

AIM 2: Artificial Intelligence in Medicine II

Harvard - BMI 702 and BMIF 203, Spring 2026
Lecture 6: Medical Imaging II

Vision foundation models. Development and validation of medical imaging interpretation models, Model robustness and performance across diverse populations.



HARVARD
MEDICAL SCHOOL

Lecture Outline

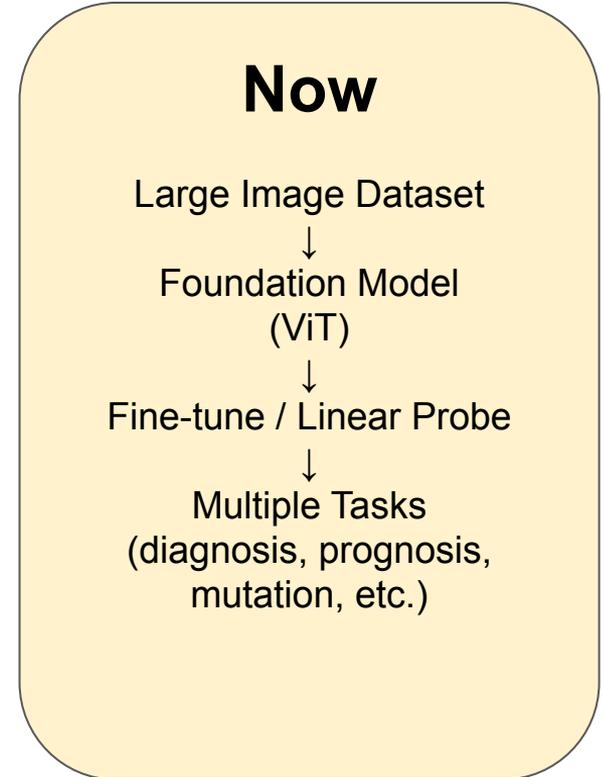
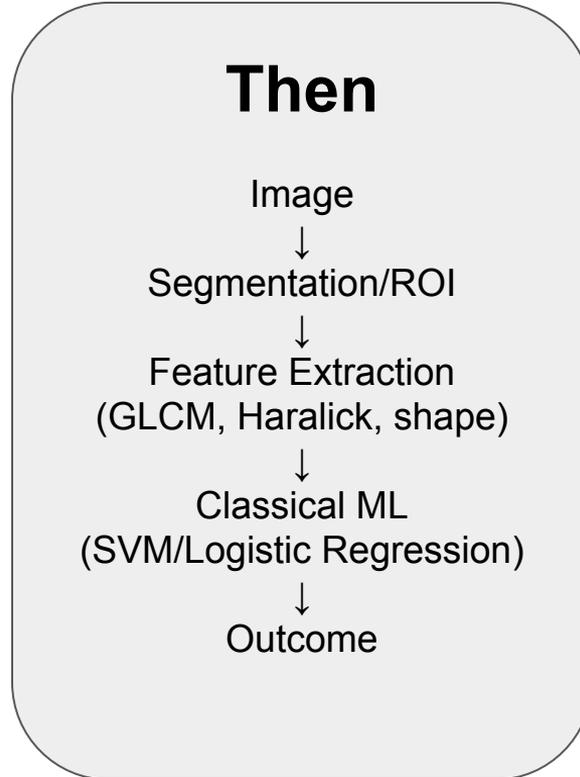
1. Foundation Models
2. Pathology Foundational Models
3. Model Interpretability
4. Model Generalizability

