# AI-based assessment
# of cardiac allograft rejections
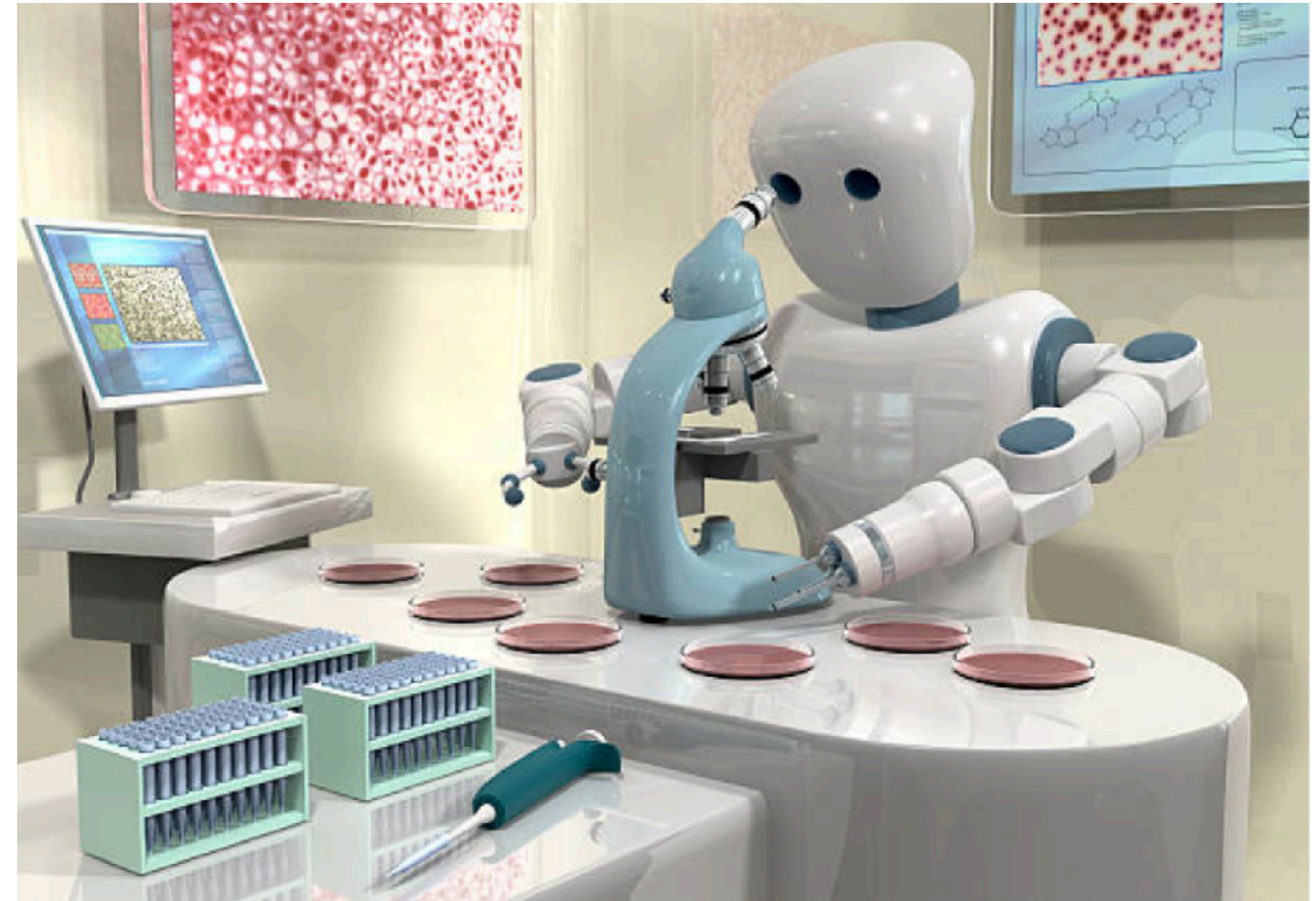


Jana Lipkova

# Overview

▶ **Background:**

- Histopathology data

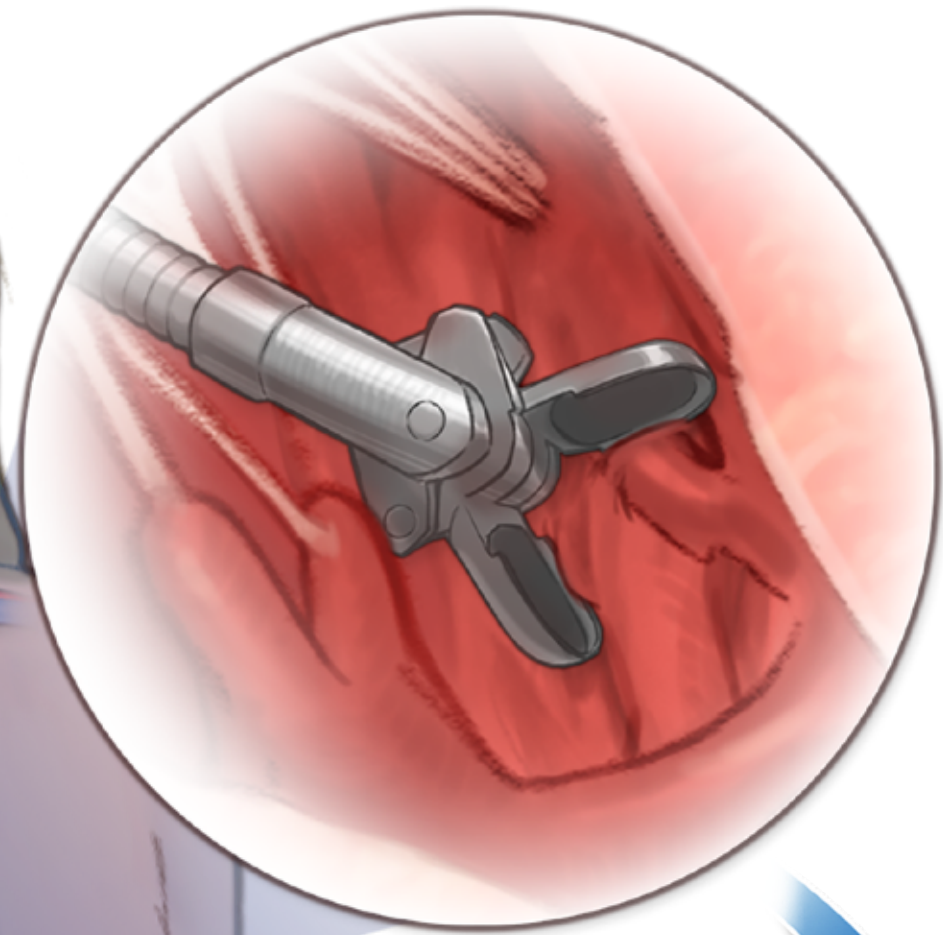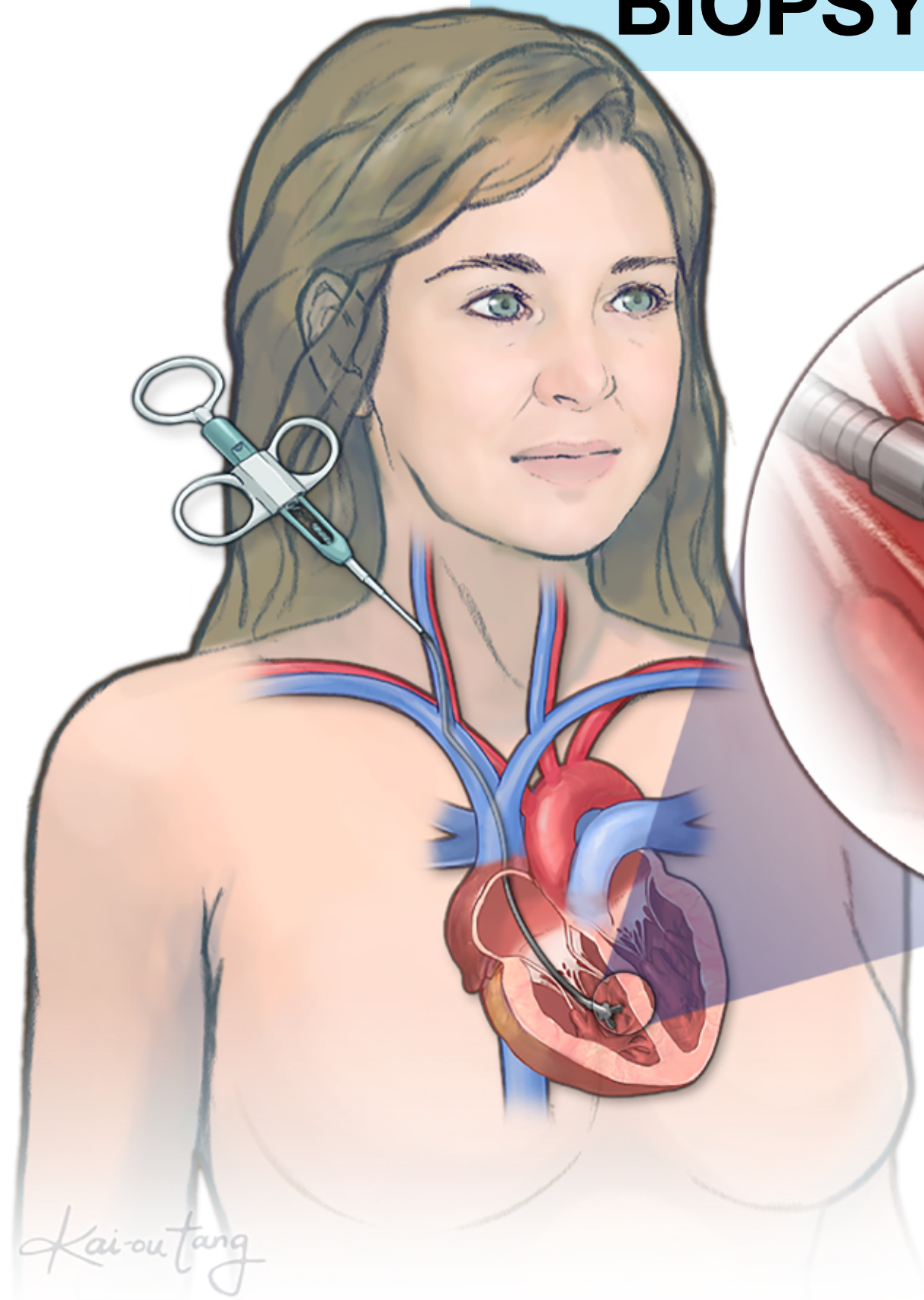- Cardiac Allograft Rejections
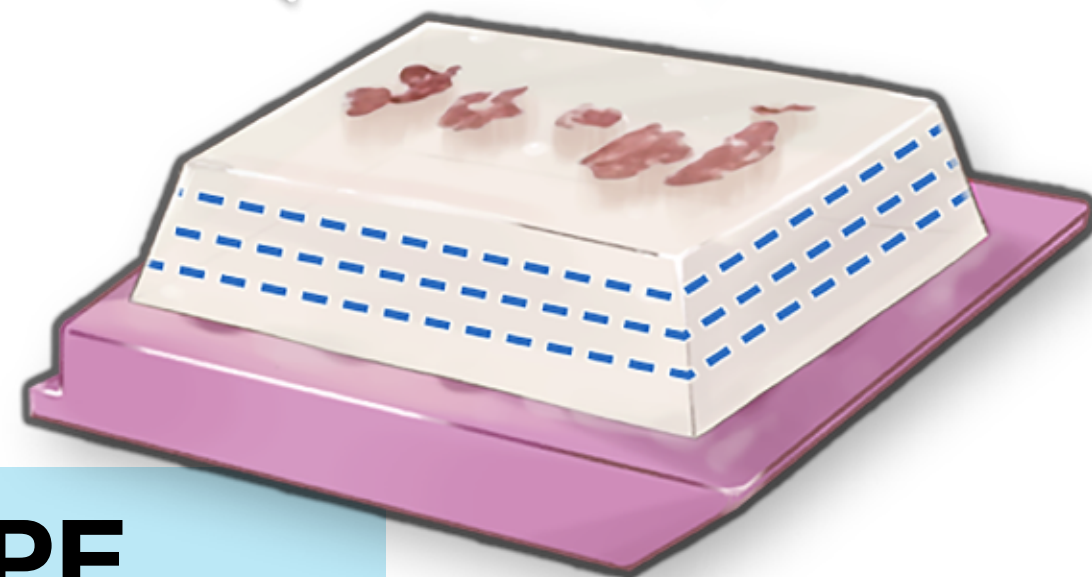
▶ **AI-based assessment of allograft rejections**

CRANE

# Histology 101



**BIOPSY**

BIOPSY

**WSI** *Whole-Slide Image*

WSI

Level 1
Level 1 2
Level 2 3
Level 3

Hematoxylin & Eosin (H&E) stain

DIAGNOSIS

\#
\#
Scanner
Magnification
Micron/pixel

**FFPE**

FFPE

*(Formalin-fixated, paraffin embedded)*

**Hematoxylin**: stains cell nuclei
**Eosin**: the extracellular matrix and cytoplasm

# Cancer detection/classification 101

## Normal tissue (kidney)



▸ Symmetric regular structure
▸ One nuclei per cell
▸ Cell/nuclei - regular shape

## Chromophobe renal carcinoma



▸ Enlarged nuclei
▸ Double nuclei per cell
▸ Irregular shape

## Papillary renal carcinoma



▸ Papillary cores lined by neoplsatic cells
▸ Tubulopapillary architecture

# Fun fact: Also pigeons can detect cancer



*R. Levenson et al: Pigeons (Columba livia) as Trainable Observers of Pathology and Radiology Breast Cancer Images, PloS one, (2015)*

# Digital Pathology: Whole Slide Images (WSIs)

- High resolution scan of an entire tissue section (0.25 - 0.5 microns per pixel)
- Gigapixel image: 100,000 x 100,000 pixels
- **100 WSI have cca same amount of pixels as whole ImageNet**
- Different Stains: H&E, IHC

# Medical Data

## Radiology



*(MRI head scan)*

- 3D images
- gray-scale images
- resolution: ~1 mm
- size: 256x256x256 voxels

## Photography



*(Fundus / skin photography)*

- 2D images
- RGB
- 10 µm - 1 mm
- size: ~1,700x1,700 pixels

## Histology



*(H&E tissue)*

- 2D images
- RGB
- scale: ~0.1µm
- 100,000x100,000 pixels

(varies with magnification, tissue size etc)

