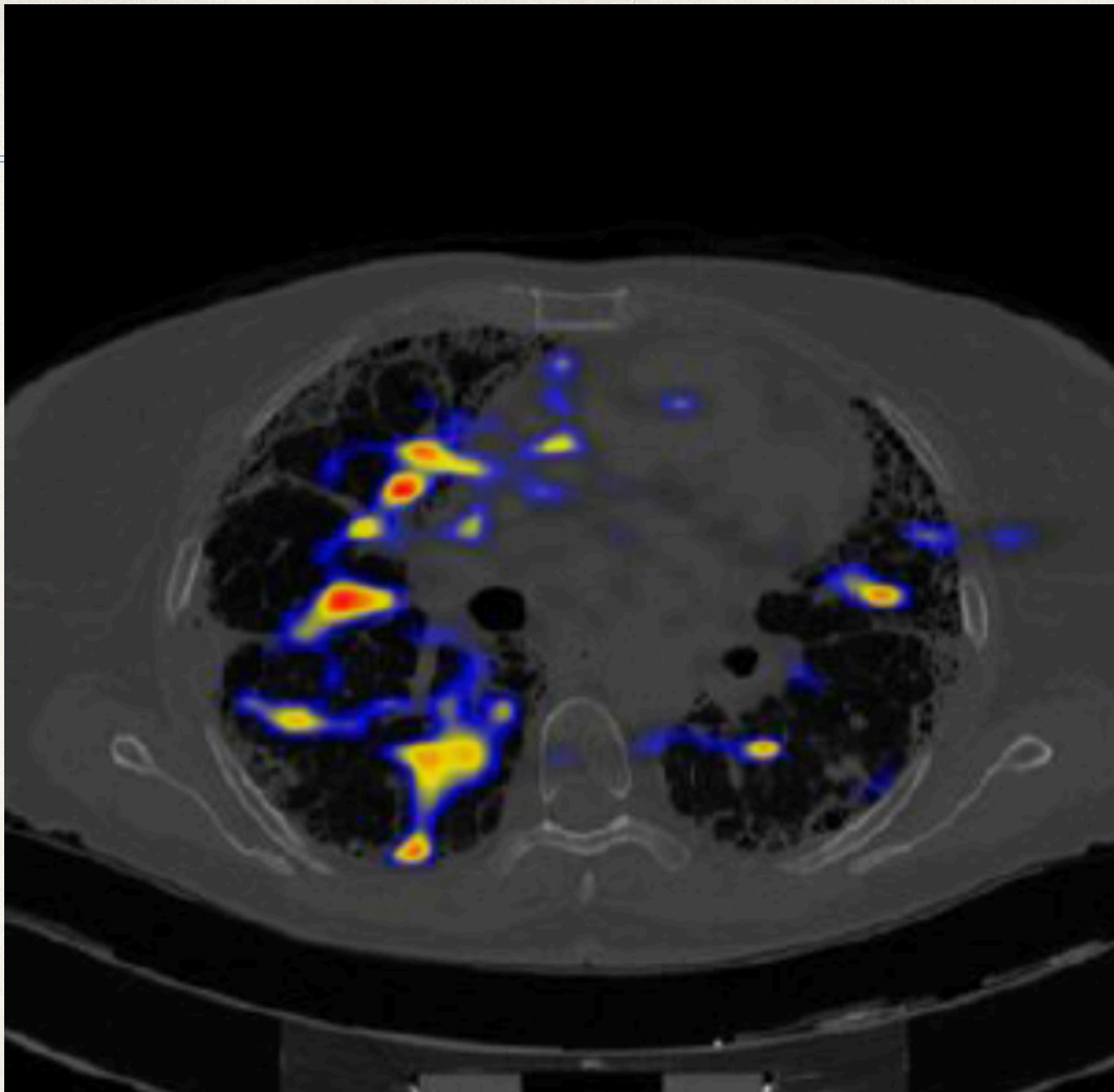**In real radiology rooms, in realistic settings!**
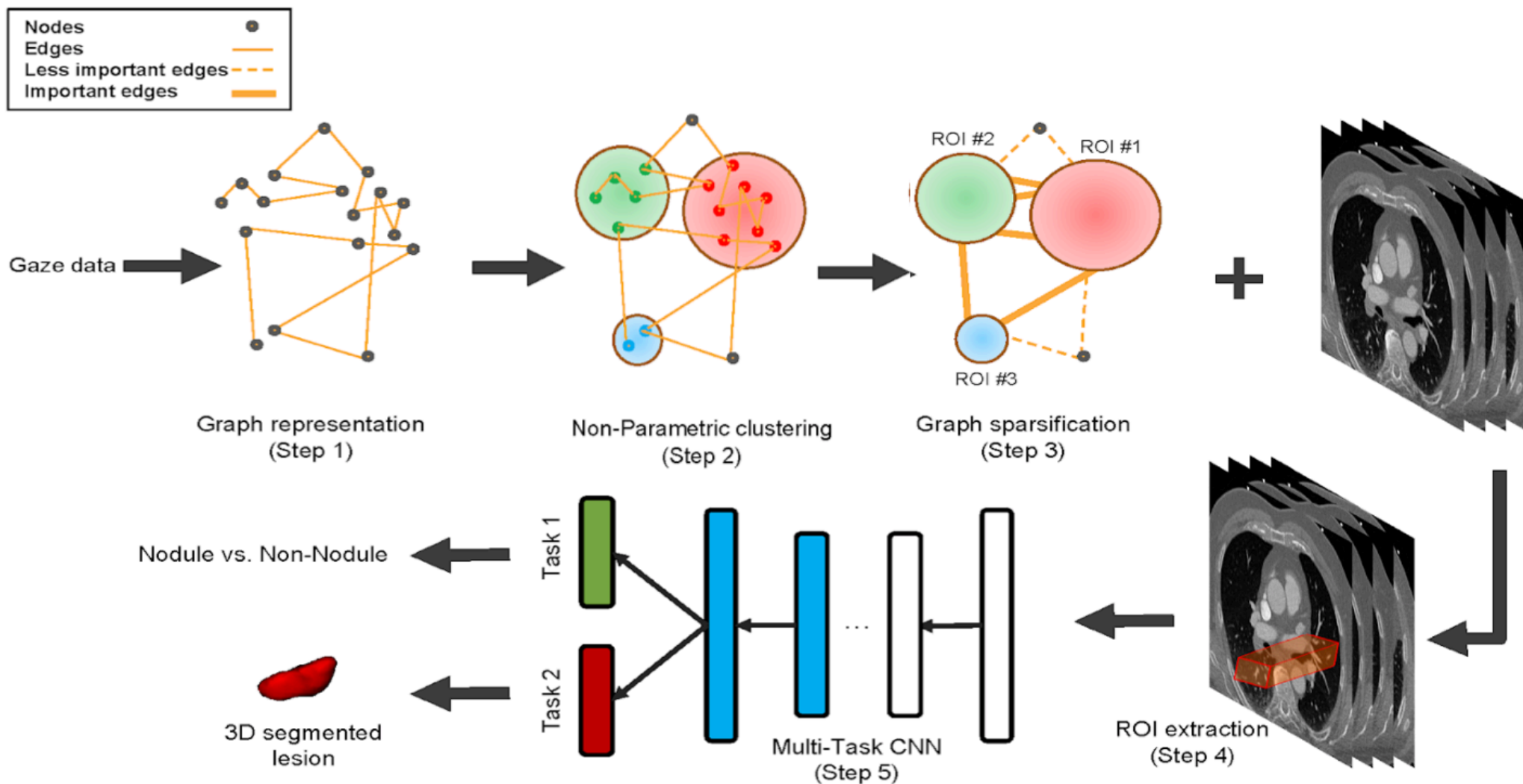