# Human-centric interpretability of deep learning for digital pathology

## Mara Graziani

PhD student, Hes-so Valais and UniGe

# Who am I ?

**PhD focus:**
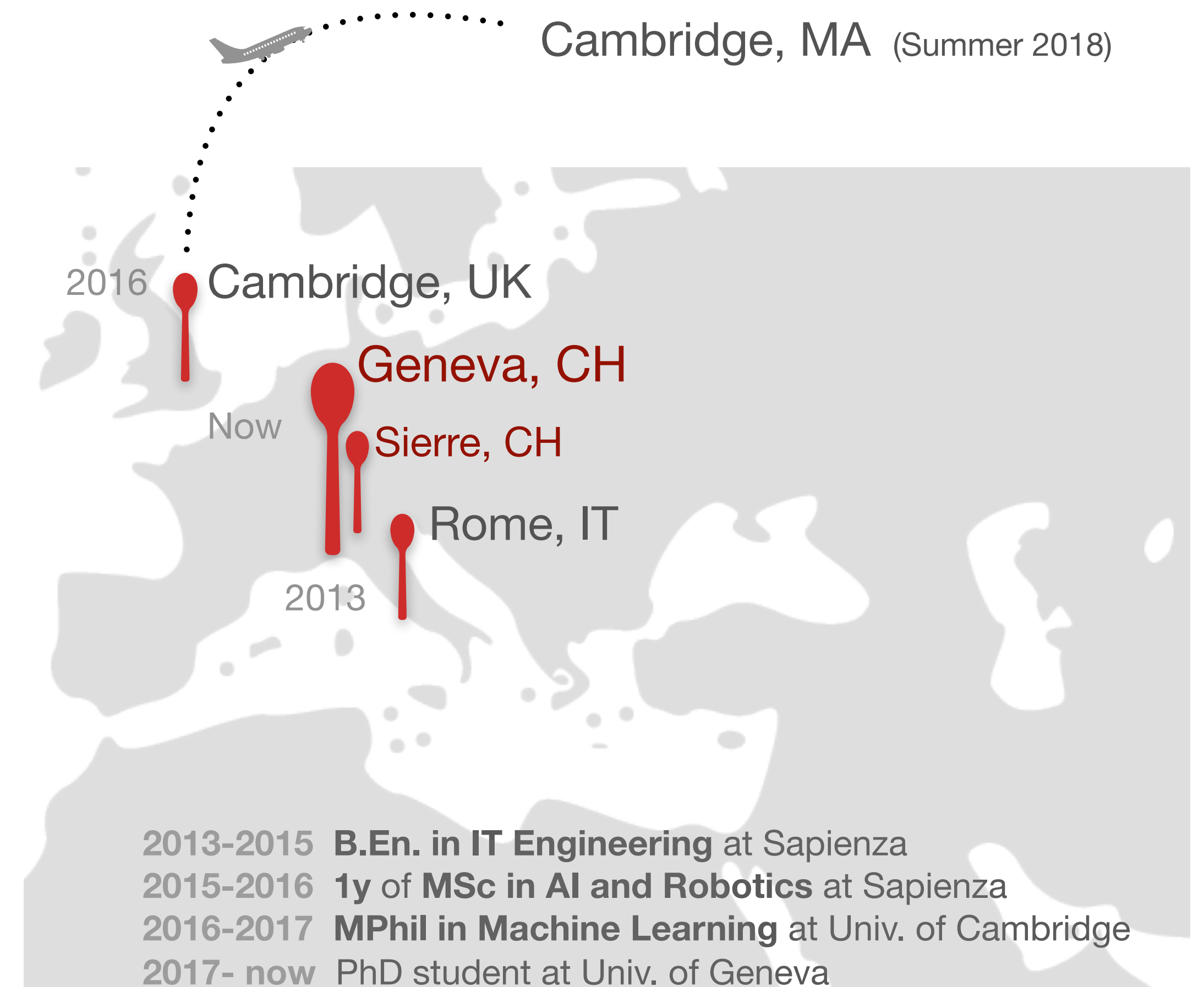
*Interpretability of Deep Learning for Medical Imaging*

**Started in:**
*November 2017*

**Funded by:**
*EU H2020 PROCESS*

mara.graziani@hevs.ch

Cambridge, MA (Summer 2018)

2016 Cambridge, UK

Geneva, CH

Now

Sierre, CH

Rome, IT

2013

2013-2015 **B.En. in IT Engineering** at Sapienza
2015-2016 **1y** of **MSc in AI and Robotics** at Sapienza
2016-2017 **MPhil in Machine Learning** at Univ. of Cambridge
2017- now PhD student at Univ. of Geneva

Can we generate human-centric explanations of deep learning and can we use them to improve model performance?

# Can we generate human-centric explanations of deep learning and can we use them to improve model performance?

## Motivation:

Ease the interaction, improve models with little extra complexity, debug models, GDPR* right for explainability, improve trust and accountability, remove bias or data memorization, generate answers to "why" questions on model behaviour and decisions.

* General Data Protection Regulation

# Outline

* Introduction and definition of **human-centric interpretability for deep learning**
* Presentation of research in this direction:
  * Evaluation of visualization tools
  * Concept-based interpretability with Regression Concept Vectors
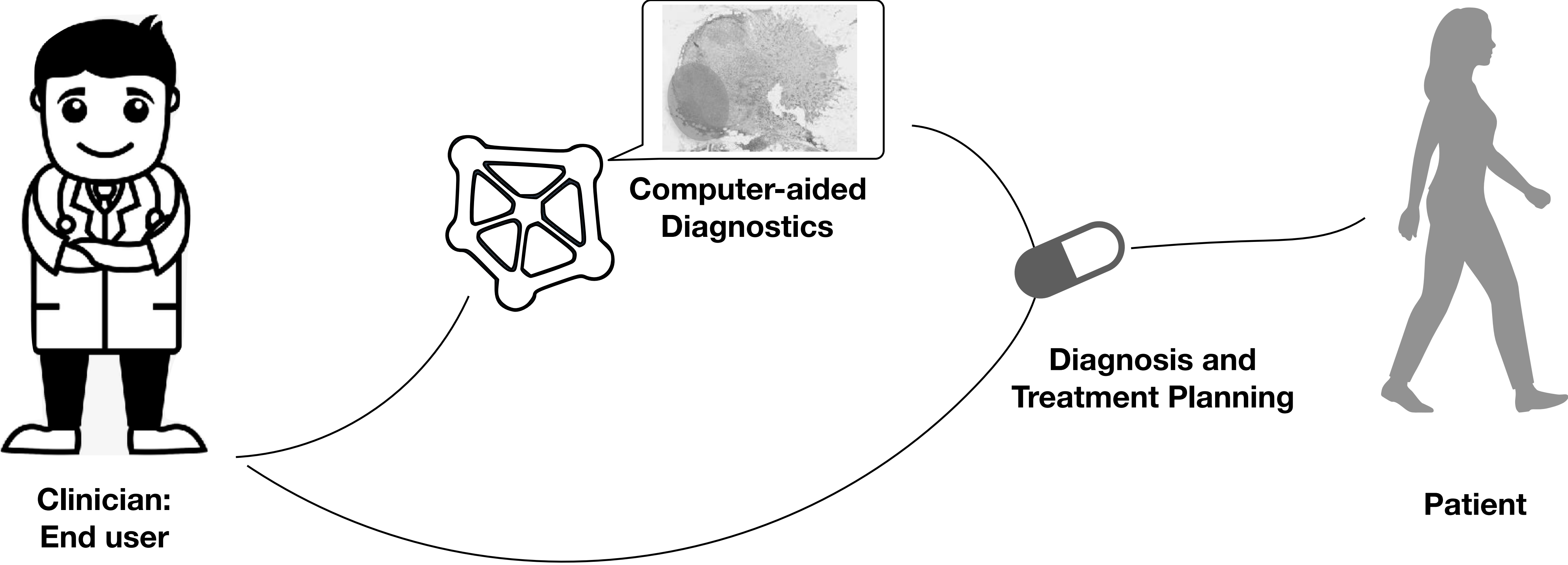  * Guiding CNNs with user-defined features
* Remarks
* Conclusions