## Genomics



- 1D array
- float (e.g. '0' wild type, '1' mutation)
- scale: 1µm - 1nm
- ~20,000 protein-coding genes

Lipkova et al. Nature Medicine (2022)

# BACKGROUND

**Heart Failure**

**Heart Transplant**

**Immune Response**

**Allograft Rejection**



- Leading cause of hospitalization in USA/EU
- 26 million cases / year

- Patients with end-stage failure
- 5000 transplants / year

- Immunosuppressives
- Patient-specific set-up

- Main complication & main cause of death
- 40% recipients

# MOTIVATION

**APPLICATION:**

▸ Early stages of rejections are **asymptomatic** → surveillance **Endomyocardial biopsy** (EMB)
▸ Gold-standard: manual assessment H&E-stained biopsies:
  ▸ **detection** and **subtyping** of rejections (*acute cellular*, *antibody-mediate*, *benign mimickers*) and **grading** (I-III)
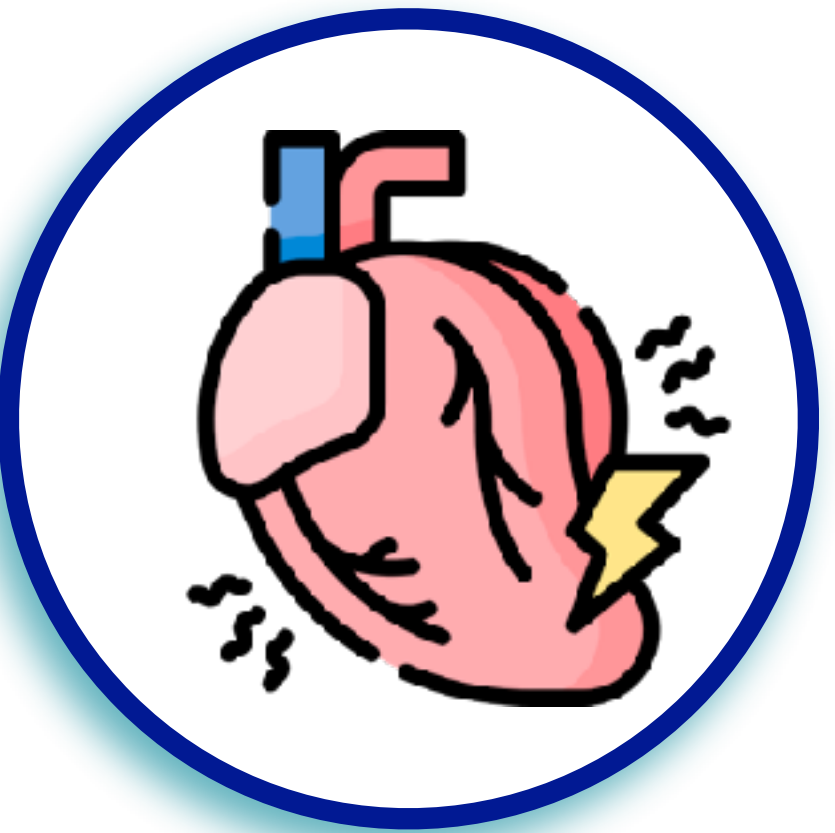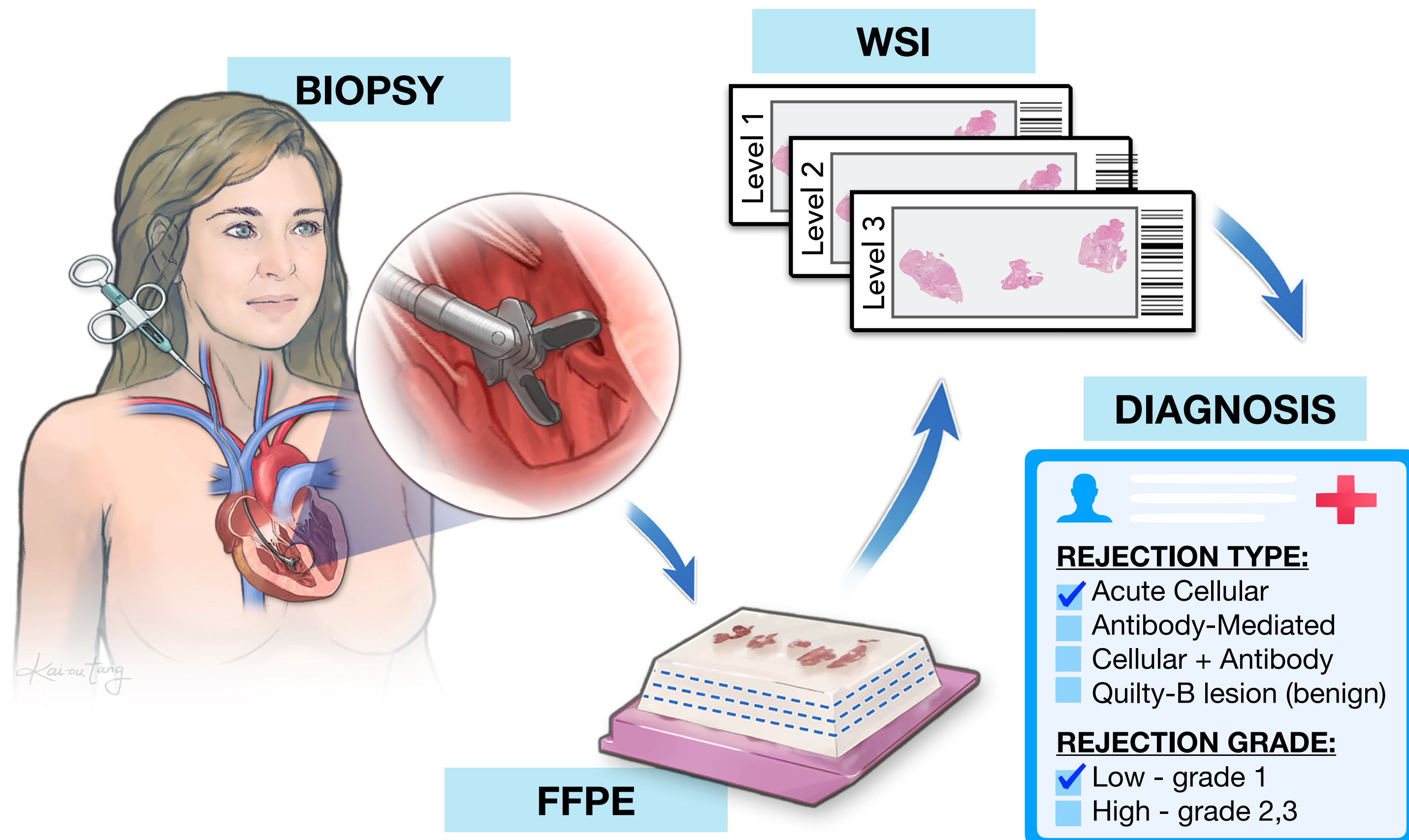▸ Rejection type & grade **determines the immunosuppressive treatment regime**



**BIOPSY**

**WSI**

Level 1
Level 2
Level 3

**FFPE**

**DIAGNOSIS**

**REJECTION TYPE:**
☑ Acute Cellular
☐ Antibody-Mediated
☐ Cellular + Antibody
☐ Quilty-B lesion (benign)

**REJECTION GRADE:**
☑ Low - grade 1
☐ High - grade 2,3

**TRAIN**

| # | 4059 |
|---|---|
| # | 1354 |
| Scanner | **HAMAMATSU** |
| Magnification | 40x |
| Micron/pixel | 0.2206 |

**TEST**

| # | |
|---|---|
| # | |
| Scanner | **HA** |
| Magnification | |
| Micron/pixel | |

| # | |
|---|---|
| # | |
| Scanner | *Le* |
| Magnification | |
| Micron/pixel | |

| # | |
|---|---|
| # | |
| Scanner | **3D** |
| Magnification | |
| Micron/pixel | |

*Lipkova et al. Nature Medicine (2022)*

▸ [1] Concordance among pathologists in the second cardiac allograft rejection gene expression observational study (CARGO II) In: *Transplantation* 94.11 (2012), pp. 1172–1177

# 101: Rejection Types

**Normal tissue**

**Abnormal tissue**



## Acute Cellular

- Lymphocyte infiltrates in muscle tissue
- Homogenous structure
- Comprised of T-cells

## Antibody Mediated

- Increased extracellular space + **edema**
- **Capillary** endothelial **changes**
- Increase cell damage
- More macrophages and necrosis

## Quilty B Lesions

- Benign lesions
- Mixed B and T-cells, macrophages and plasma cells
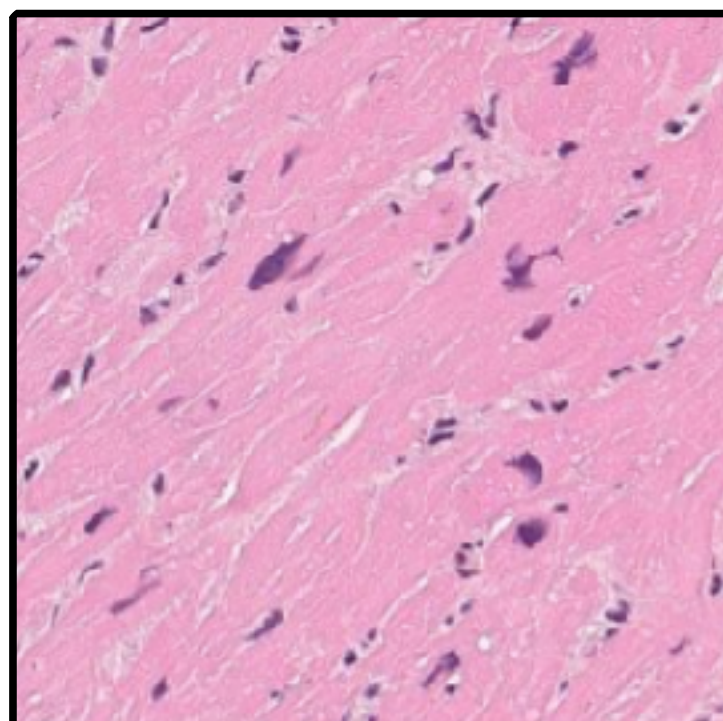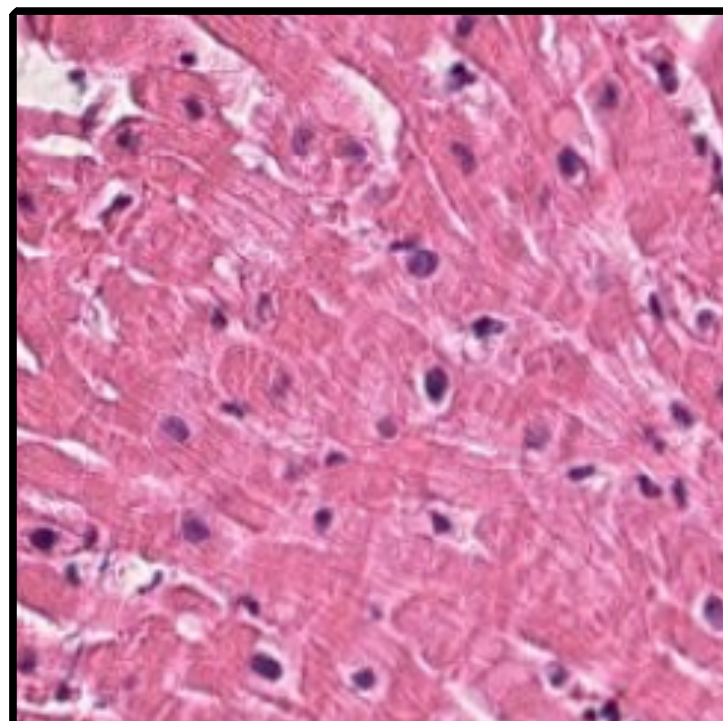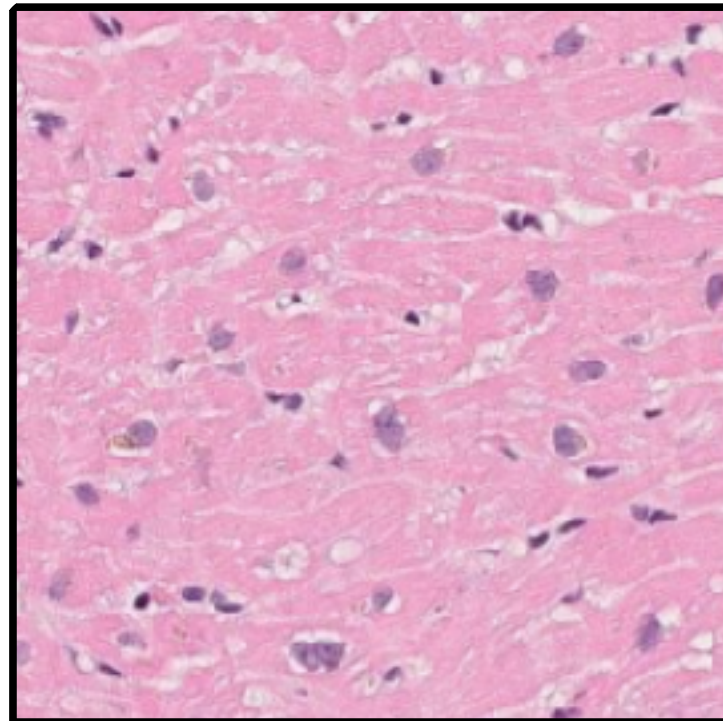- commonly mistaken for cellular rejections