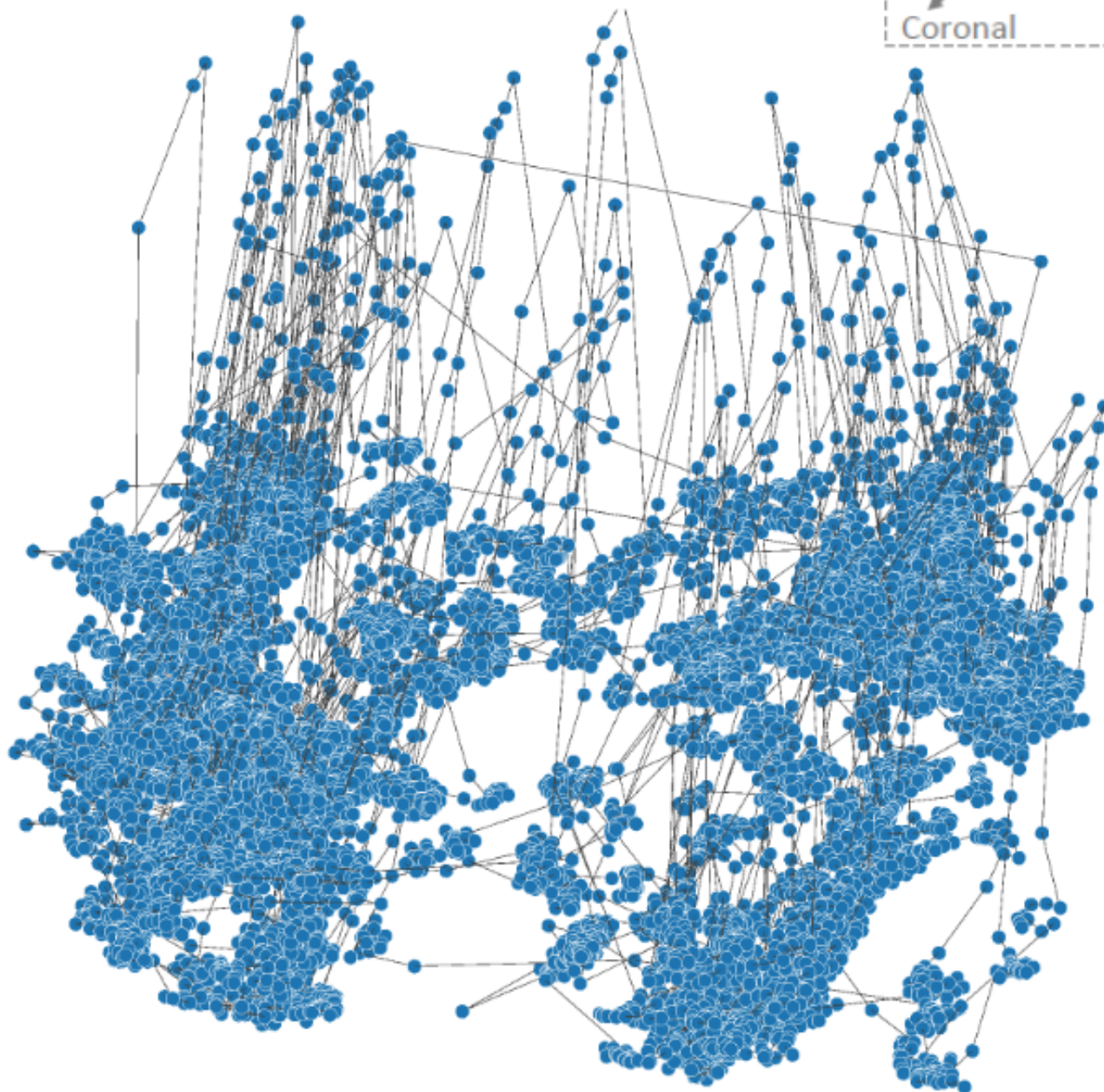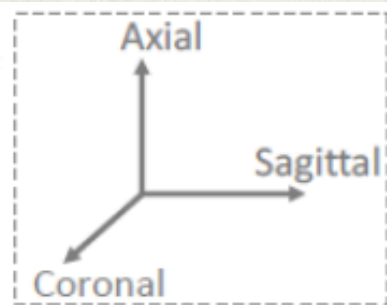
# Eye-Tracker / Device Info.