# Interpretability: What and why?

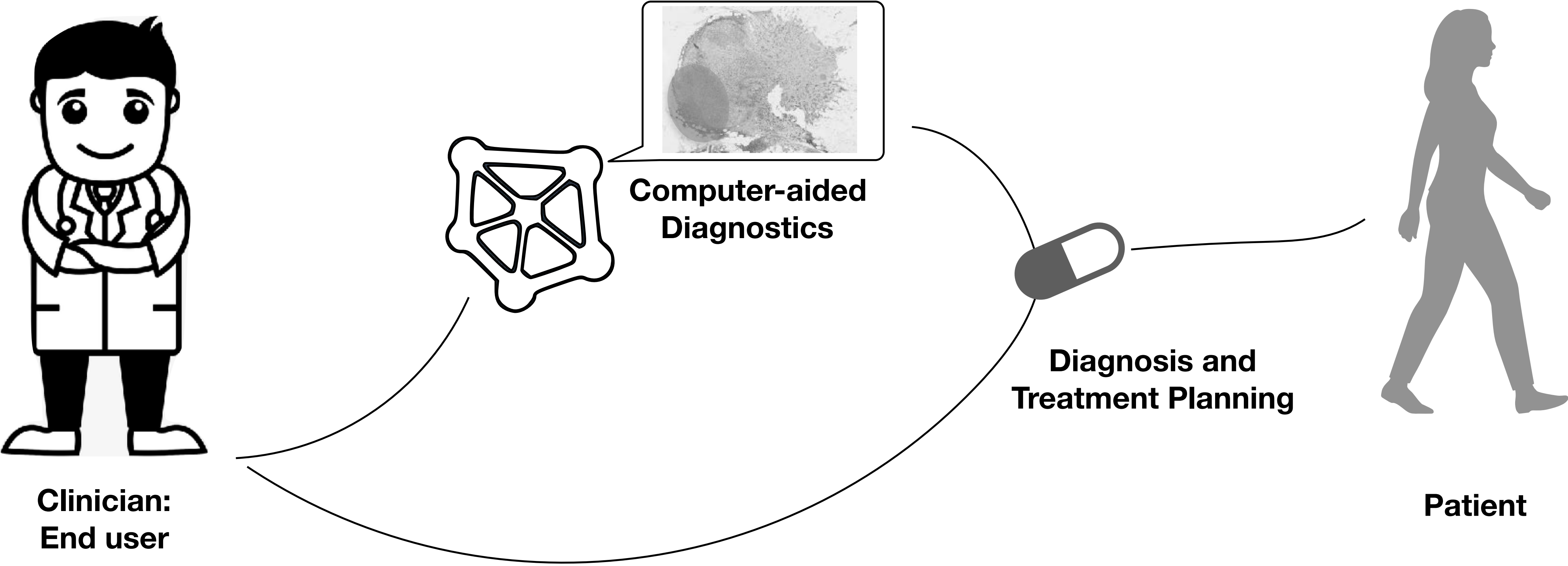"Interpretability is the ability to explain or to present in understandable terms to a human*."

[Kim et al., 2018]

* not all humans are familiar with Machine Learning

# *Human-Centric Interpretability for Cancer Diagnosis*



**Computer-aided Diagnostics**

**Diagnosis and Treatment Planning**

**Clinician: End user**

**Patient**

**Example:**
*Human-Centric Interpretability for Cancer Diagnosis*

**Computer-aided Diagnostics**

**Clinician: End user**

**Diagnosis and Treatment Planning**

**Patient**

# *1. CNN for localization*

CNN

This is a high-grade tumor region!