# MOTIVATION

- Early stages of rejections are **asymptomatic** → surveillance **Endomyocardial biopsy** (EMB)
- Gold-standard: manual assessment H&E-stained biopsies:
  - **detection** and **subtyping** of rejections (*acute cellular*, *antibody-mediate*, *benign mimickers*) and **grading** (I-III)
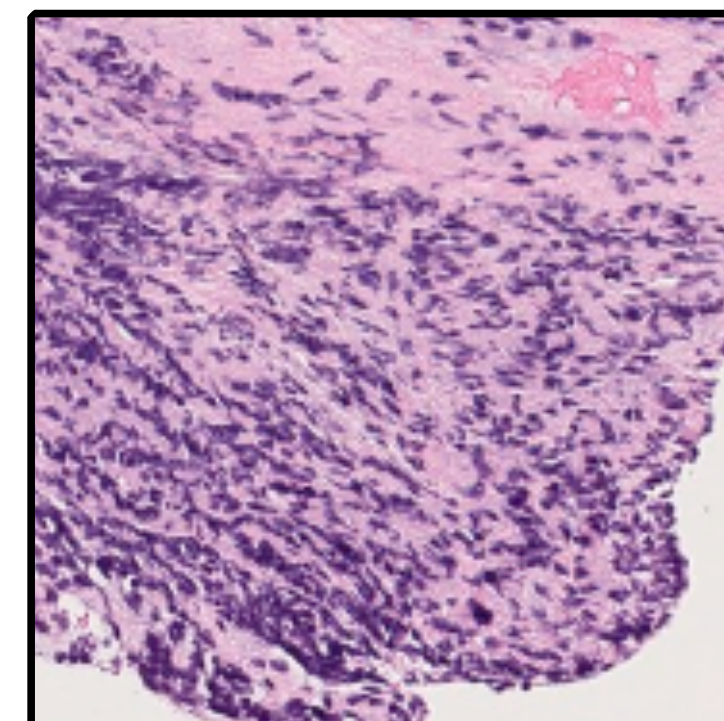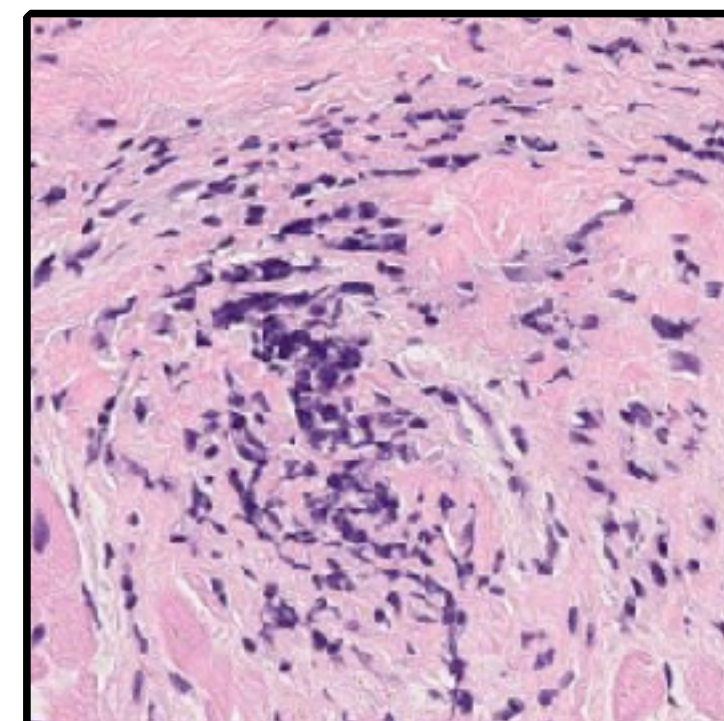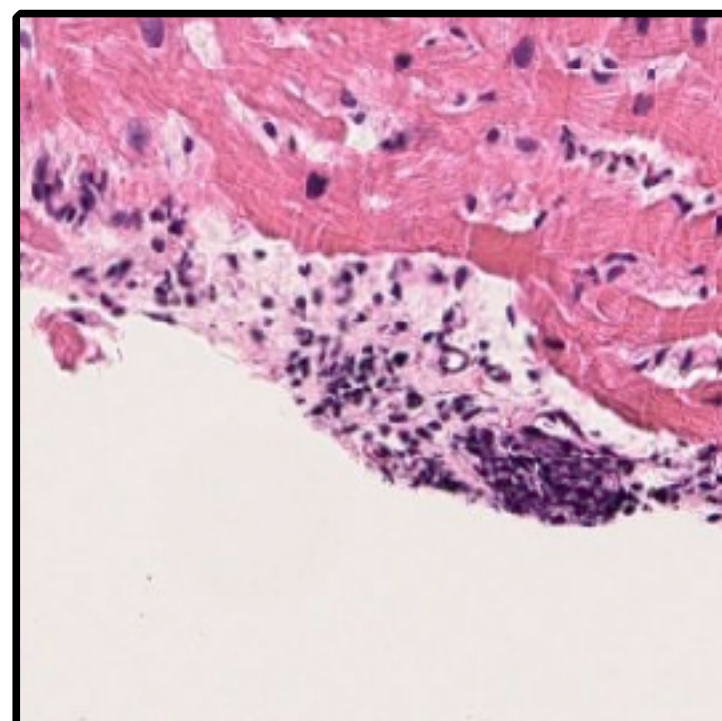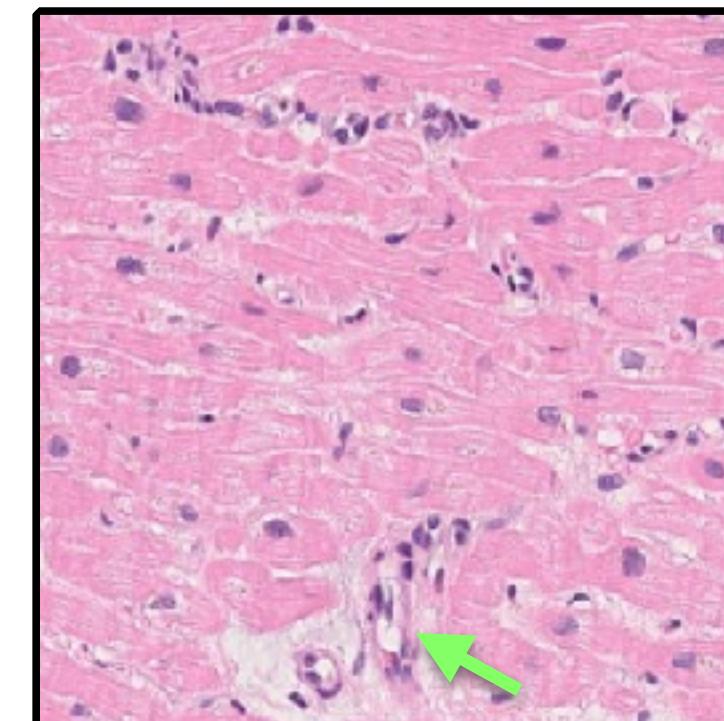- Rejection type & grade determines the immunosuppressive treatment regime

**BIOPSY**

**WSI**

Level 1
Level 2
Level 3

**DIAGNOSIS**

**REJECTION TYPE:**
- ✔ Acute Cellular
- ☐ Antibody-Mediated
- ☐ Cellular + Antibody
- ☐ Quilty-B lesion (benign)

**REJECTION GRADE:**
- ✔ Low - grade 1
- ☐ High - grade 2,3

**FFPE**
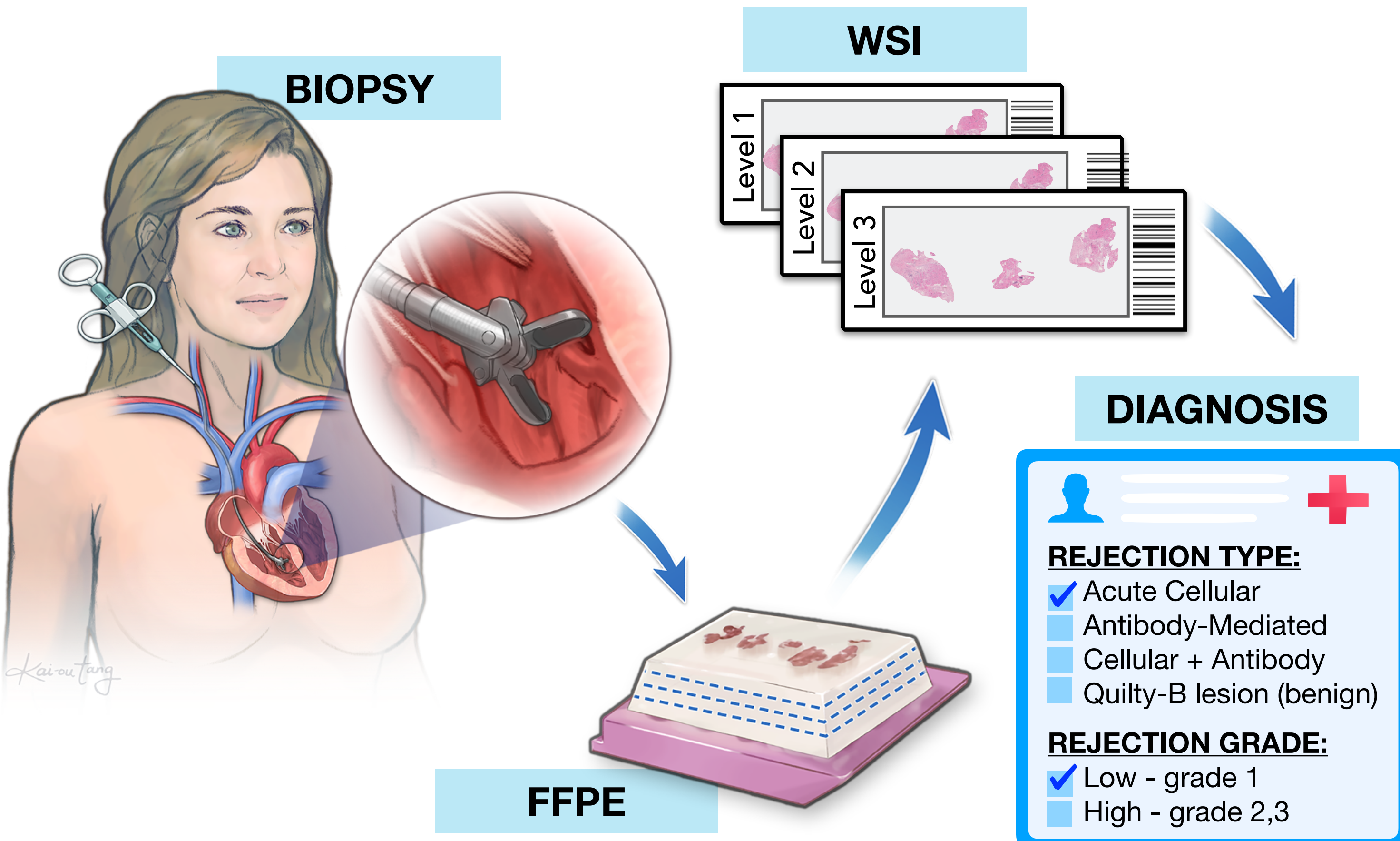
## CHALLENGES:

- **Substantial inter-rater variability** [1]:
  - <71 % agree if recipient is rejecting the heart
  - <28 % agree on the grade of advance rejections
  - 19 % unable to reach majority agreement

- **Misinterpretation**:
  - under/over treatment with immunosuppressives
  - unnecessary follow-up biopsies
  - worse outcomes

| Scanner | HAMAMATSU |
| --- | --- |
| Magnification | 40x |
| Micron/pixel | 0.2206 |

## AIM:

- **Objective and automated EMBs assessment**

**TEST**

tion
kel

tion

Micron/pixel

\#
\#

Magnification
Micron/pixel

[1] Concordance among pathologists in the second cardiac allograft rejection gene expression observational study (CARGO II) In: *Transplantation* 94.11 (2012), pp. 1172–1177

**PREDICTIONS**

# Cardiac Rejection Assessment Neural Estimator

CRANE

- ▸ **Input:** H&E-stained EMBs whole-slide-images (WSIs)
- ▸ **Multi-task, multi-label model:** simultaneously identifies **presence** and **type of the rejection** (cellular, antibody, and/or quilty lesions). Separate classifier estimate **rejection grade**
- ▸ **Multiple-instance learning:** use **patient diagnosis** as only **label**
  - − (avoid pixel-level annotations, supports large-scale deployment)
- ▸ **Attention scores,** reflecting relevance of each biopsy region, enable **visual interpretation** of the model's predictions



*Lipkova et al. Nature Medicine (2022)*

# Digital Pathology: Whole Slide Images (WSIs)

- High resolution scan of an entire tissue section (0.25 - 0.5 microns per pixel)
- 1 WSI ~ 1 billion pixels !!!
- **100 WSI has more pixels than <u>whole</u> ImageNet**
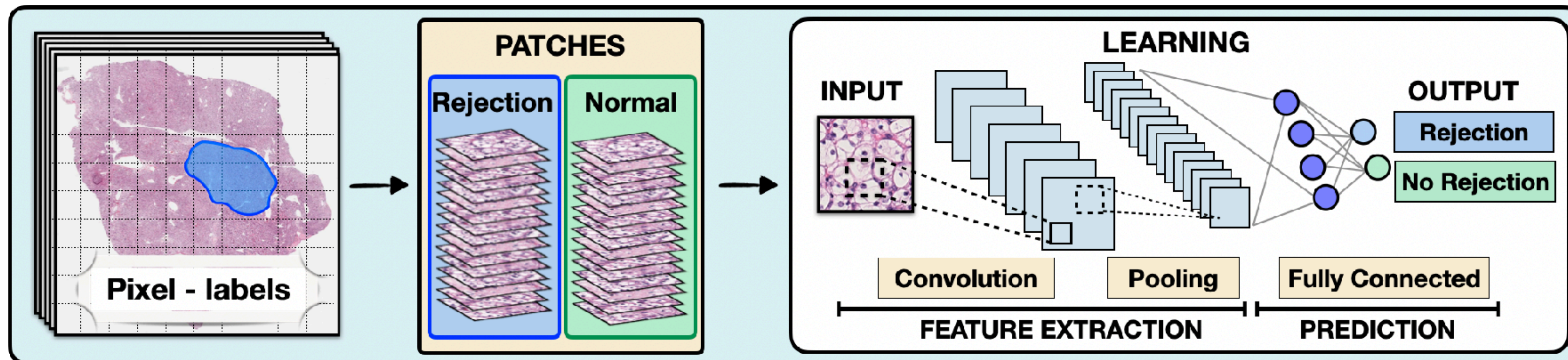- Difficult to train AI directly on WSI