- **Fovio™ Eye Tracker** remote system.

- **System Type:** Remote (contact-free)

- **Sampling Rate:** 60Hz

- **Method:** Proprietary Algorithm
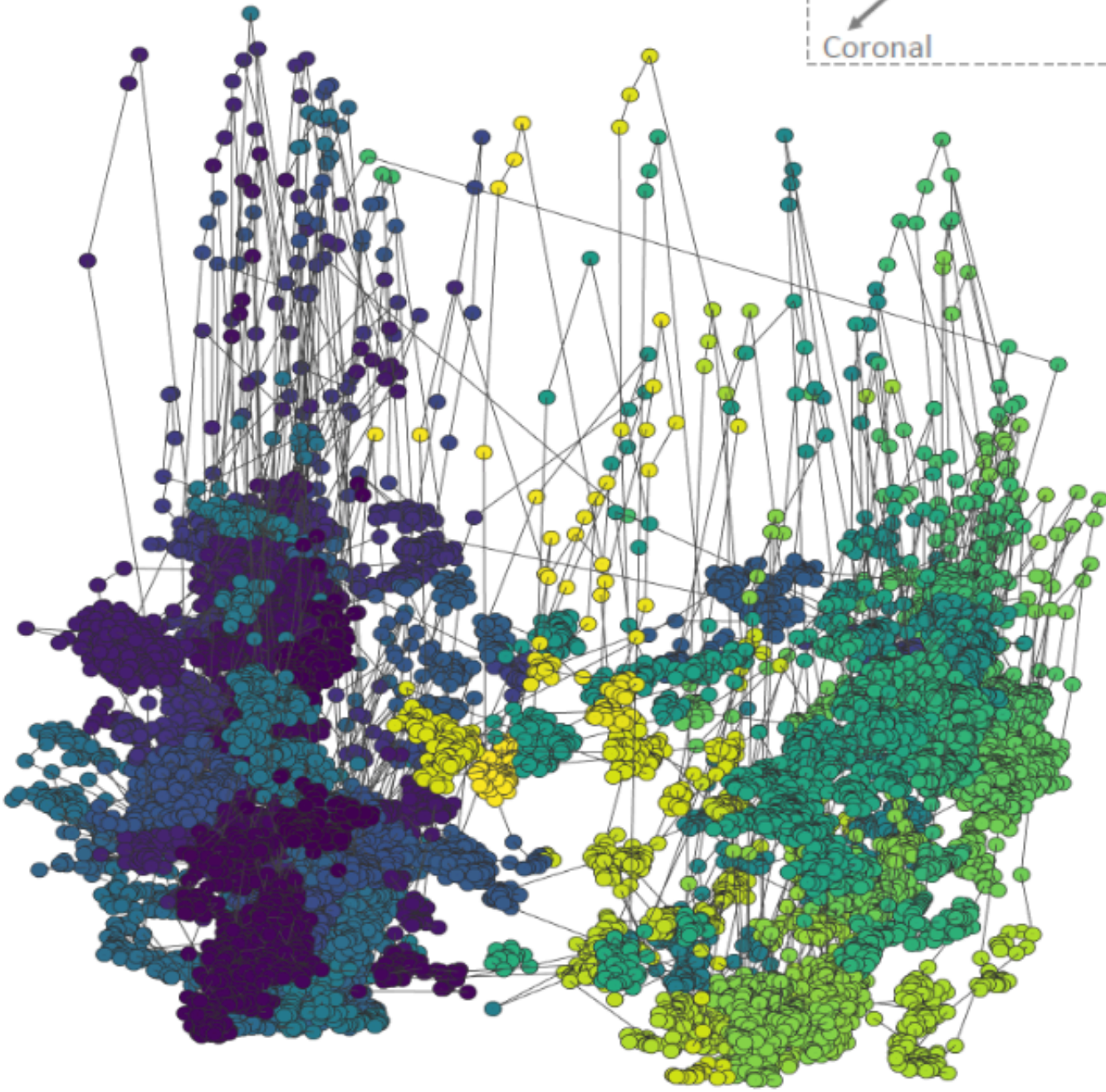
- **Binocular Tracking:** Yes