**Clinician: End user**

## *2. CNN for localization with explanations of abnormalities*
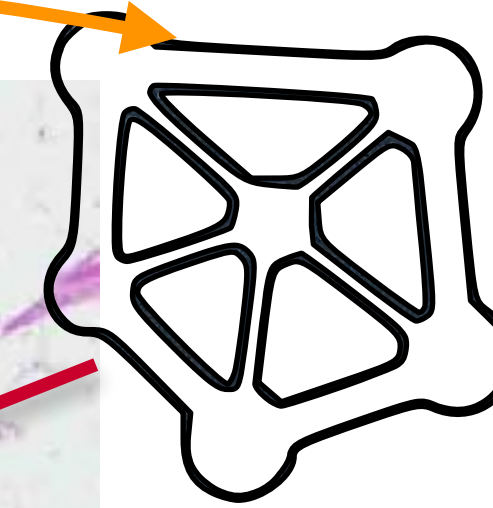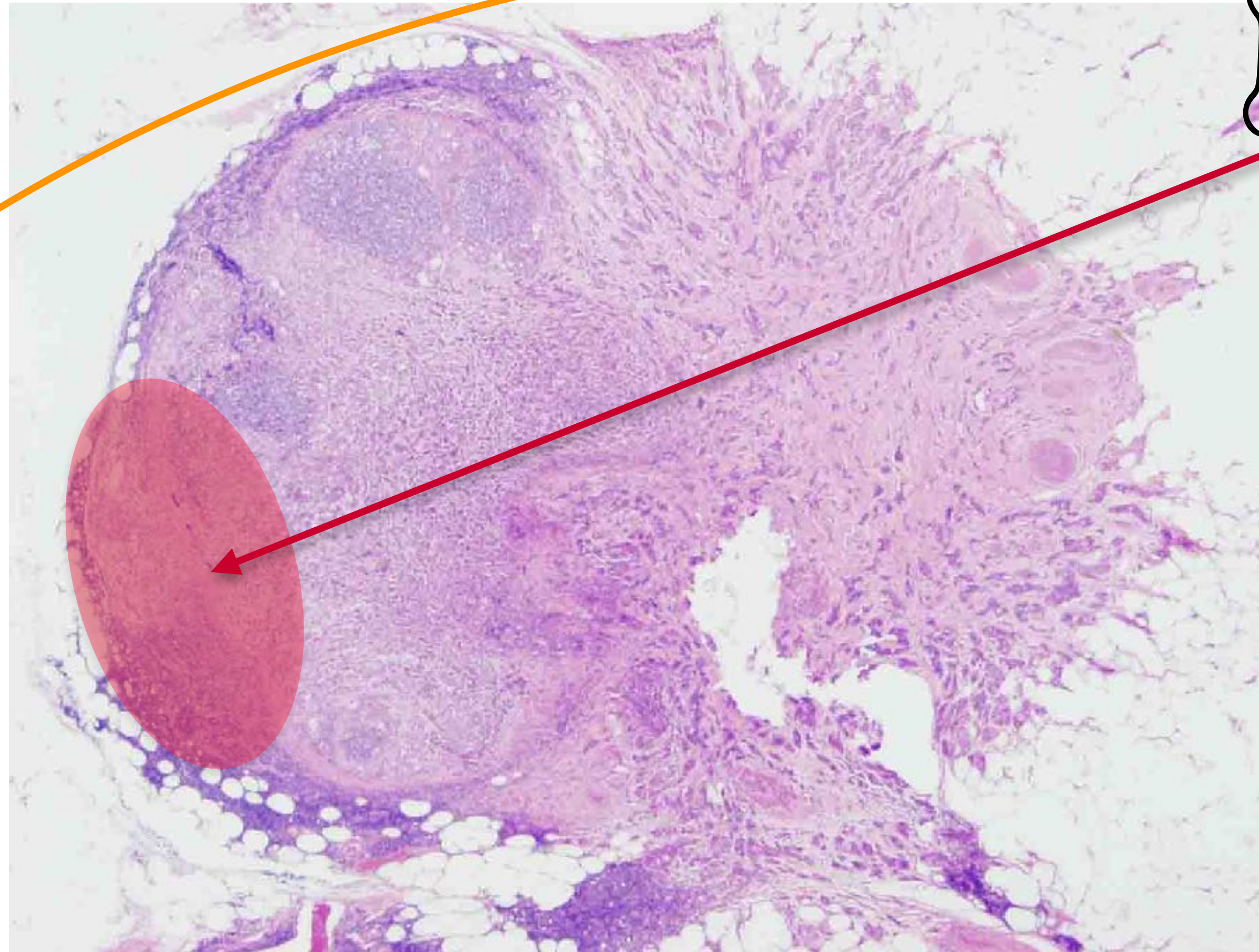


CNN
With
Interpretability

This is a high-grade tumor region:

1. The nuclei are 30% larger than non-tumor average
2. The nuclei texture appears vesicular (contrast is 40% larger than average)