# Typical Deep Learning for Pathology



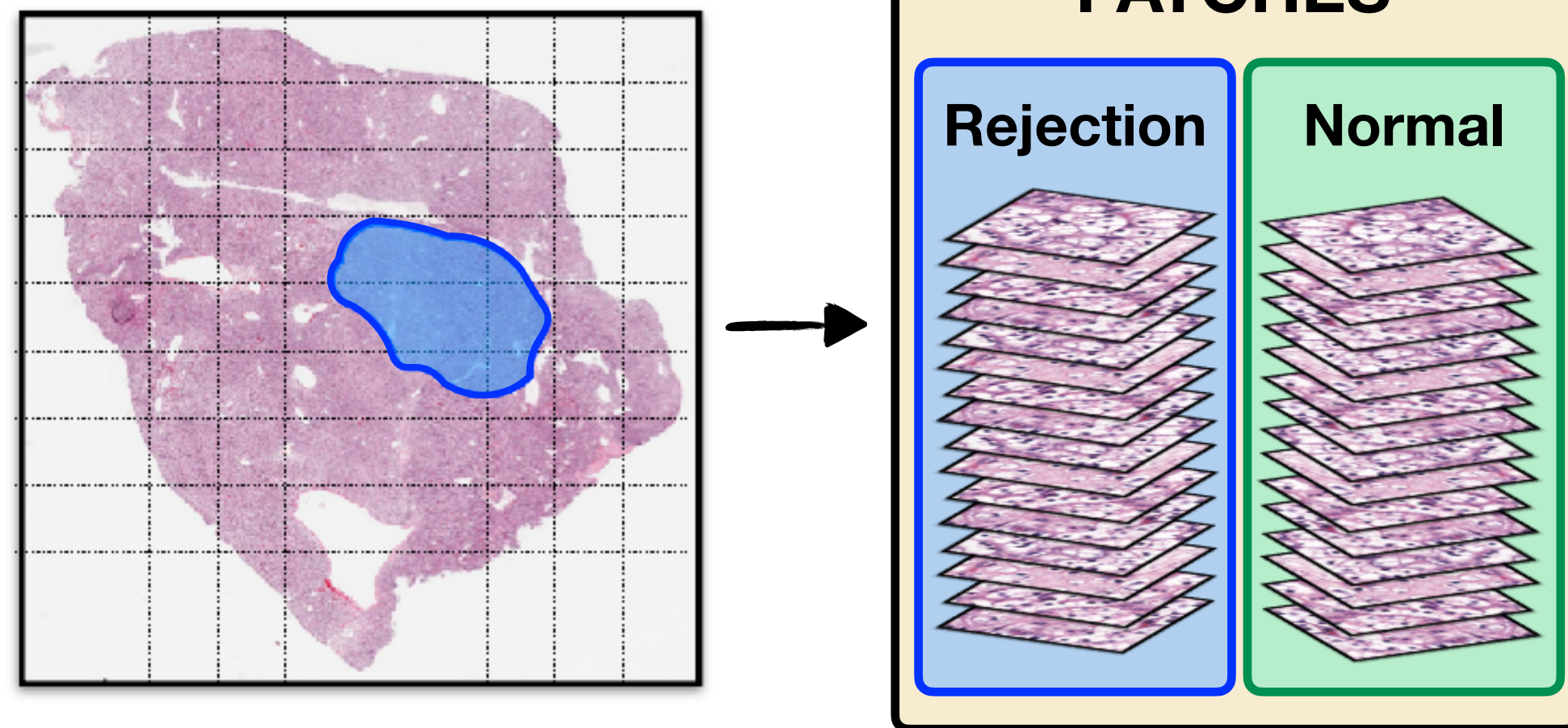- ▸ **Laborious** and **time consuming** to annotate gigapixels large histology images
- ▸ Disease borders not always well defined → **inter-rater variability** → **bias**
- ▸ **Predictive regions** for some tasks (e.g. treatment response) **might be unknown**
- ▸ **Possible data imbalance:** small proportion of image contain the disease (needle-in-haystack problem)
- ▸ **Image annotation is not part of standard clinical practice**

# Strong vs Weak Supervision



**STRONG LABELS**

**PATCH-LEVEL LABELS**

PATCHES

Rejection    Normal

**WEAK LABELS**

**PATIENT-LEVEL LABELS**

Normal

Rejection

PATCHES

➡ Model alone must discover which tissue regions and which features are predictive for rejections.

# Analogy with Natural Images

**STRONG LABELS**

▸ Label for each input

| Muffin | Chihuahua |



**WEAK LABELS**
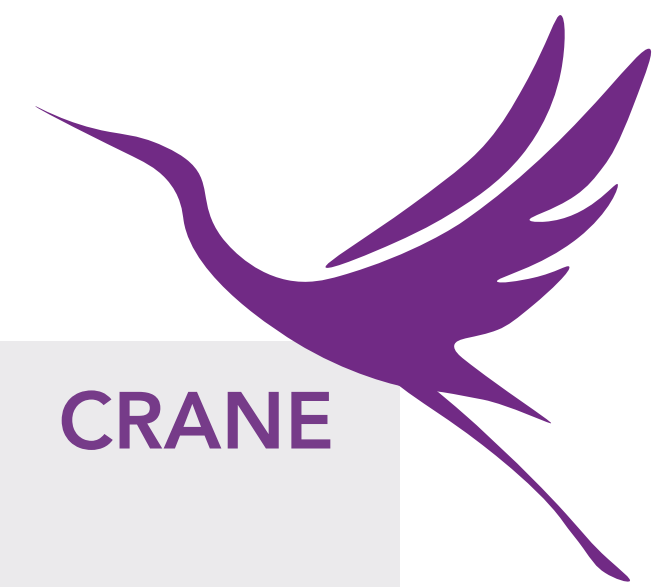
▸ Label for bag of inputs

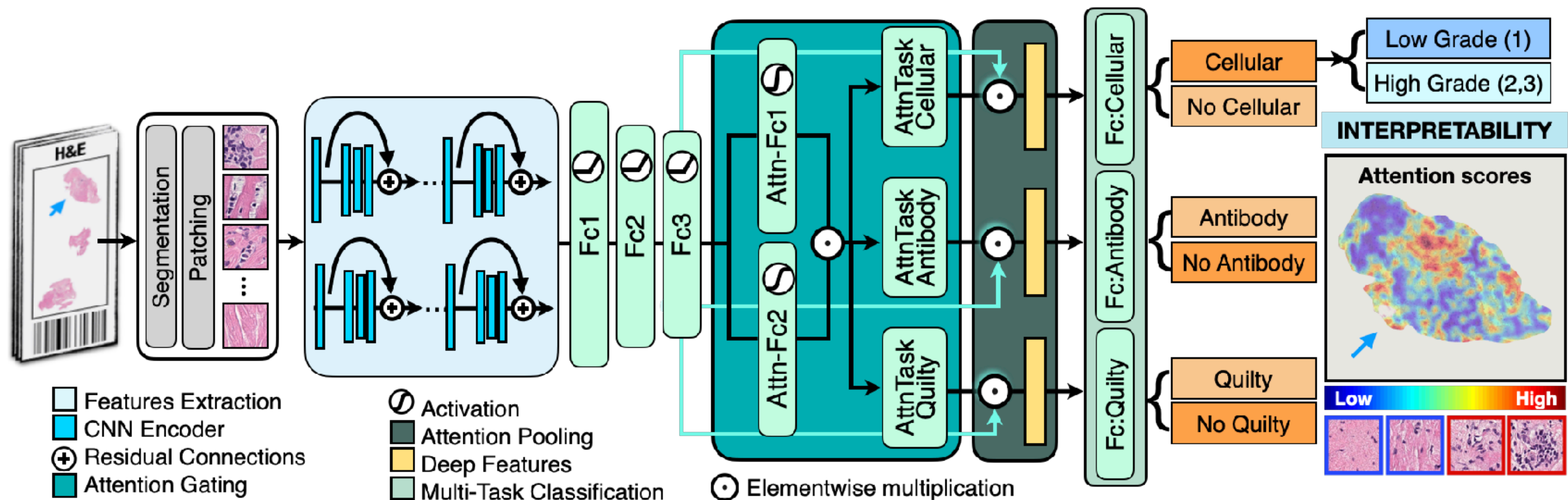| Contains Chihuahua | No Chihuahua |



➡ The model alone has to discover which image items and features correspond to chihuahua
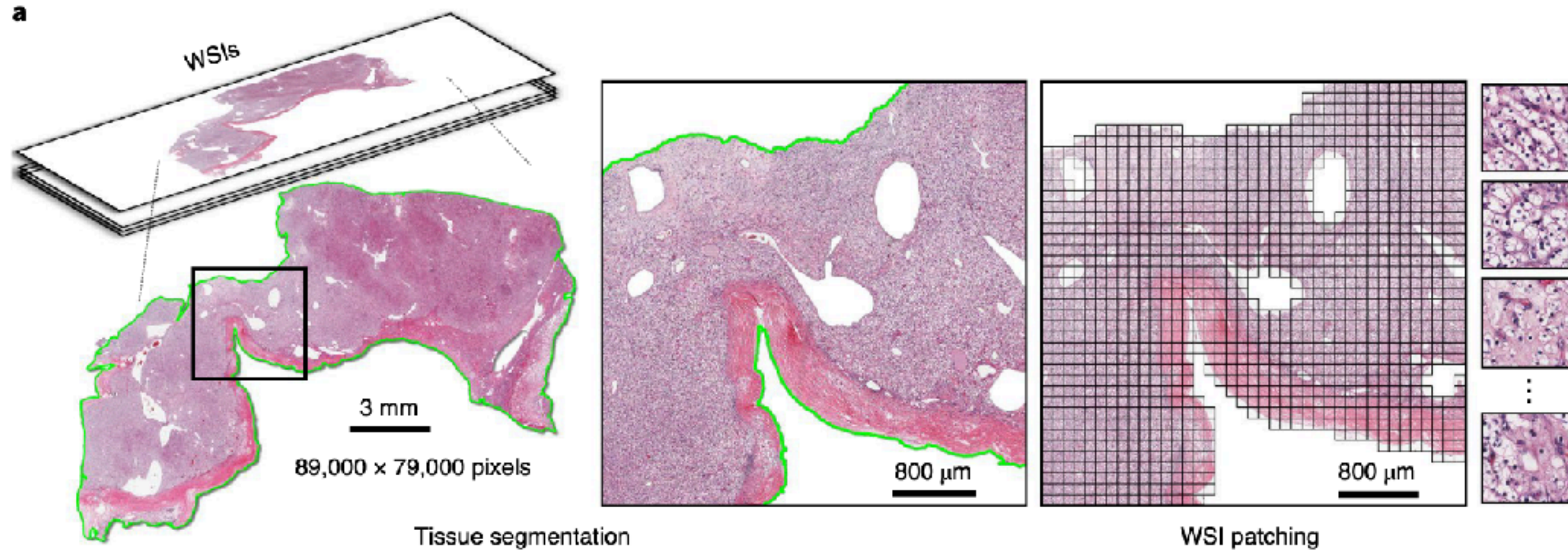
# Cardiac Rejection Assessment Neural Estimator

- ▸ **Input:** H&E-stained EMBs whole-slide-images (WSIs)
- ▸ **Multi-task, multi-label model:** simultaneously identifies **presence** and **type of the rejection** (cellular, antibody, and/or quilty lesions). Separate classifier estimate **rejection grade**
- ▸ **Multiple-instance learning:** use **patient diagnosis** as only **label**
  - – (avoid pixel-level annotations, supports large-scale deployment)
- ▸ **Attention scores,** reflecting relevance of each biopsy region, enable **visual interpretation** of the model's predictions

CRANE



*Lipkova et al. Nature Medicine (2022)*

a



Tissue segmentation

WSI patching

**≈ 1 Billion Pixels!**

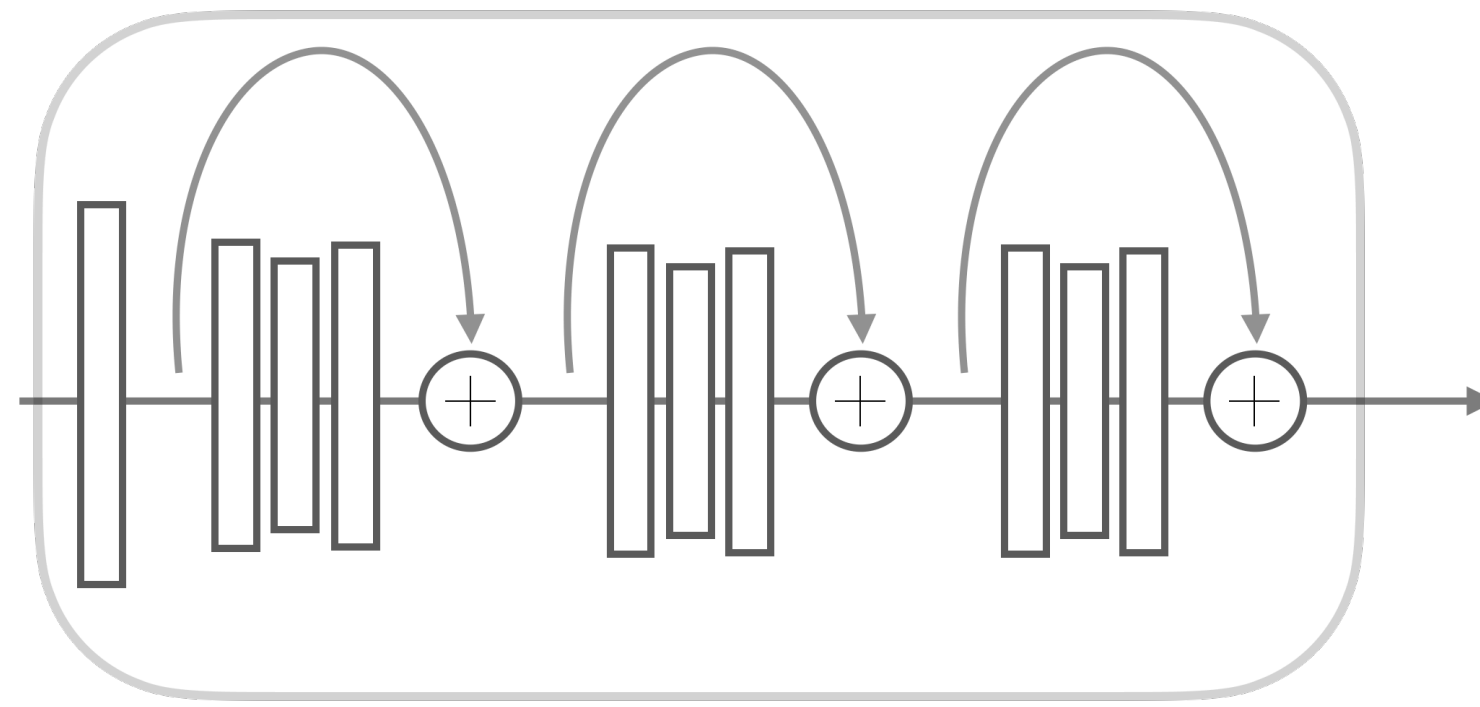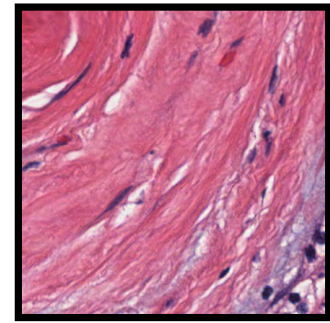▸ Patch-level representation of patch k from {1,...,K}

$$X = \{\mathbf{x}_1, \mathbf{x}_2, \ldots, \mathbf{x}_K\}$$

Input $\mathbf{x}_k$ :
$256 \times 256 \times 3$



Pretrained Encoder

$$f(\cdot; \theta) : \mathbb{R}^{C \times H \times W} \rightarrow \mathbb{R}^{1024}$$