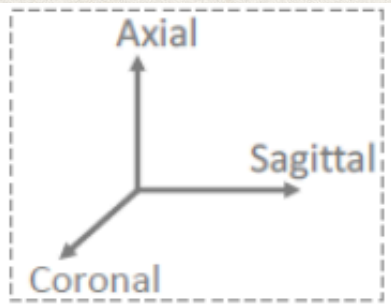- **Accuracy:** 0.78 Degrees (Mean) 0.59 (Std. Dev.) angular error

- **Head Box:** 31cm x 40cm @ 65cm - range 40-80cm

- **Additional Details:** Large head box, robust to glasses and ambient light, multi-display tracking

# Reduce Gaze Points for Faster Analysis!



Raw data           Clustered data           Graph Sparsification

**Radiologist #1**

Dense

Time analysis

N=511, Tc=8.51sec, T=121.06sec

Sparse

|  | Radiologist #1 | Radiologist #2 |
|---|---|---|
| Dense | | |
| Time analysis | | |
| | N=511, Tc=8.51sec, T=121.06sec | N=867, Tc=14.45sec, T=143.3sec |
| Sparse | | |

|  | Radiologist #1 | Radiologist #2 | Radiologist #3 |
|---|---|---|---|
| Dense | | | |
| Time analysis | | | |
| | N=511, Tc=8.51sec, T=121.06sec | N=867, Tc=14.45sec, T=143.3sec | N=958, Tc=15.96sec, T=337.63sec |
| Sparse | | | |

|  | Raw data | Clustered data | Graph build on cluster centers data | Sparsified data |
|---|---|---|---|---|
| Radiologist #1 | | | | |
| Radiologist #2 | | | | |
| Radiologist #3 | | | | |

| FPs/scan | 0.125 | 0.25 | 0.5 | 1 | 2 | 4 | 8 | Average |
|---|---|---|---|---|---|---|---|---|
| Sensitivity | 0.773 | 0.870 | 0.924 | 0.941 | 0.962 | 0.980 | 0.986 | **0.919** |

# Diagnosis via Visual Explanations

# Visual Explanations via Attributes

Every lung nodule is associated with 6 attributes provided by the radiologist:

–Calcification
–Sphericity
–Margin
–Lobulation
–Spiculation
–Texture

**Interpretable to Whom? A Role-based Model for Analyzing Interpretable Machine Learning Systems**

**Richard Tomsett** [1]   **Dave Braines** [1,2]   **Dan Harborne** [2]   **Alun Preece** [2]   **Supriyo Chakraborty** [3]

# Visual Explanations via Attributes

Every lung nodule is associated with 6 attributes provided by the radiologist:

–Calcification
–Sphericity
–Margin
–Lobulation
–Spiculation
–Texture



- *We explore the significance of these attributes to determine malignancy*
- *We concatenated these attribute score with 4096 dimension feature vector of CNN and perform Gaussian Progress regression*
- *LIDC-IDRI data base was used (1018 CT scans), multiple radiologists annotated the data sets.*

# Multi-Task Learning of Visual Attributes

| Methods | Accuracy | Mean Score Diff |
|---------|----------|-----------------|
| GIST+LASSO | 76.83% | 0.6753 |
| 3D CNN MTL + Trace | 80.08% | 0.6259 |
| **Proposed approach** | **91.26%** | **0.4593** |

# Drawbacks

* The requirement of large scale well annotated data

* Lack of object-part relationship with typical CNNs

* Fragile nature of the CNN systems (easily fooled!)

# CapsNet (Capsule Networks)

Two Simple Changes from CNNs:

▶ Features are now represented as **vectors** rather than scalars.

   ▶ Vectors **store orientation information** about the input.

▶ Agreement between feature "predictions" is computed to **weight the presence and orientation of higher-level features**.

# CapsNet (Capsule Networks)

$$W_{ij} = [8 \times 16]$$

## DRAWBACKS OF POOLING

▶ A form of routing, just an unintelligent one.