**Clinician: End user**

# 3. *CNN for localization with explanations of abnormalities and with guided feature learning by user-input*



**Guided** CNN
With
Interpretability

**Clinician:
End user**
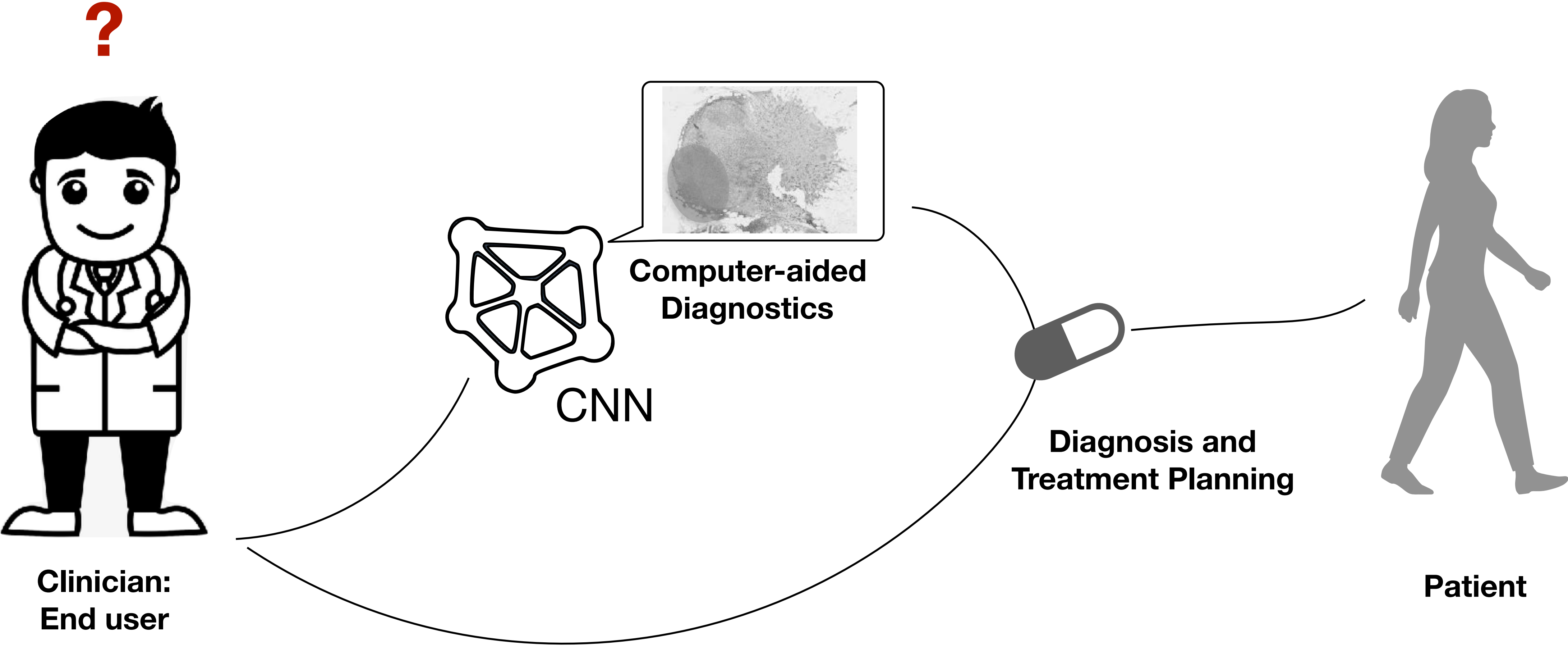
This is a high-grade tumor region:

1. The cells are 30% larger than non-tumor average
2. The nuclei texture appears vesicular (contrast is 40% larger than average)
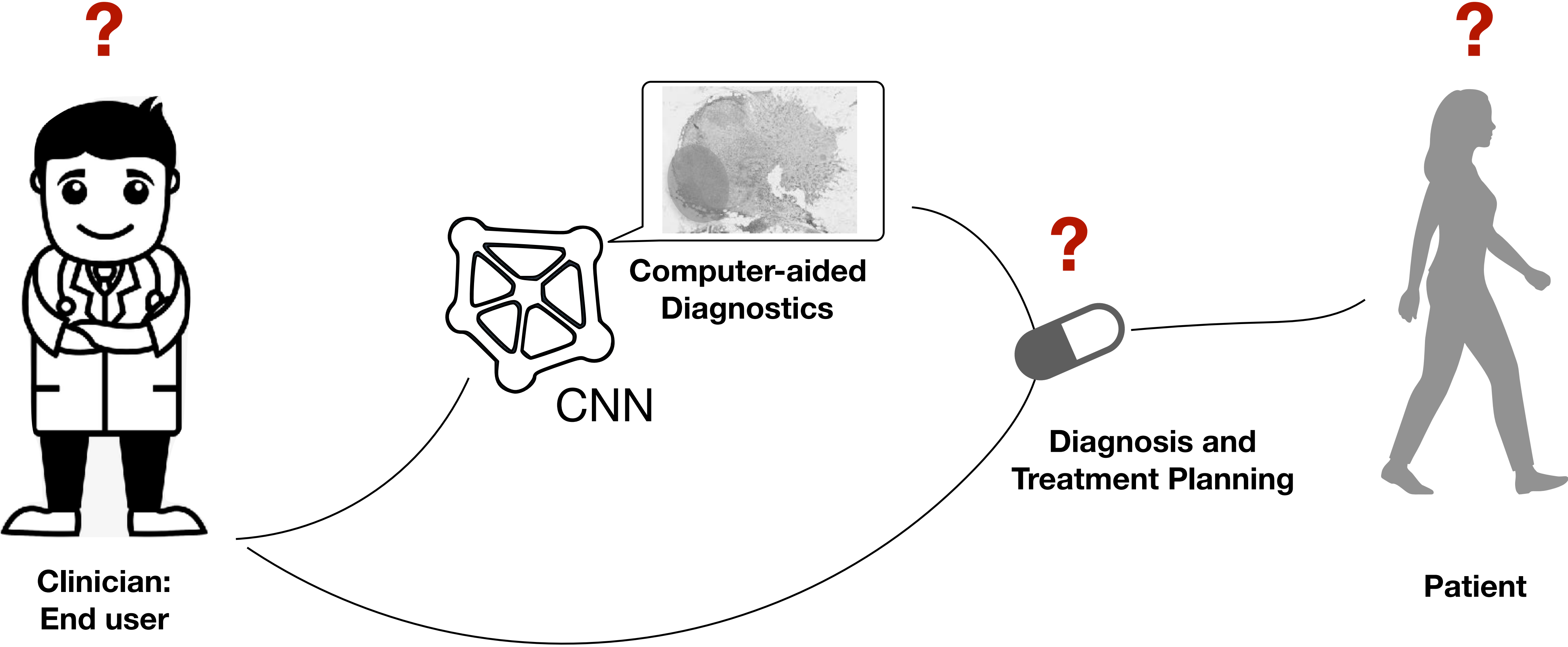
# Let's analyse the three scenarios

# *1. CNN for localization*

# 1. CNN for localization
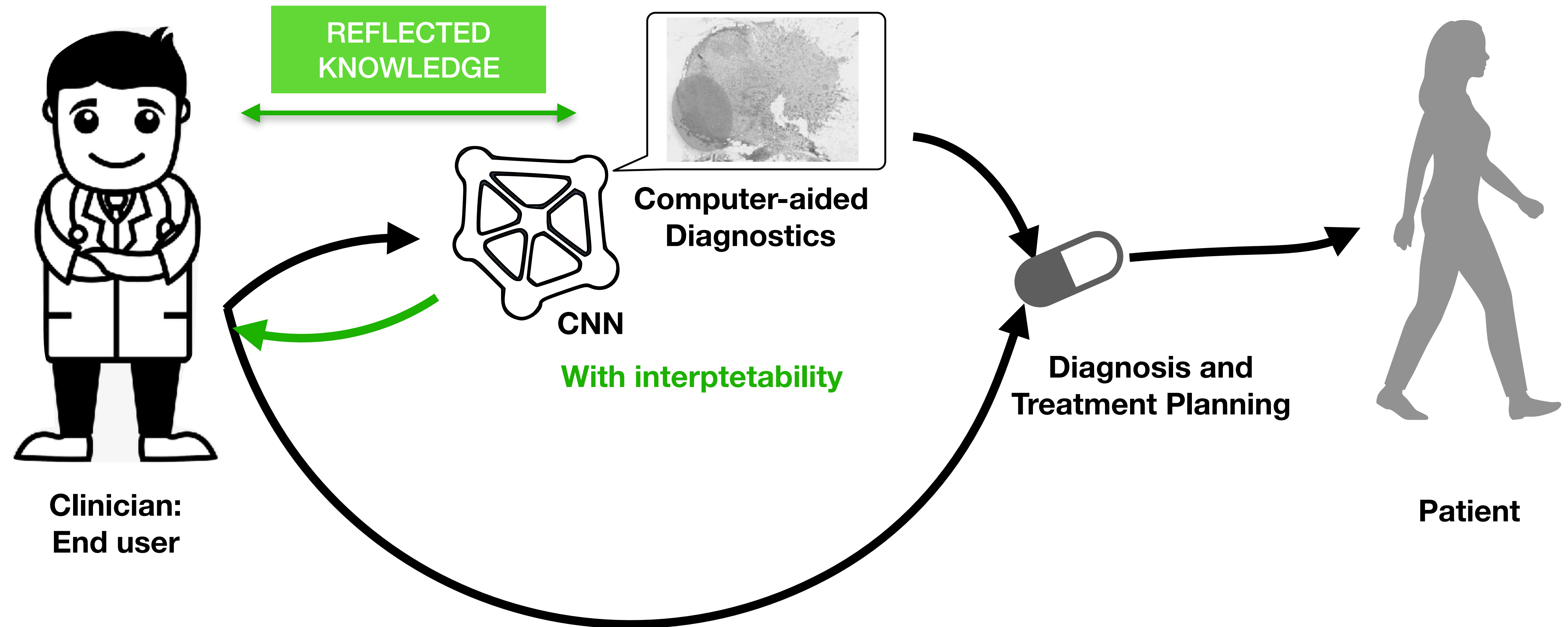


**Computer-aided Diagnostics**

CNN

**Diagnosis and Treatment Planning**

**Clinician: End user**

**Patient**

# 2. CNN for localization with explanations of abnormalities



REFLECTED KNOWLEDGE

Computer-aided Diagnostics

CNN

With interptetability

Clinician: End user

Diagnosis and Treatment Planning

Patient

– Explaining the decisions of a complex model in understandable terms by doctors eases the interaction with AI and improves the quality of the diagnosis [Carrie J.C. et al., 2019].

# 3. CNN for localization with explanations of abnormalities and with guided feature learning by user-input

**Control Targets:**
– Nuclei size is **relevant**
– Image domain is **not relevant**

REFLECTED KNOWLEDGE

*Control targets*

Computer-aided Diagnostics

*Guided* CNN
**With interptetability**

Clinician: End user

Diagnosis and Treatment Planning

Patient

Hes·so// VALAIS WALLIS

# To summarize

* CNNs for tumor localization can support pathologists in the diagnosis, but may leave them with **unanswered questions about the output** (scenario 1)
* **Interpretability should help** the clinician verify that the CNN decision making respects the guidelines and knowledge in the domain (scenario 2).
* The expertise of clinicians is a valuable input for the network training, that **could** be **guided** to ensure that certain visual features are taken into account and others are not (scenario 3).