Embedding $z_k$:
1024

▸ ResNET50 features $\mathbf{z}_k \in \mathbb{R}^{1024}$
▸ Three FC layers:
 - FC1: $\mathbf{W}_1 \in \mathbb{R}^{768 \times 1024}$
 - FC2: $\mathbf{W}_2 \in \mathbb{R}^{512 \times 768}$
 - FC3: $\mathbf{W}_3 \in \mathbb{R}^{512 \times 512}$



**INTERPRETABILITY**

**Attention scores**

Low    High

□ Features Extraction
■ CNN Encoder
⊕ Residual Connections
■ Attention Gating

⊘ Activation
■ Attention Pooling
□ Deep Features
■ Multi-Task Classification

⊙ Elementwise multiplication

# ATTENTION LEARNING

**Attention score** *(for patch k and task t):*

$$a_{k,t} = \frac{\exp\left\{ \mathbf{W}_{a,t} \left( \tanh\left( \mathbf{V}_a \mathbf{h}_k \right) \odot \mathrm{sigm}\left( \mathbf{U}_a \mathbf{h}_k \right) \right) \right\}}{\sum_{j=1}^{N} 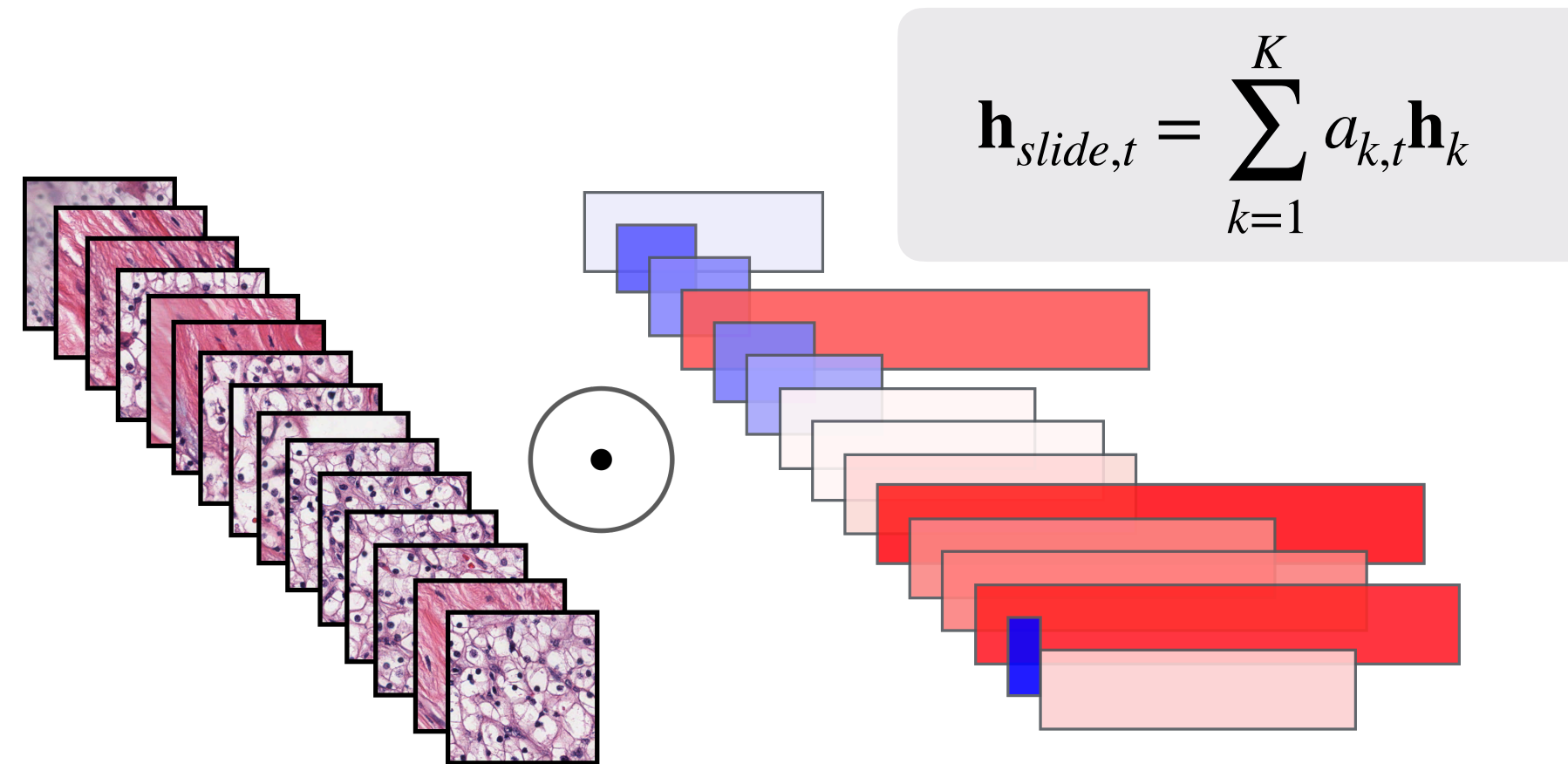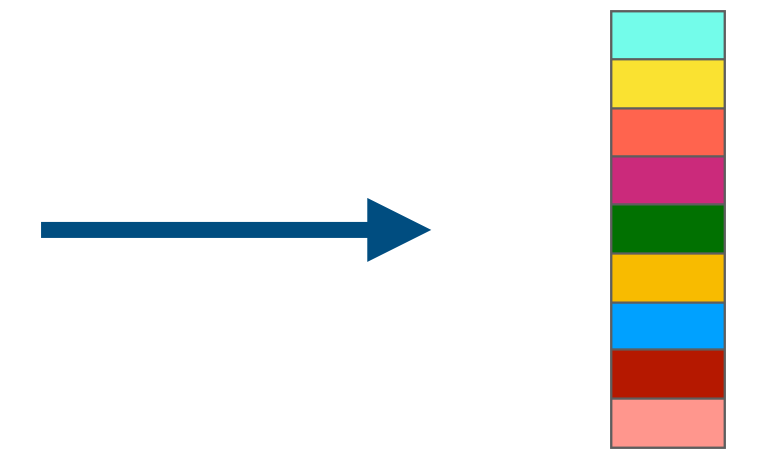\exp\left\{ \mathbf{W}_{a,t} \left( \tanh\left( \mathbf{V}_a \mathbf{h}_j \right) \odot \mathrm{sigm}\left( \mathbf{U}_a \mathbf{h}_j \right) \right) \right\}}$$
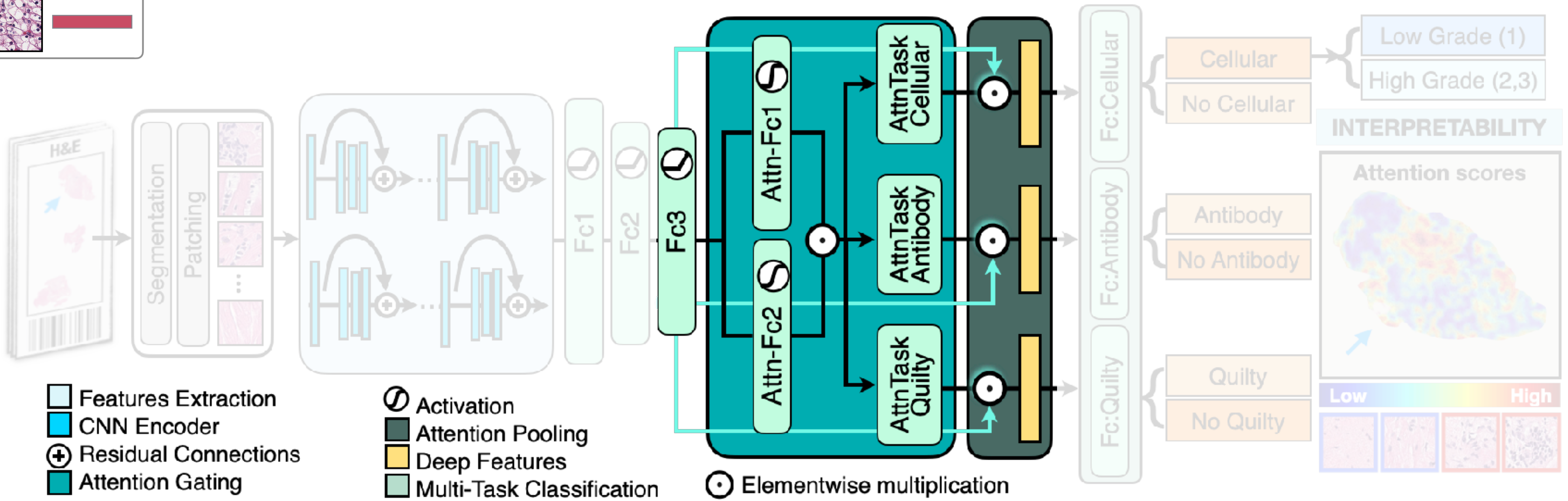
**Attention-based pooling**

$$\mathbf{h}_{slide,t} = \sum_{k=1}^{K} a_{k,t} \mathbf{h}_k$$

**Learned WSI Embedding**

*(Ilse et al. ICML 2018)*

Not Informative

Highly Informative

Attn-Fc1

Attn-Fc2

AttnTask Cellular

AttnTask Antibody

AttnTask Quilty

Fc1

Fc2

Fc3

H&E

Segmentation

Patching

Fc:Cellular

Fc:Antibody

Fc:Quilty

Cellular

No Cellular

Antibody

No Antibody

Quilty

No Quilty

Low Grade (1)

High Grade (2,3)

**INTERPRETABILITY**

**Attention scores**

Low    High

Features Extraction
CNN Encoder
Residual Connections
Attention Gating

Activation
Attention Pooling
Deep Features
Multi-Task Classification

⊙ Elementwise multiplication

## MULTI-TASK CLASSIFIER

**REJECTION GRADE**

- Same MIL model, just single-task

**Learned WSI Embedding**

**INTERPRETABILITY**

- WSI attention heatmaps

Fc:Cellular — Cellular / No Cellular

Fc:Antibody — Antibody / No Antibody
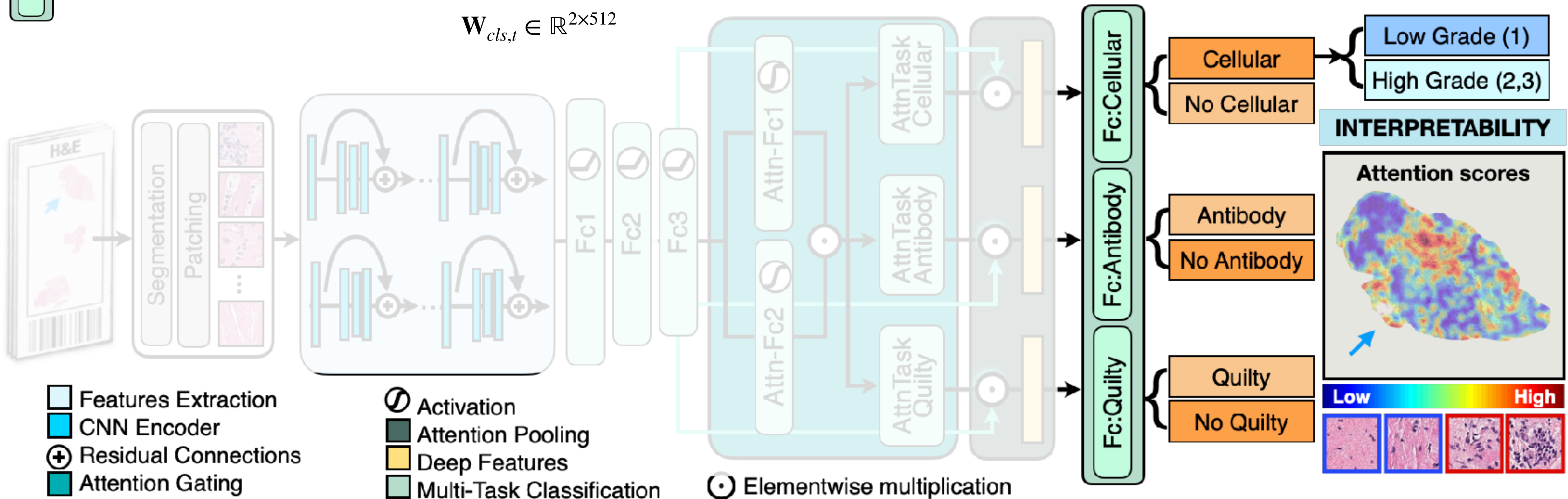
Fc:Quilty — Quilty / No Quilty

- Slide-level representations for task $t$:

$$\mathbf{h}_{slide,t} = \sum_{k=1}^{K} a_{k,t} \mathbf{h}_k$$

- Slide-level predictions for task t:

$$\mathbf{p}_t = \text{Softmax}(\mathbf{W}_{cls,t} \mathbf{h}_{slide,t} + \mathbf{b}_{cls,t})$$

$$\mathbf{W}_{cls,t} \in \mathbb{R}^{2\times512}$$

Fc:Cellular — Cellular / No Cellular → Low Grade (1) / High Grade (2,3)