Lecture Outline

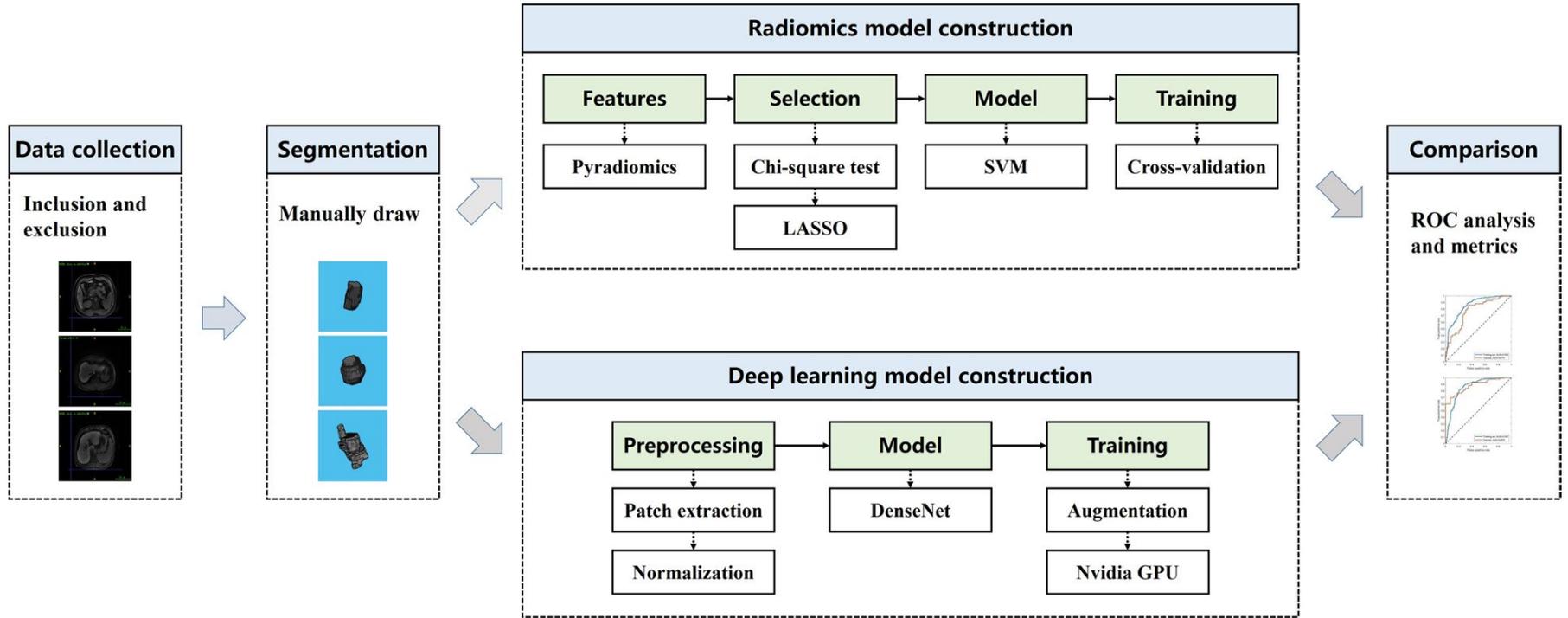
1. Foundation Models
2. Pathology Foundational Models
3. Model Interpretability
4. Model Generalizability

From Engineered Features to Foundation Models

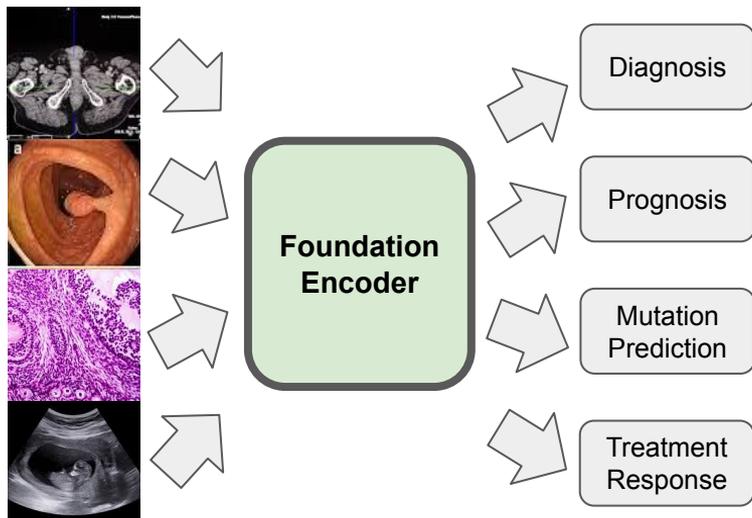
- Then:
engineered features +
small supervised
models
- Now:
large-scale pretraining
+ adaptation
- Constant:
data limits, domain
shift, clinical reality