# Human-centric DL interpretability

**=**

A tool that **supports** the pathologists in making decisions by providing **explanations** and allowing the introduction of **feedback** to refine training

# Our work in this direction

* **Evaluation of visualization methods for histopathology**
* Concept-based interpretability of CNNs
* Guidable CNNs

**Feature-attribution:**
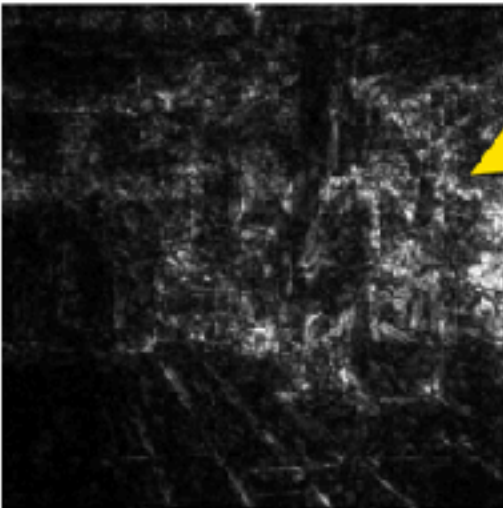*Evaluation of visualization tools*

class
ATTENTION
occlusion
deconvolution
values lime network
NETWORK Shapley
relevant maximisation
saliency
feature
dissection
attention maps
activation
propagation
integrated

**e.g saliency**

$$\frac{\partial output}{\partial input}$$

One of the most popular interpretability methods for images:

Saliency maps

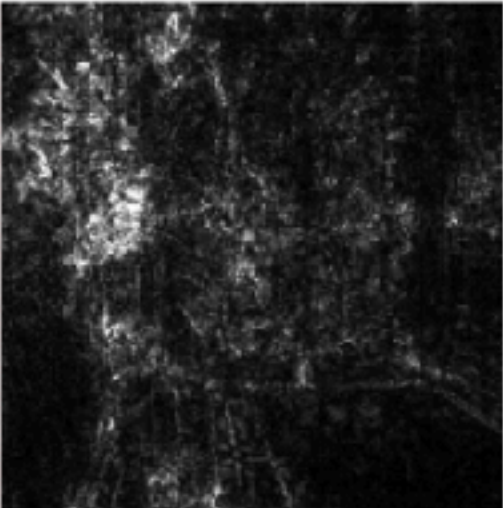a logit $\rightarrow \frac{\partial p(z)}{\partial x_{i,j}}$
pixel i,j $\rightarrow$

prediction:
Cash machine

Why correct?
Why incorrect?

prediction:
Sliding door
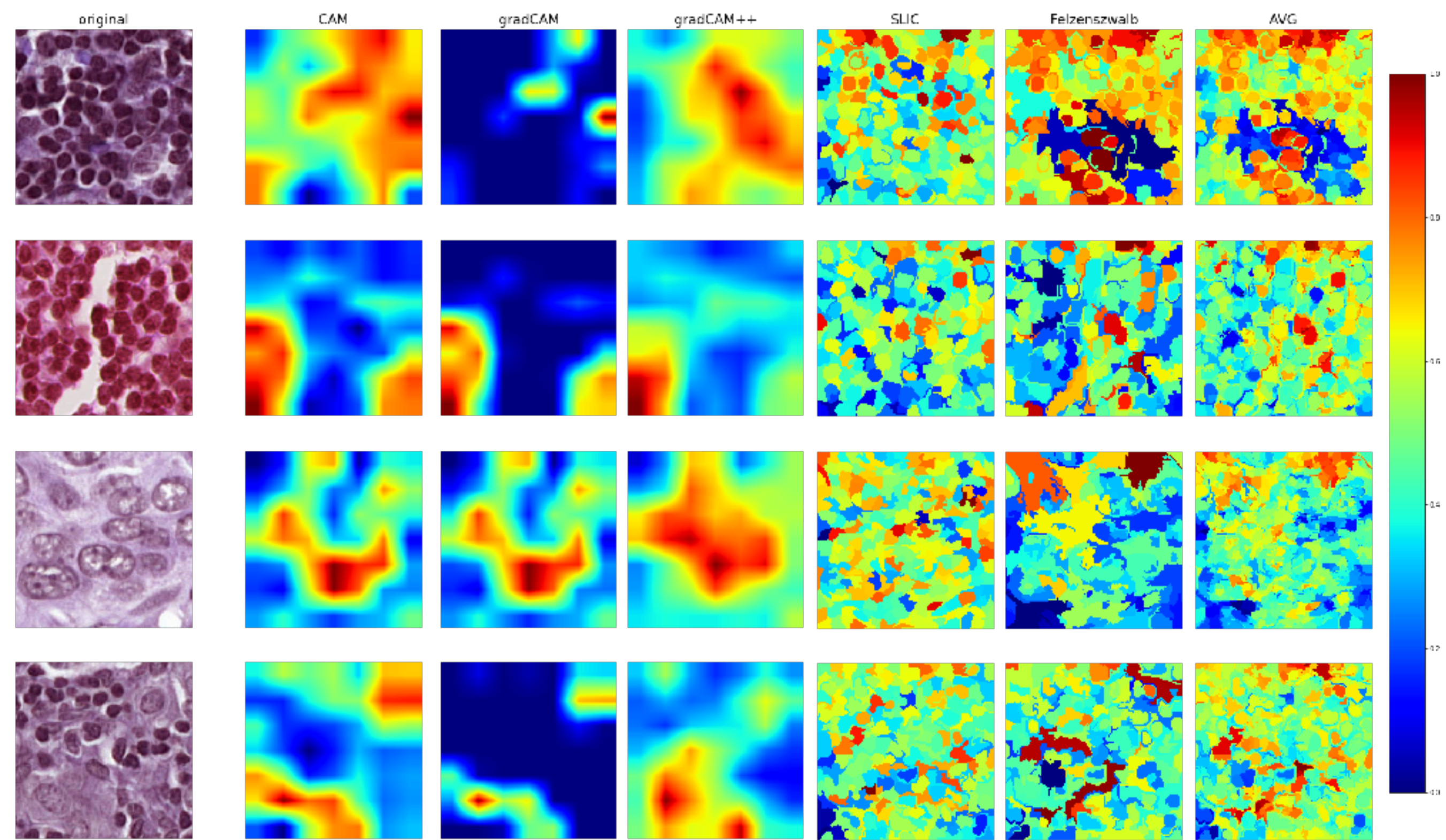
Slide credits: B. Kim