**INTERPRETABILITY**

**Attention scores**

Low — High

Fc:Antibody — Antibody / No Antibody

Fc:Quilty — Quilty / No Quilty

H&E · Segmentation · Patching

Fc1 · Fc2 · Fc3

Attn-Fc1 · Attn-Fc2

AttnTask Cellular · AttnTask Antibody · AttnTask Quilty

- Features Extraction
- CNN Encoder
- ⊕ Residual Connections
- Attention Gating
- ⊘ Activation
- Attention Pooling
- Deep Features
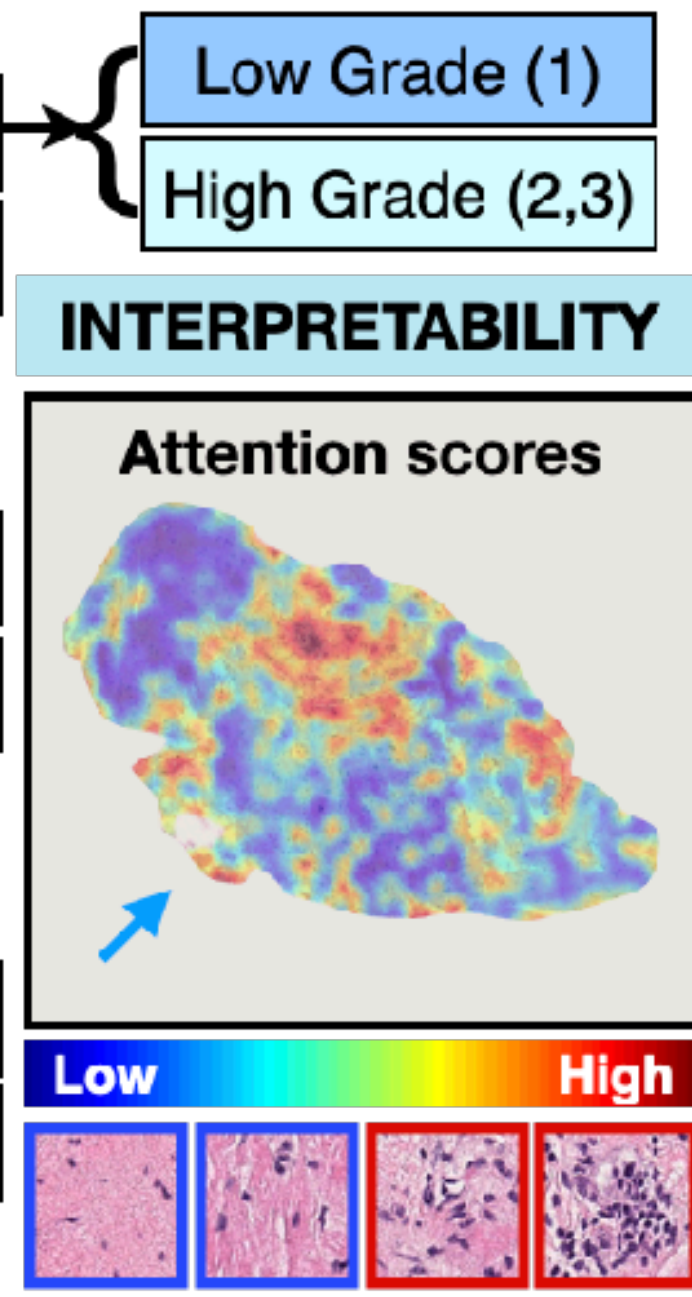- Multi-Task Classification
- ⊙ Elementwise multiplication

# Study Design

## WSI
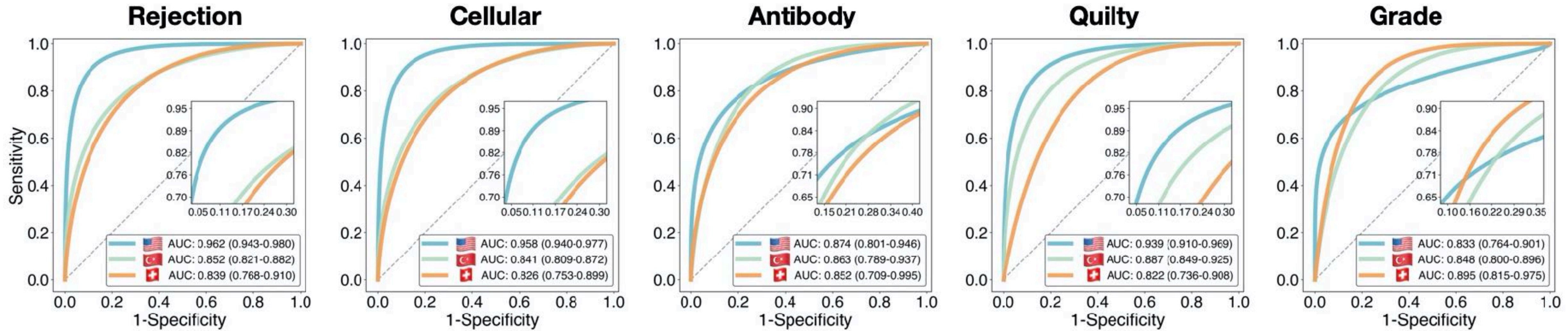
▸ Over 7000 WSI from 3 **independent centers**

▸ **Large diversity:**

- *population (i.e., pediatric vs adult),*
- *scanners,*
- *biopsy protocols,*
- *staining (manual vs automated),*
- *noise,*
- *micron/pixel,*
- *etc*

▸ The model is trained on subset of data collected in USA

- 70/10/20% split (balance diagnosis)
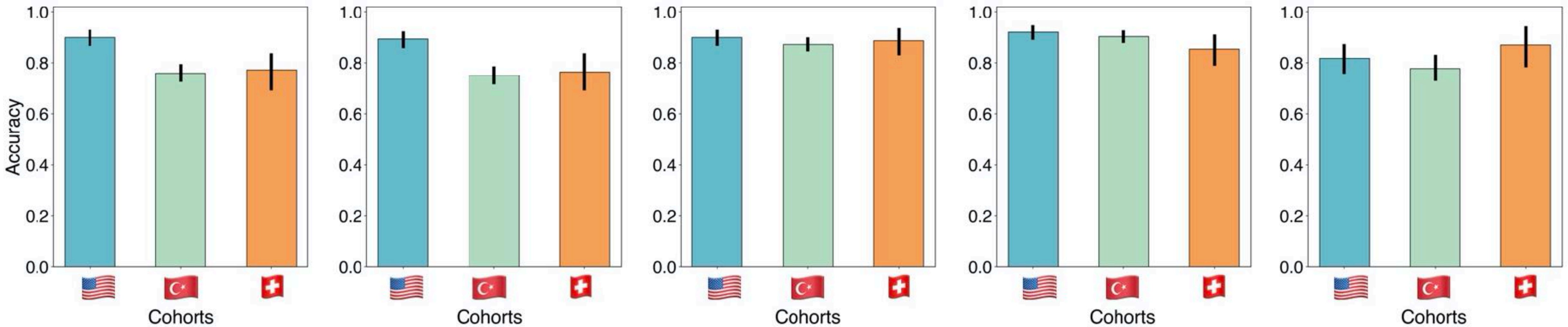- **...tion** to external cohorts
- **...domain-specific adaptations**

## DIAGNOSIS

**TRAIN**

| # 🔲 | 4059 |
|---|---|
| # 👤 | 1354 |
| Scanner | **HAMAMATSU** |
| Magnification | 40x |
| Micron/pixel | 0.2206 |

**TEST**

**2004-2021**

| # 🔲 | 995 |
|---|---|
| # 👤 | 336 |
| Scanner | **HAMAMATSU** |
| Magnification | 40x |
| Micron/pixel | 0.2206 |

**2002-2020**

| # 🔲 | 1717 |
|---|---|
| # 👤 | 585 |
| Scanner | *Leica* **APERIO** |
| Magnification | 40x |
| Micron/pixel | 0.2523 |

**2014-2020**

| # 🔲 | 123 |
|---|---|
| # 👤 | 123 |
| Scanner | **3DHISTECH** |
| Magnification | 20x |
| Micron/pixel | 0.2431 |

## PREDICTIONS

*Lipkova et al. Nature Medicine (2022)*

Low Grade (1)

Level 1

Level 2

Level 3

# Evaluation & Results



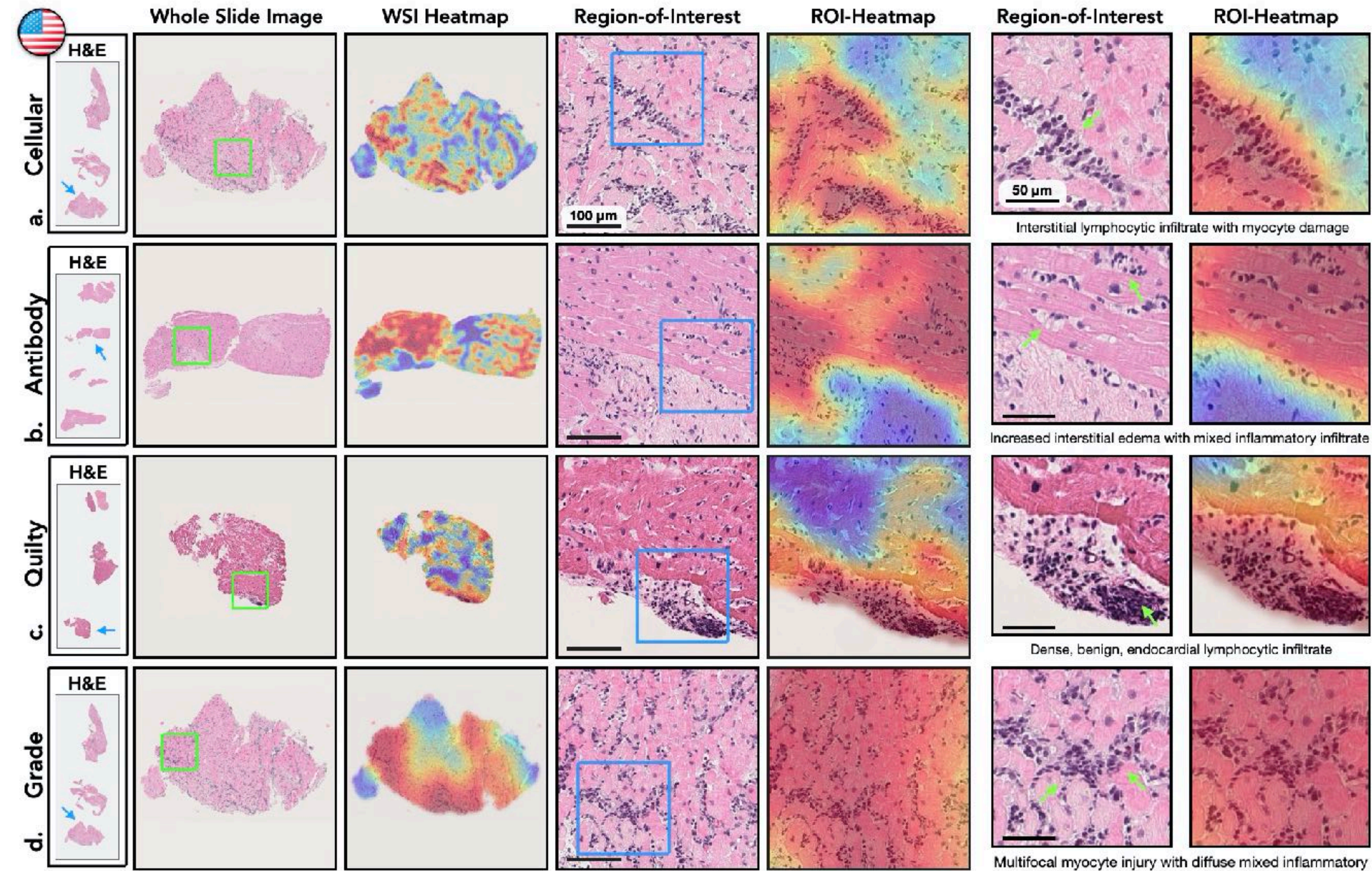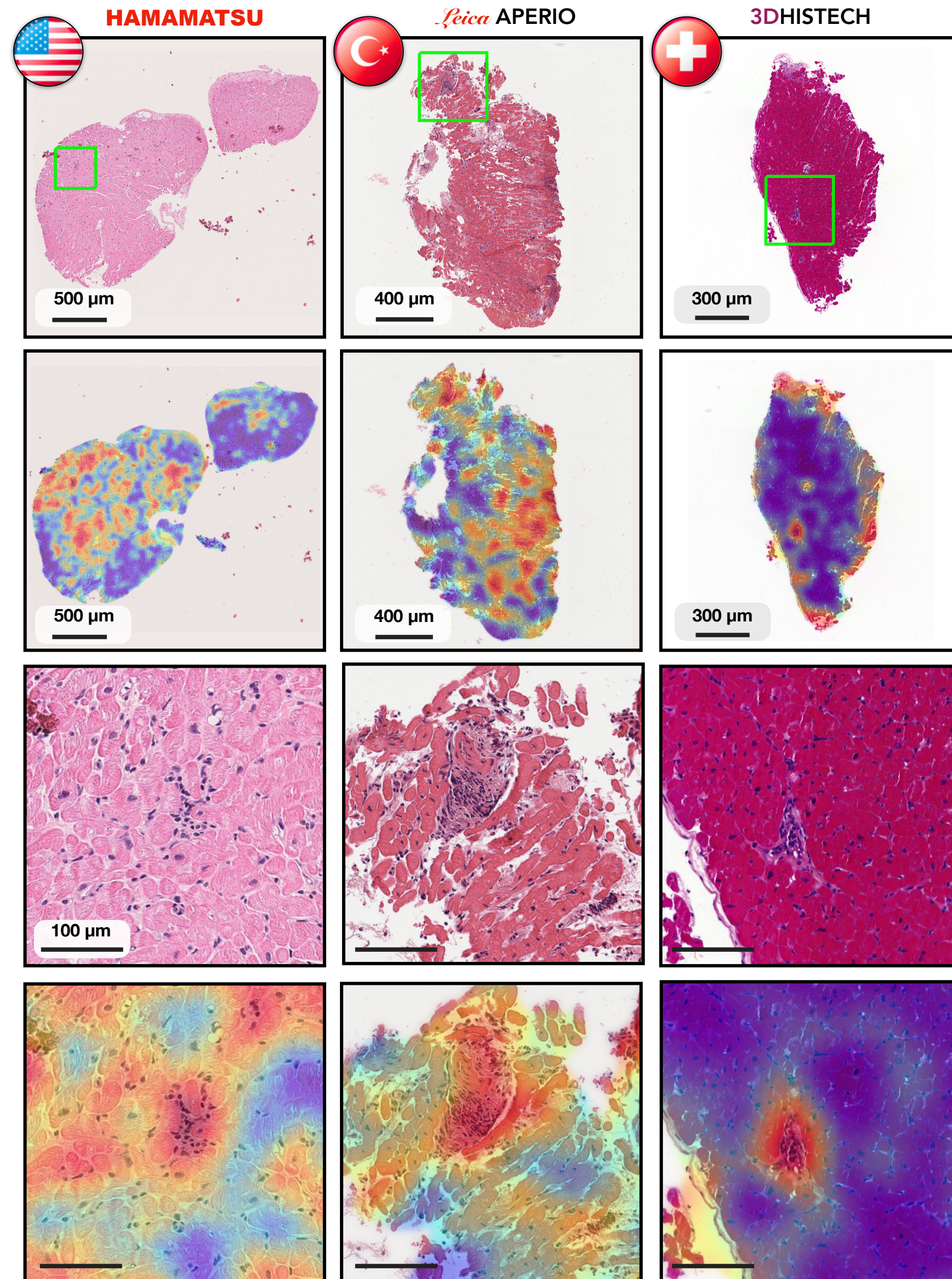*Lipkova et al. Nature Medicine (2022)*
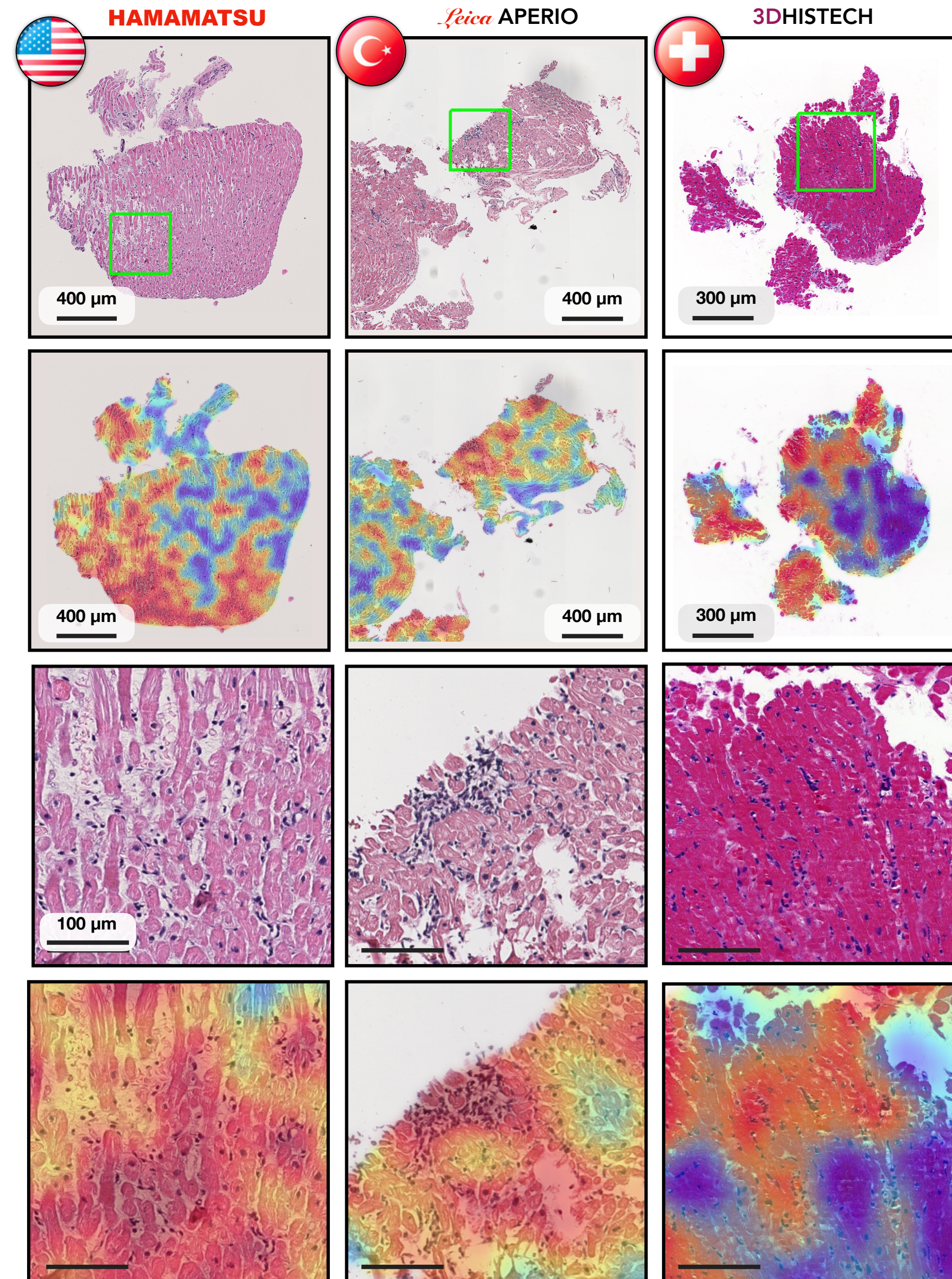
# Interpretability