▶ **Pro:** Some spatial-invariance.

▶ **Pro:** Reduces memory burden.

▶ **Con:** Throws away information without regard to importance/ heterogeneity of the region.

▶ **Capsules** use strided-overlapping convolutions and dynamic routing.

# CapsNet (Capsule Networks)

*Sabour, et al. NeurIPS 17*



- ▶ Requires less training data for good generalization.
- ▶ Preserves part-whole relationships and shape information.
- ▶ Capsule vectors encode information about input.



Figure 1: Stacked Capsule Autoencoder (SCAE): (a) *part* capsules segment the input into parts and their poses. The poses are then used to reconstruct the input by affine-transforming learned templates. (b) *object* capsules try to arrange inferred poses into objects, thereby discovering underlying structure. SCAE is trained by maximizing image and part log-likelihoods subject to sparsity constraints.

# X-Caps

## X-Caps: Explainable Capsule Networks

# DX-Caps



DX-Caps: Deep Explainable Capsule Networks

Table 2: Prediction accuracy of visual attribute learning with capsule networks. Dashes (-) represent values which the given method could not produce. *X-Caps* significantly outperforms the state-of-the-art explainable method (*HSCNN*) at attribute modeling (the main goal of both studies), while also producing higher malignancy prediction scores, approaching state-of-the-art non-explainable methods performance.