Radiomics vs Learned Representations

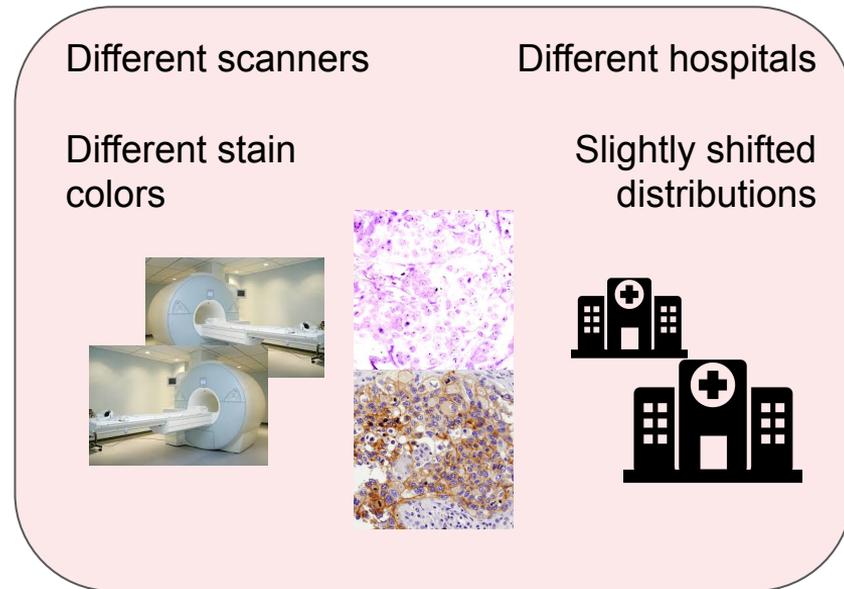


The Promise



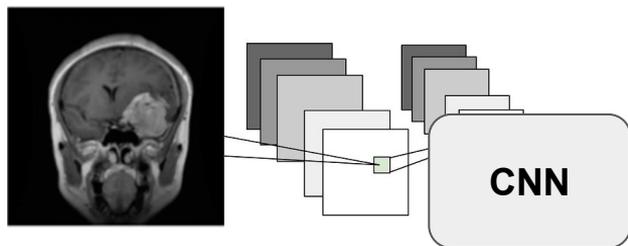
One encoder | Many tasks | Less task-specific data

The Problem



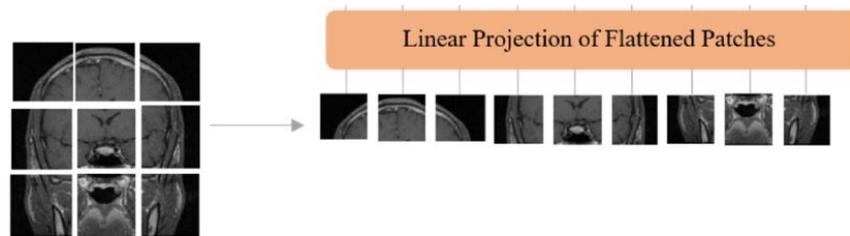
Heterogeneous modalities | Noisy labels
Site artifacts & shortcuts | Calibration & reliability

From Convolutions to Tokens



CNNs

- Local receptive fields
- Translation equivariance

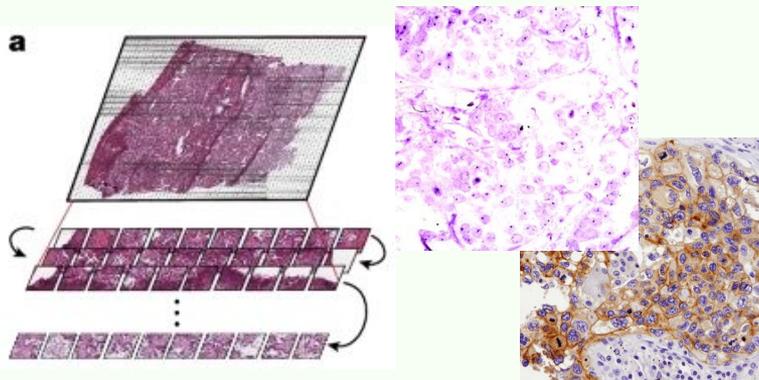


ViT

- Patchify image
- Global self-attention

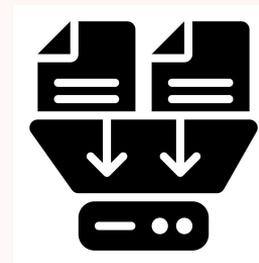
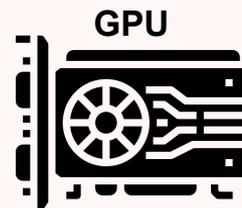
Shift: local filters → global token interactions