- ▸ **High-attention (red)** regions **correspond to rejection morphology used by** pathologist for diagnosis
- ▸ **Low-attention (blue)** scores are assigned mostly to **benign tissue**

*Lipkova et al. Nature Medicine (2022)*

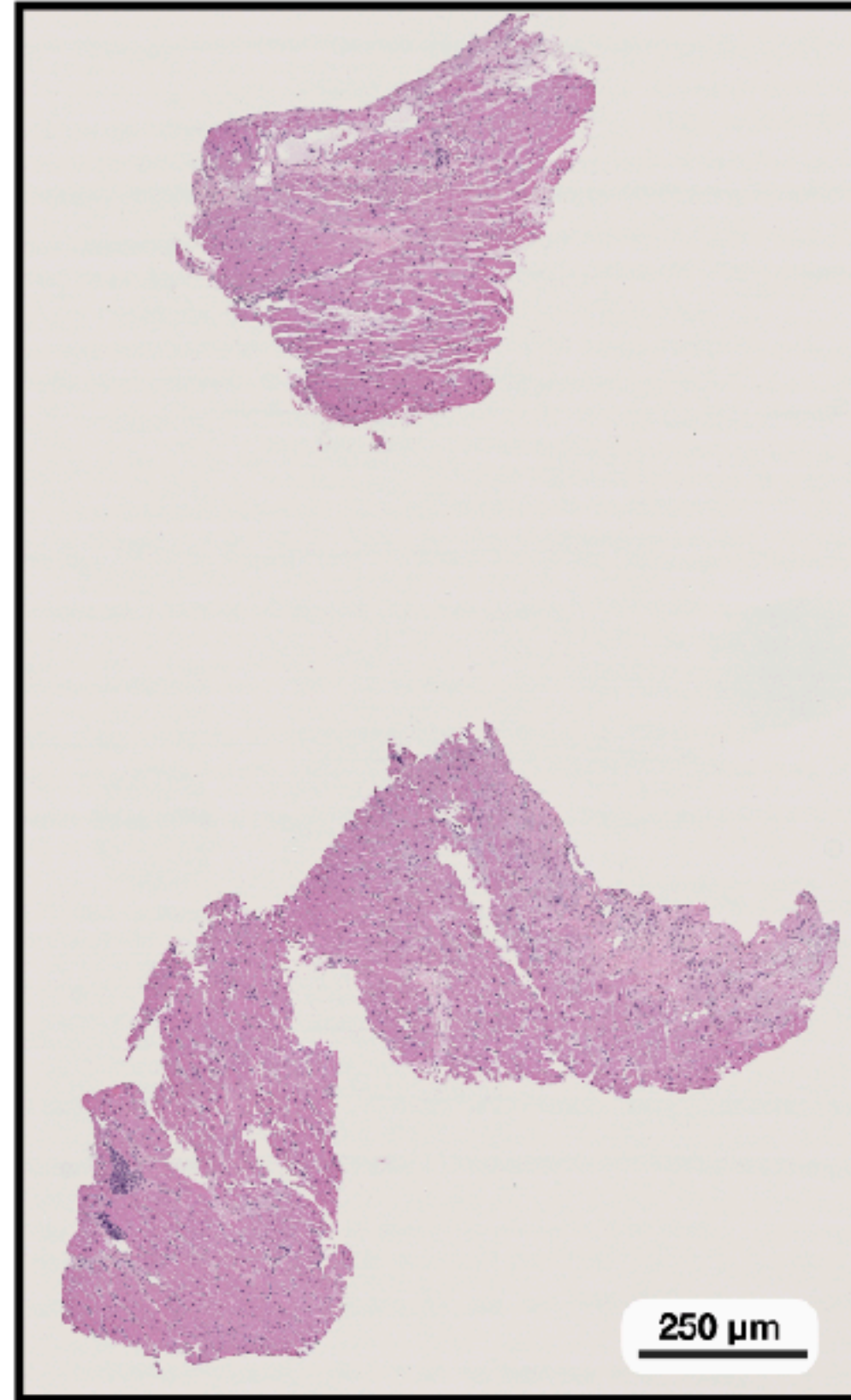# Assessment of Failure Cases



**a.** Model: Normal    True: Cellular

**b.** Model: Normal    True: Antibody
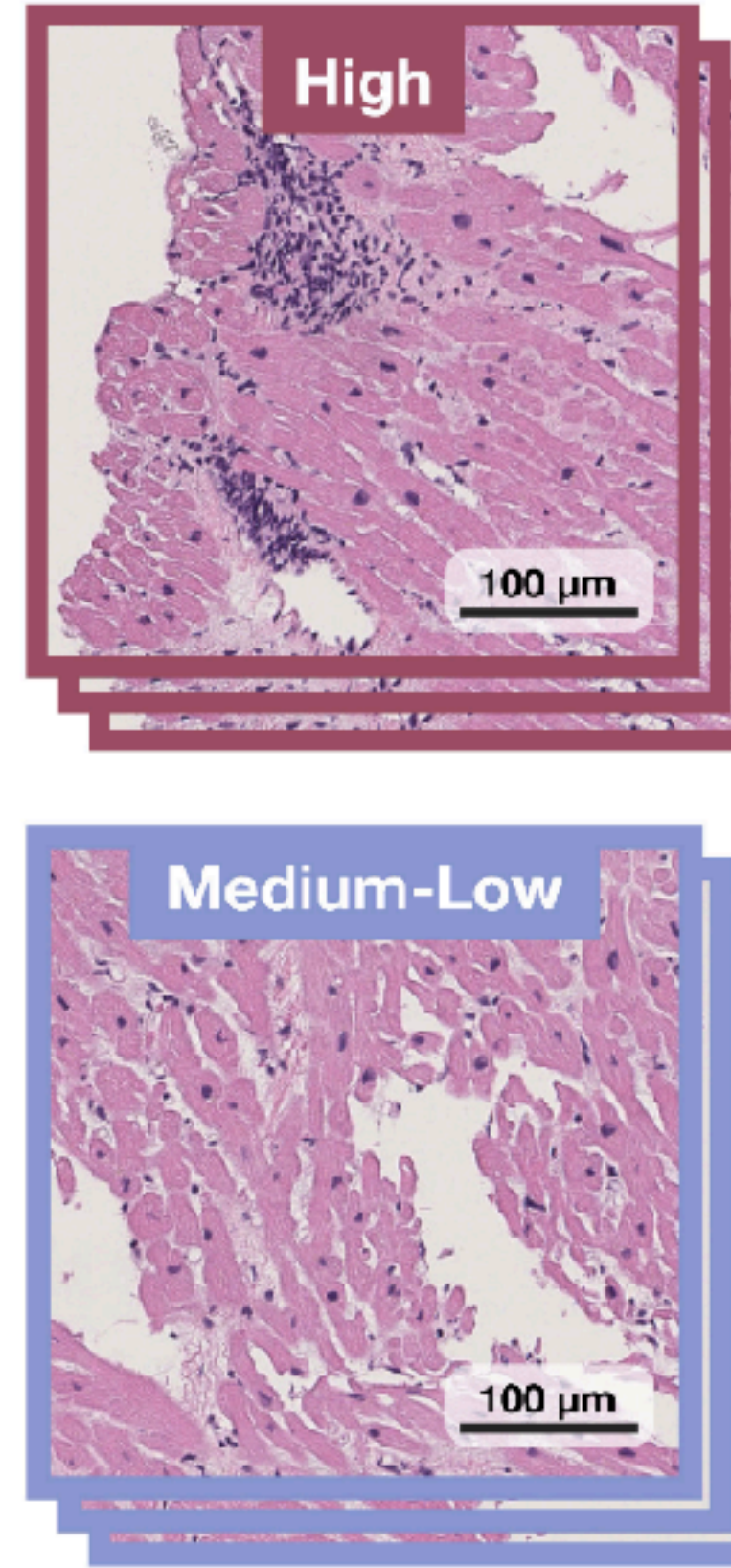
# Quantitative Assessment of Interpretability



a. Whole Slide Image

b. WSI Heatmap

c. Patches

High

100 μm

Medium-Low

100 μm

250 μm

d. Diagnostic Relevance

AUC ROC: 0.880 (0.850-0.910)

e. Patch-Level Scores

| Tasks: | Accuracy | F1 | $\kappa$ |
|--------|----------|------|------|
| All | 0.873 | 0.855 | 0.744 |
| Cellular | 0.925 | 0.914 | 0.848 |
| Antibody | 0.902 | 0.911 | 0.802 |
| Quilty | 0.809 | 0.729 | 0.596 |

f. Pathologist annotation

200 μm

g. High-attention regions

h. Slide-Level Scores

| Tasks: | Detection rate |
|--------|----------------|
| All | 0.922 |
| Cellular | 0.942 |
| Antibody | 0.901 |
| Quilty | 0.924 |

INTERACTIVE DEMO

crane.mahmoodlab.org

# Comparison with human readers



**a**

AI Model Trained on US Cases → 150 Cases from Turkish Cohort used for observer study