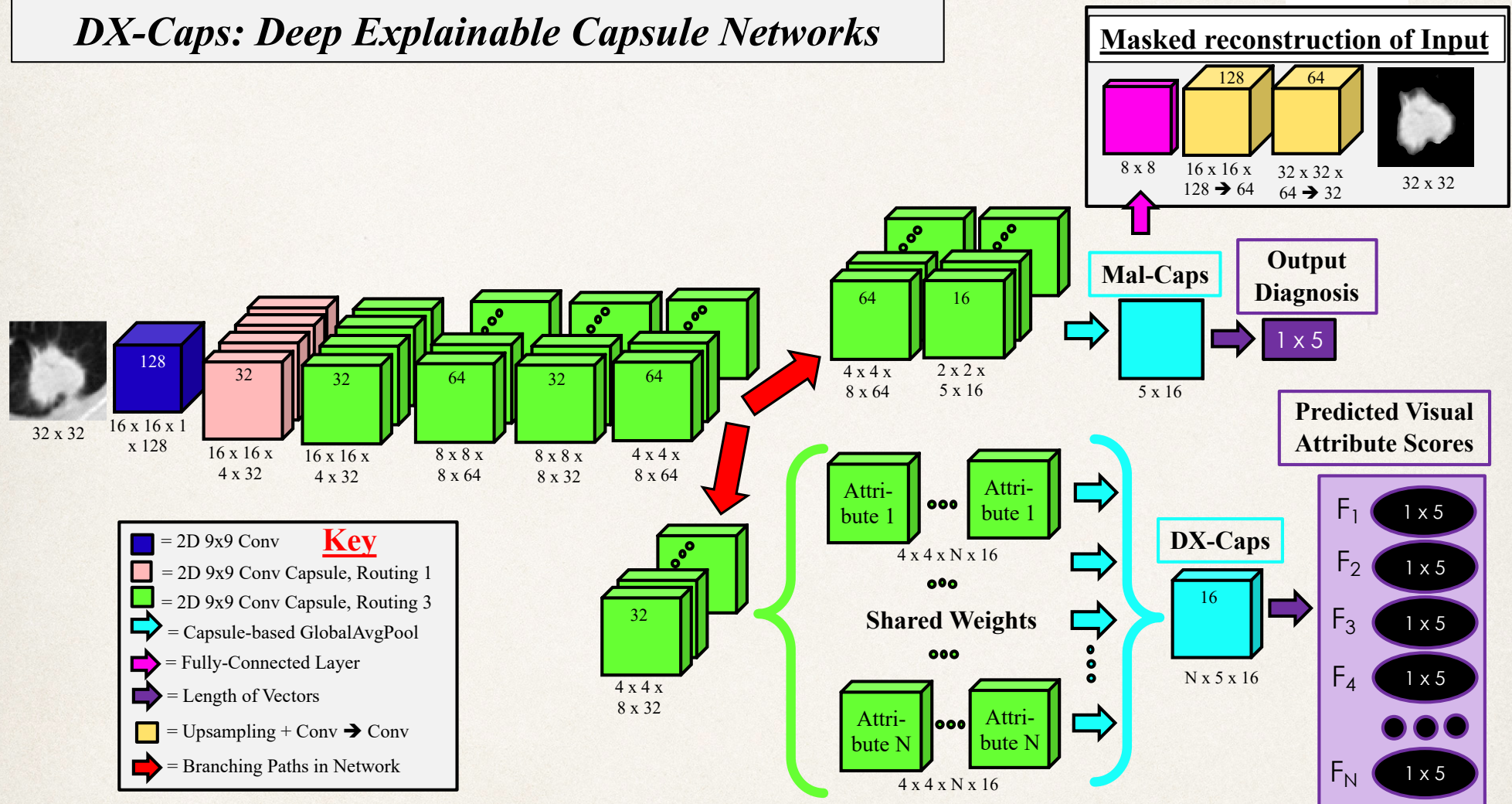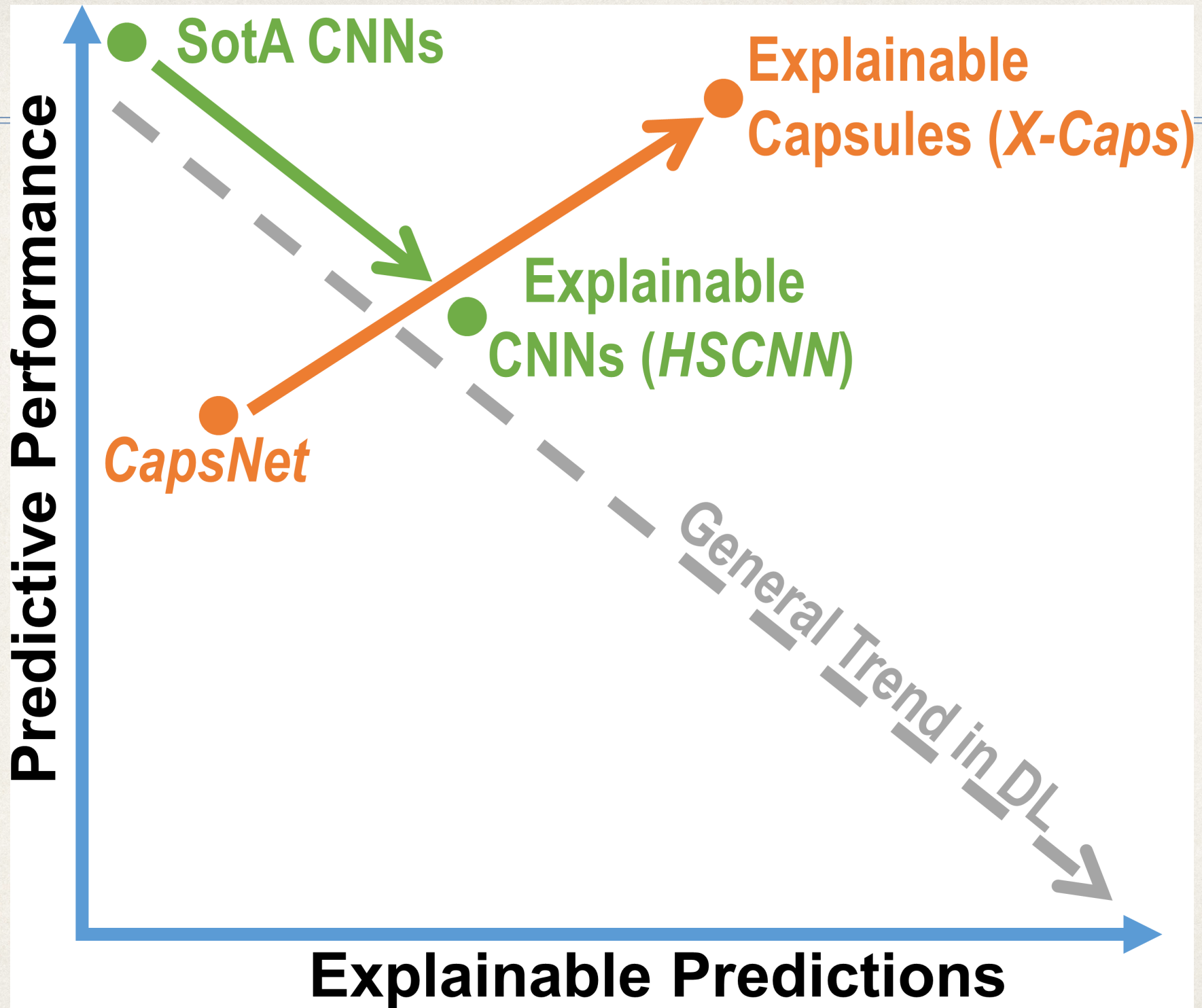| | Attribute Prediction Accuracy % | | | | | | Malignancy Accuracy % |
|---|---|---|---|---|---|---|---|
| | subtlety | sphericity | margin | lobulation | spiculation | texture | |
| **Non-Explainable Methods** | | | | | | | |
| 3D Multi-Scale + RF [31] | - | - | - | - | - | - | 86.84 |
| 3D Multi-Crop [32] | - | - | - | - | - | - | 87.14 |
| 3D Multi-Out-DenseNet [6] | - | - | - | - | - | - | 90.40 |
| 3D Dual-Path GBM [38] | - | - | - | - | - | - | 90.44 |
| *CapsNet* [28] | - | - | - | - | - | - | 77.04 |
| **Explainable Methods** | | | | | | | |
| 3D Dual-Path-Dense *HSCNN* [30] | 71.9 | 55.2 | 72.5 | - | - | 83.4 | 84.20 |
| **Proposed *X-Caps*** | **90.39** | **85.44** | **84.14** | **70.69** | **75.23** | **93.10** | **86.39** |