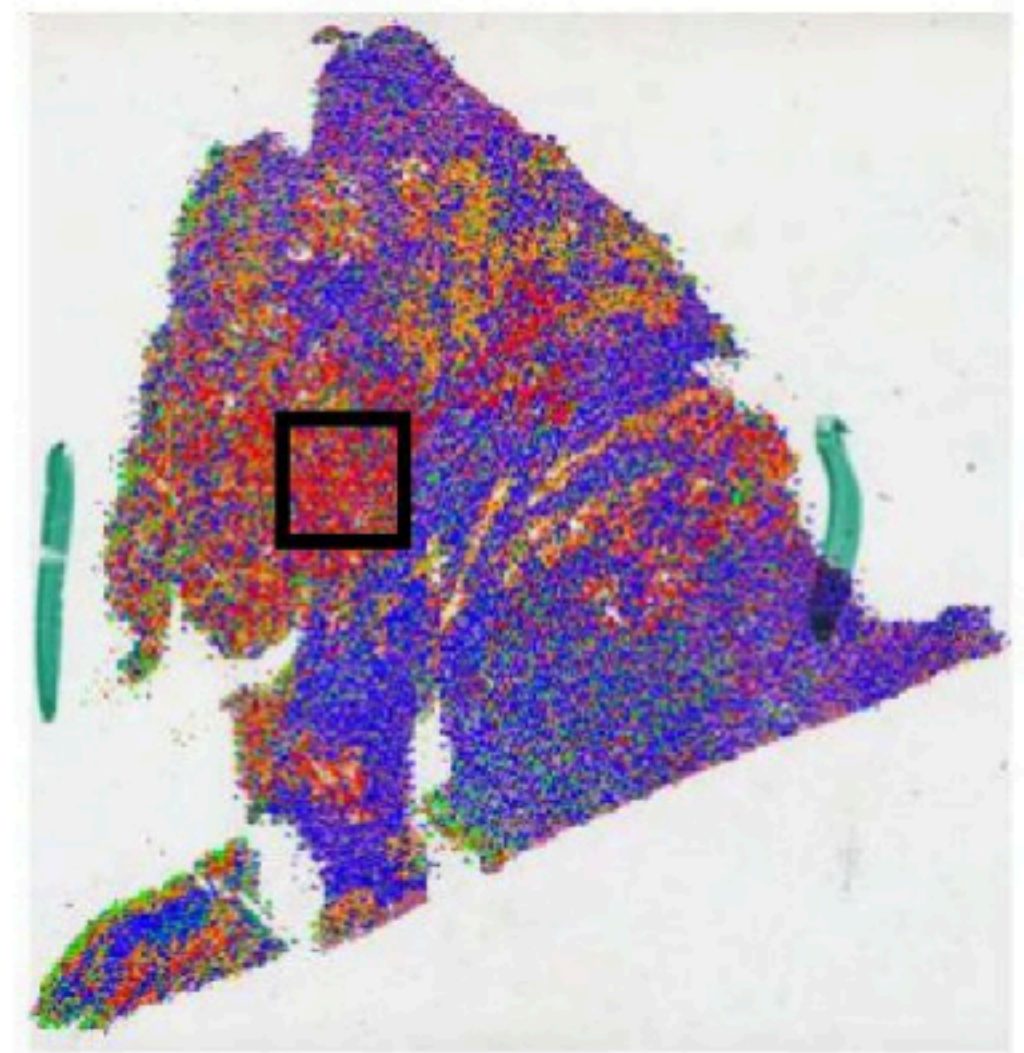**Issues:**
– Difficult Abstraction
– Sometimes Ambiguous [Kim et al., 2018]
– Consistency issues [Adebayo et al., 2018]

Hes·so// VALAIS WALLIS

# *Evaluation of visualization tools*
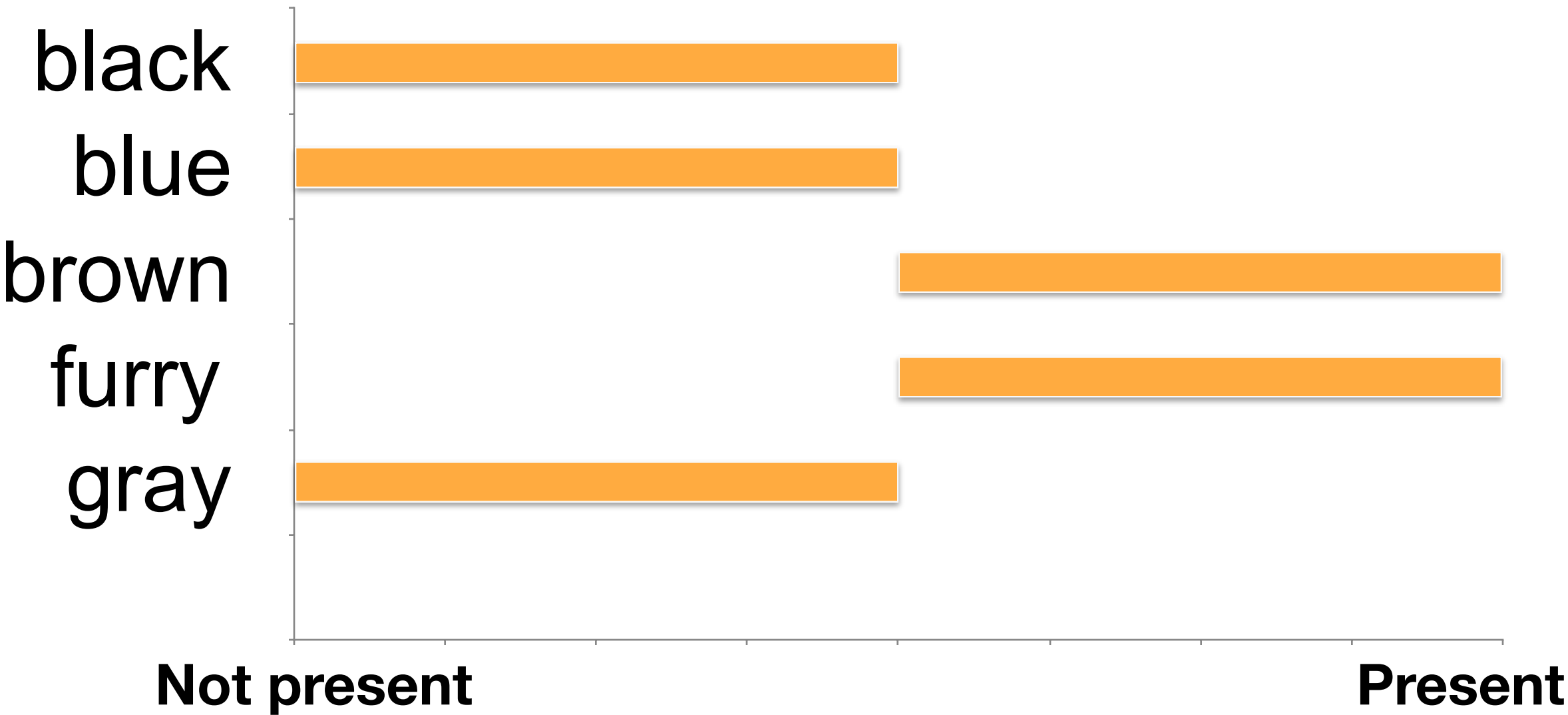


https://jgamper.github.io/PanNukeDataset/

# Our work in this direction

* Evaluation of visualization methods for histopathology
* **Concept-based interpretability of CNNs**
* Guidable CNNs

# *Concept attribution with Regression Concept Vectors*

Taking inspiration from [Kim et al., 2018] on interpreting CNN activations with human-friendly binary concepts (presence vs absence).