**b**

Level 1 / Level 2 / Level 3

AI Read

P1 P2 P3

P4 P5

**150 Cases:** 91 ACR; 23 AMR (14 ACR+AMR); 50 Normal

**(avg. 10.5 years of experience)**

▸ **Cohens' κ** (-1 to 1): inter-observer agreement:
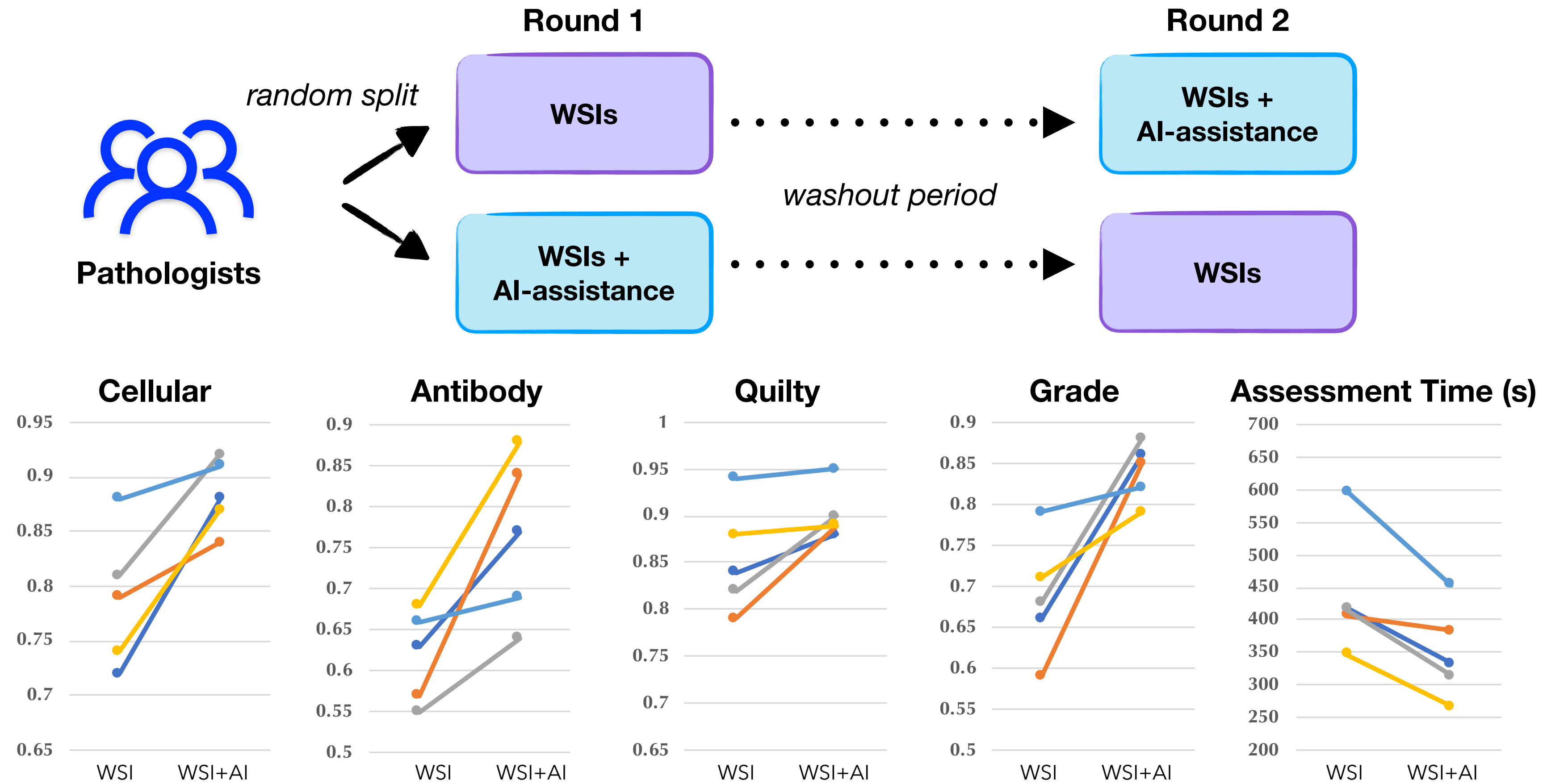▸ Agreement between expert is comparable to previous studies

▸For all tasks **AI-predictions are <u>not inferior</u> to human experts:**
  ▸ avg. agreement on **rejection** between **pathologists** **κ** = 0.537 (moderate agreement)
  ▸ avg. agreement between **pathologists and model** **κ** = 0.639 (substantial agreement)

**c**

Legend: Rejection, Cellular, Antibody, Quilty, Grade

P1 - P2, P1 - P3, P1 - P4, P1 - P5, P2 - P3, P2 - P4, P2 - P5, P3 - P4, P3 - P5, P4 - P5

Cohen's κ (0.0 – 1.0)

**d**

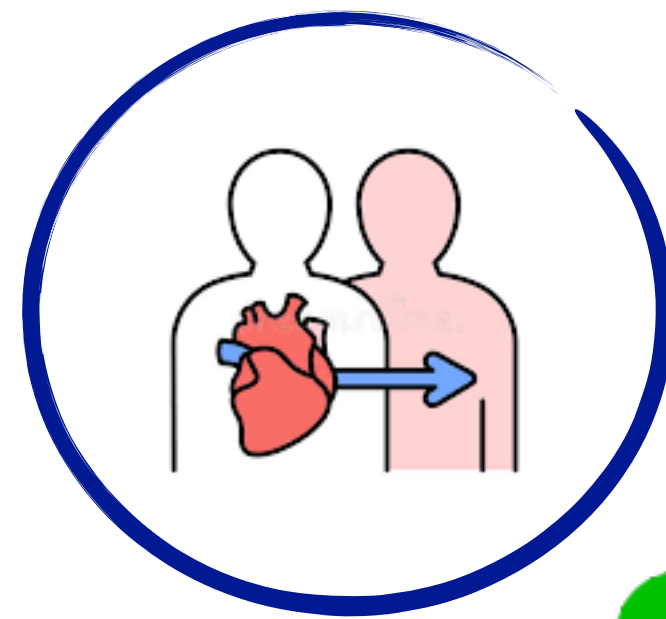P1 - AI, P2 - AI, P3 - AI, P4 - AI, P5 - AI

Cohen's κ (0.0 – 1.0)

# Clinical Potential



- ▶ **Ground-truth labels**: consensus of readers from the first study
- ▶ **AI-assistance**: attention heatmaps as semi-transparent layer at the top of H&E slide

- ▶ **For all readers:**
  - – **Increase accuracy**
  - – (i.e. reduce inter-rater variability)
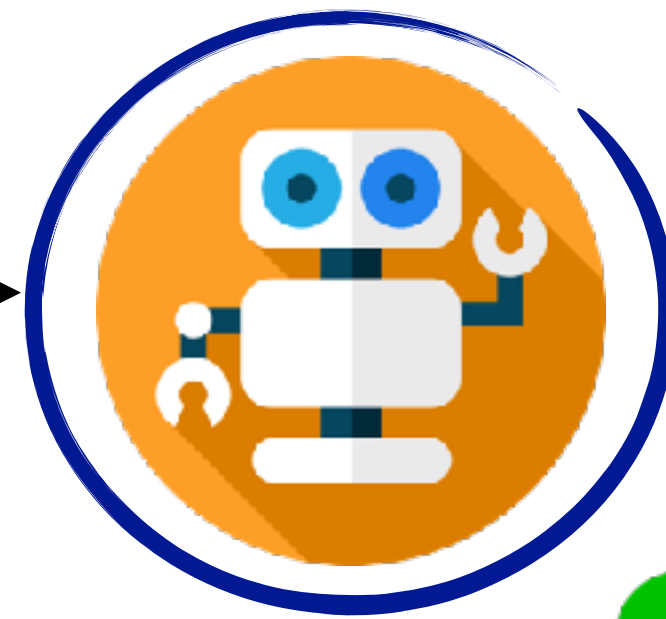  - – **Decrease assessment time**

*Lipkova et al. Nature Medicine (2022)*

# Study Design Flow Chart



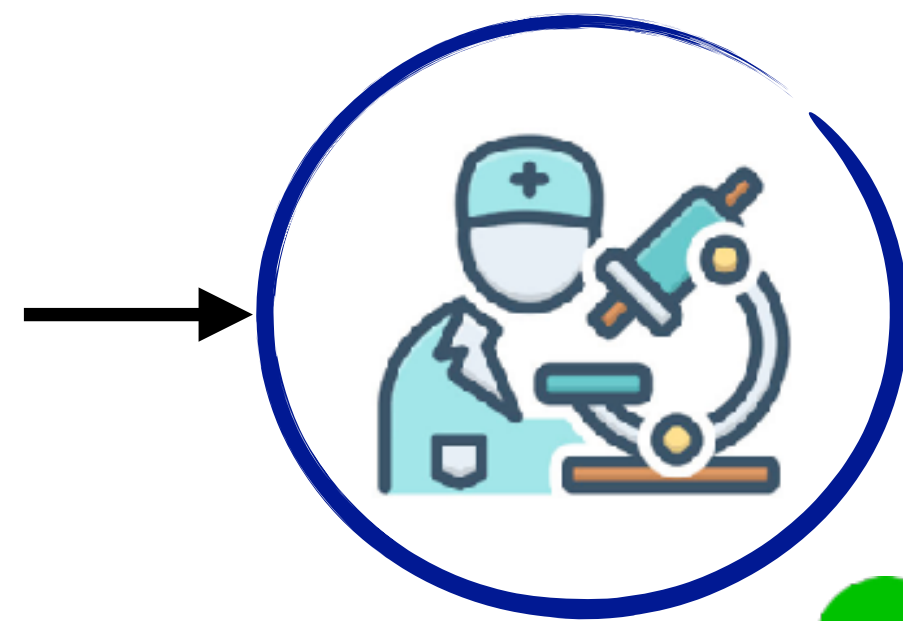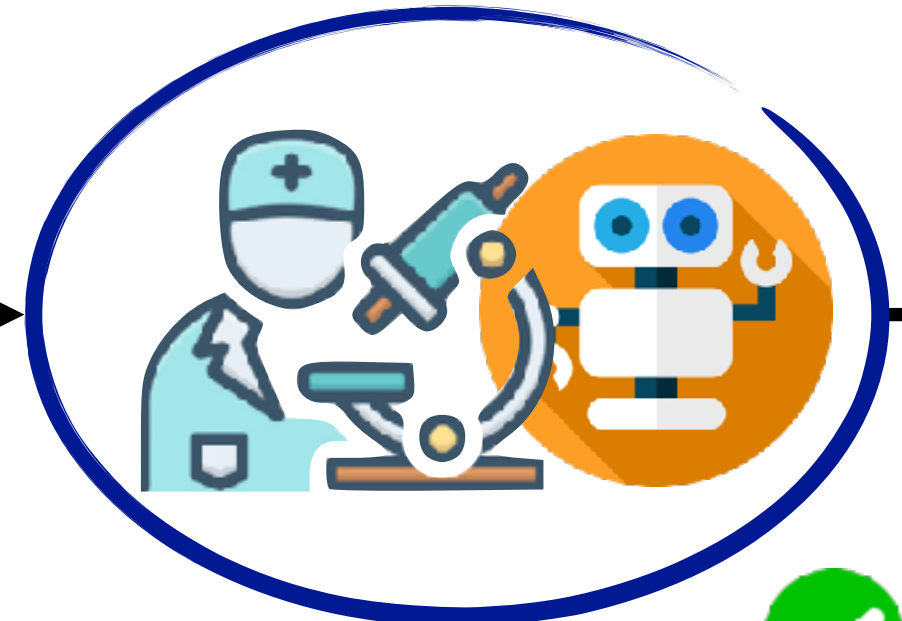Problem definition ✓ → Data collection ✓ → Label preparation ✓ → Model development ✓ → Interpretation ✓ → Robust evaluation ✓

→ Comparison with human ✓ → Clinical potential for humans ✓ → Peer-review publication ✓ → Large scale clinical trial → Clinical deployment
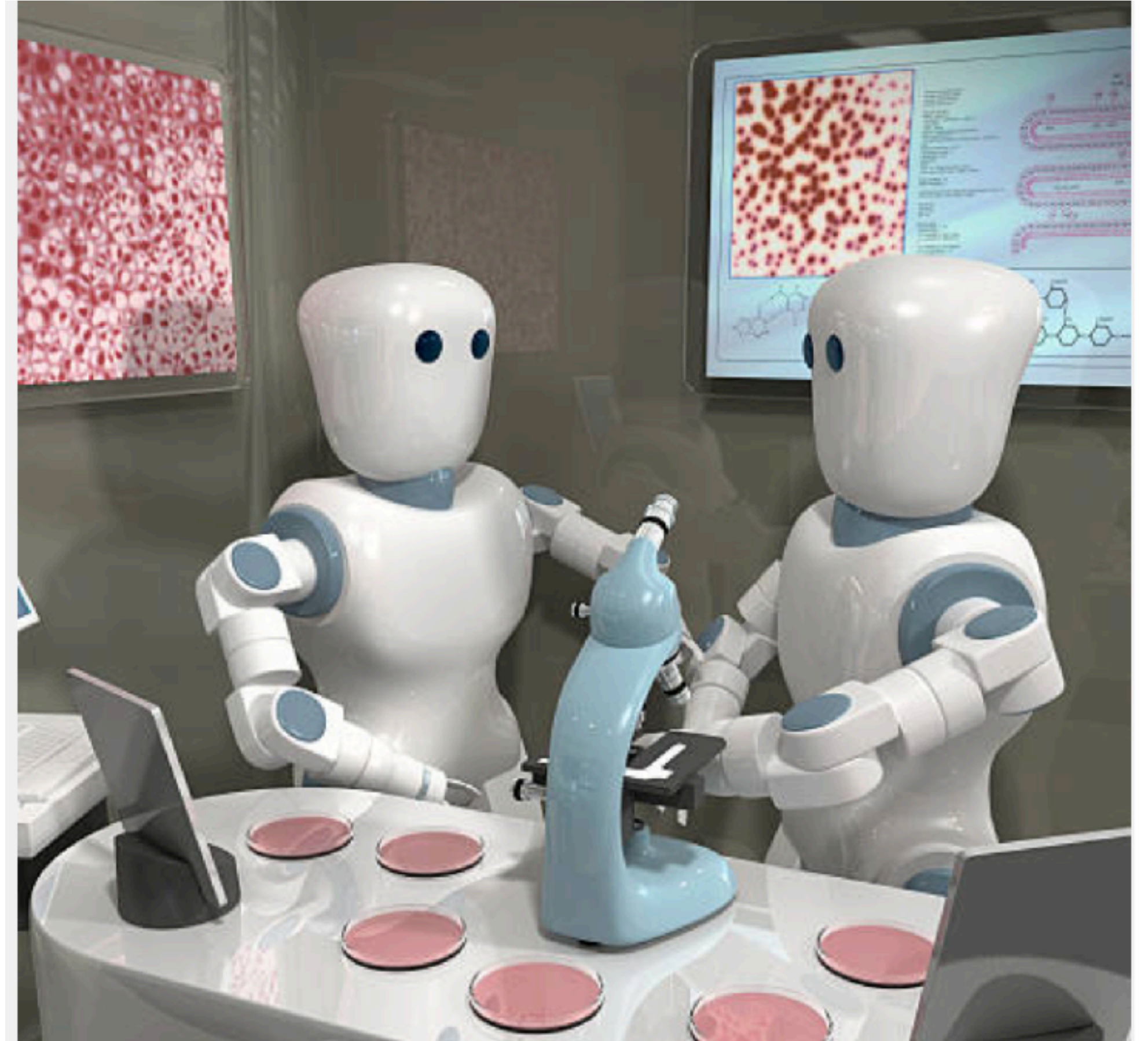
# The Mahmood Lab