# Concluding Remarks

**1** **STRENGTHEN TRUST AND TRANSPARENCE**

Without understanding the contribution of each explanatory variable to the outcome, we will have no guarantee that the model will make a relevant and fair recommendation.

# Concluding Remarks



2 **EXPLAIN DECISIONS**

An interpretable Machine Learning model allows humans to understand the proposed outcome and establish the diagnosis.

# Concluding Remarks

## 3 IMPROVE THE MODELS

Interpretability ensures data scientists that the model is good for the right reasons and wrong for the right reasons as well. Interpretability offers new possibilities for feature engineering and model debugging.

# Concluding Remarks

Visual search is gold standard; however, prone to errors, malpractice/perceptual error etc., time consuming, and sub-optimal

# Concluding Remarks

Visual search is gold standard; however, prone to errors, malpractice/perceptual error etc., time consuming, and sub-optimal

Computer (AI) helps radiologists to find pathologies that can be missed however, computer also depicts so many false positives that radiologists easily capture them with true labels!

# Concluding Remarks

Visual search is gold standard; however, prone to errors, malpractice/perceptual error etc., time consuming, and sub-optimal

Computer (AI) helps radiologists to find pathologies that can be missed however, computer also depicts so many false positives that radiologists easily capture them with true labels!

We can design new AI tools that are more intelligent and less artificial by collaborating with humans (experts)!

# Concluding Remarks

Visual search is gold standard; however, prone to errors, malpractice/ perceptual error etc., time consuming, and sub-optimal

Computer (AI) helps radiologists to find pathologies that can be missed however, computer also depicts so many false positives that radiologists easily capture them with true labels!

We can design new AI tools that are more intelligent and less artificial by collaborating with humans (experts)!

Explainability plays an important role in building thrust and robust systems; hence, increasing the chance of deploying such system in real clinic

# Thank you!