# *Concept attribution with Regression Concept Vectors*
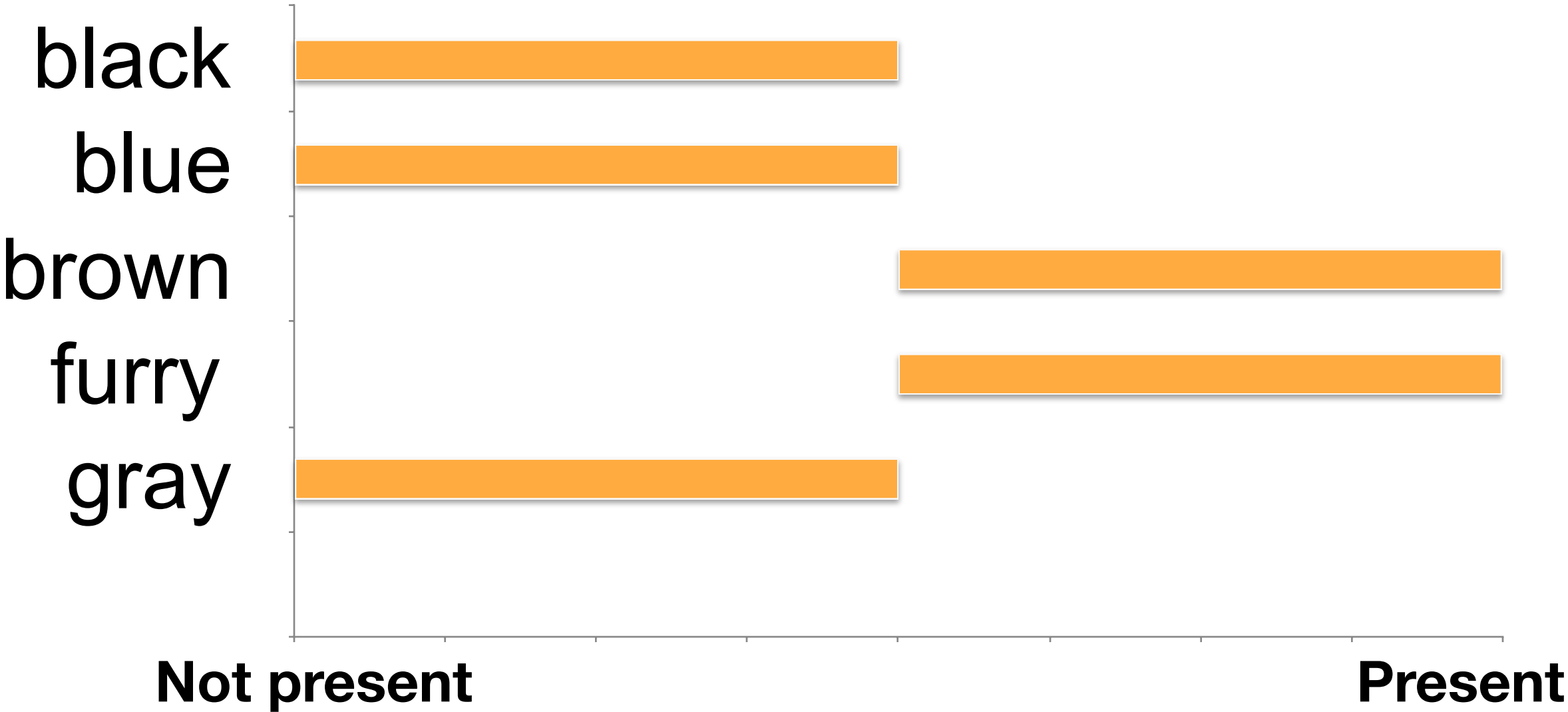
Taking inspiration from [Kim et al., 2018] on interpreting CNN activations with human-friendly binary concepts (presence vs absence).
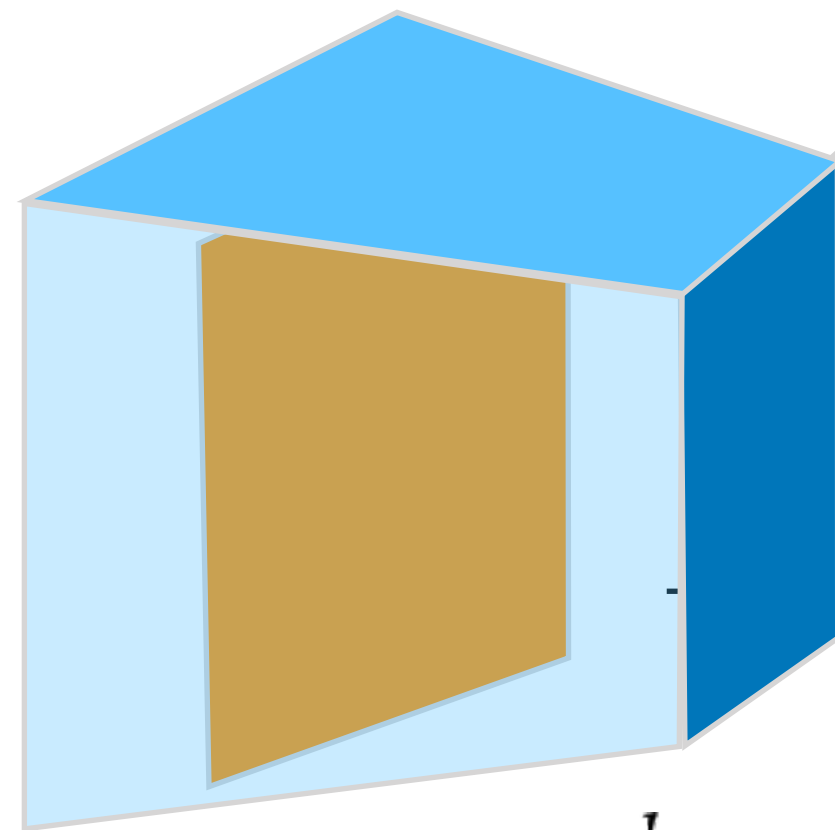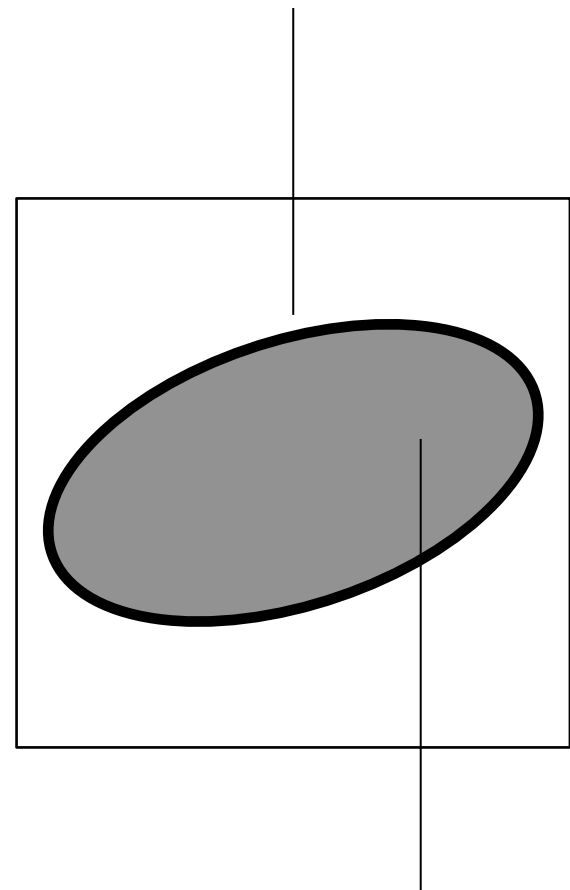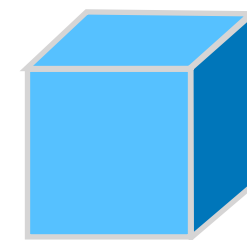


Measuring texture "**coarseness**", "**brown-ness**"

# *Concept attribution with Regression Concept Vectors**

Segmentation
(manual or
automatic)



$$\Phi^l(\mathbf{x_j})$$

Handcrafted
features, texture
descriptors, shape,
size, ...

*Concept attribution with Regression Concept Vectors**

Size of the ball = concept value corresponding to one input image

Segmentation
(manual or
automatic)

$\Phi^l(\mathbf{x_j})$

$\Phi_2$

$\vec{v_C}$

Vector of "size"

$\Phi$

Handcrafted
features, texture
descriptors, shape,
size, ...

Take the internal
activations (aggregation)

Linear regression of
measures

# *Concept attribution with Regression Concept Vectors**



Size of the ball = concept value corresponding to one input image
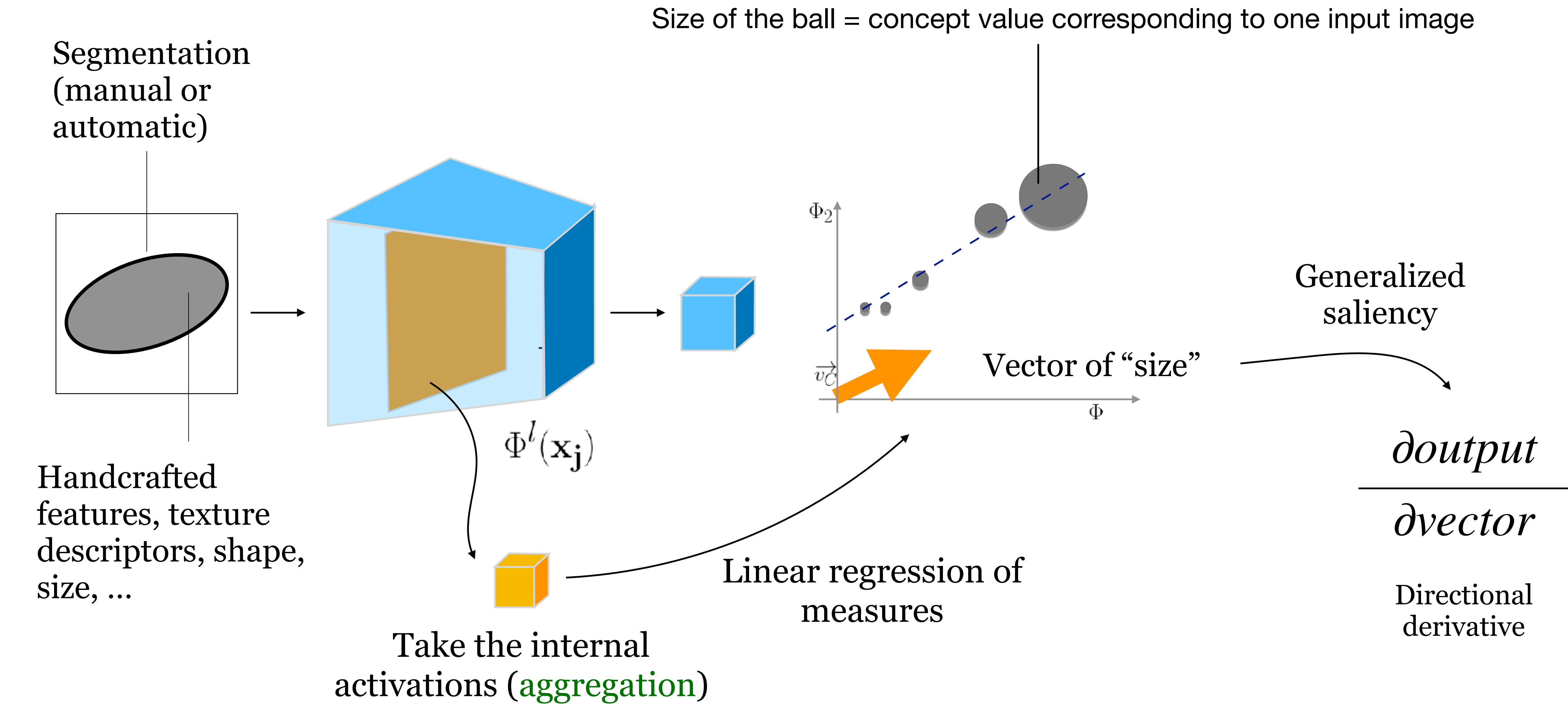
Segmentation
(manual or
automatic)

$\Phi^l(\mathbf{x_j})$

Handcrafted
features, texture
descriptors, shape,
size, ...

Take the internal
activations (aggregation)

Linear regression of
measures

Vector of "size"

$\vec{v_C}$

$\Phi_2$

$\Phi$

Generalized
saliency

$$\frac{\partial output}{\partial vector}$$

Directional
derivative

# *Regression Concept Vectors: application to histopathology*

# Remarks

* **Interpretability can be used to** verify that the CNN decision making respects clinical guidelines and knowledge in the domain
* **Visualizations of saliency heatmaps** give feedback on the relevant input pixels, while **concept-based** explanations use directly clinically relevant measures such as nuclei size and appearance.
* The expertise of clinicians can be used to **guide network training** by the combination of multitask and adversarial learning.

# Remarks

* **Interpretability can be used to** verify that the CNN decision making respects clinical guidelines and knowledge in the domain
* **Visualizations of saliency heatmaps** give feedback on the relevant input pixels, while **concept-based** explanations use directly clinically relevant measures such as nuclei size and appearance.
* The expertise of clinicians can be used to **guide network training** by the combination of multitask and adversarial learning.

# Q&A