Patches as Tokens — Power and Cost



What it buys you

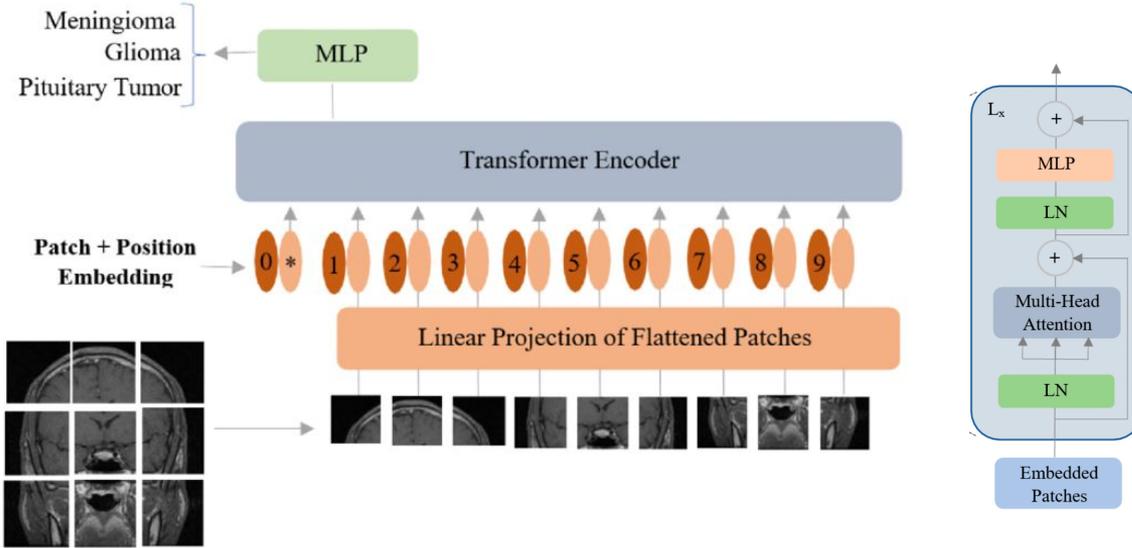
- Global context
- Flexible scale
- Modality-agnostic interface



What it costs

- Data hunger
- Compute
- Sensitivity to shift

Inside a Vision Transformer Encoder



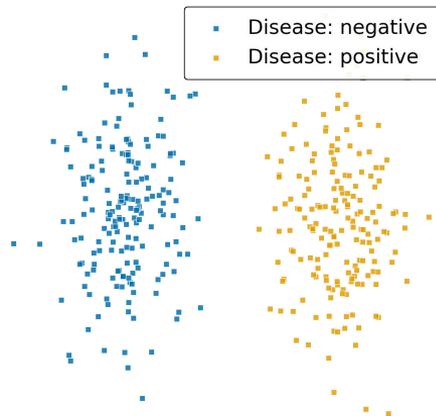
Self-Attention Mechanism

- Given Queries (\mathbf{Q}), Keys (\mathbf{K}), Values (\mathbf{V}) of dimension d :

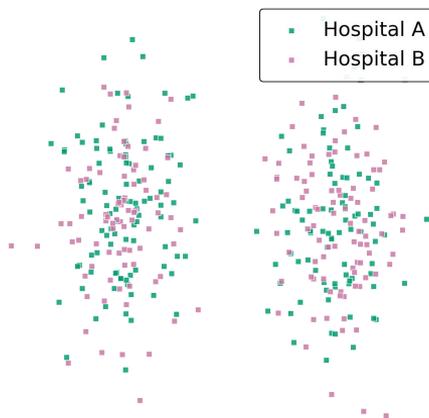
$$\text{Attention}(\mathbf{Q}, \mathbf{K}, \mathbf{V}) = \text{softmax}\left(\frac{\mathbf{Q}\mathbf{K}^T}{\sqrt{d}}\right)\mathbf{V}$$

What Makes a “Good” Representation?

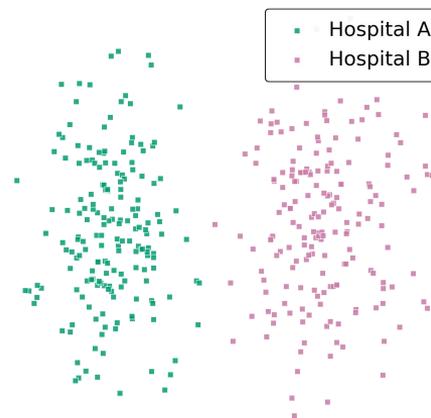
Good representation: separates disease



Good representation: hospital mixed



Bad representation: separates hospital

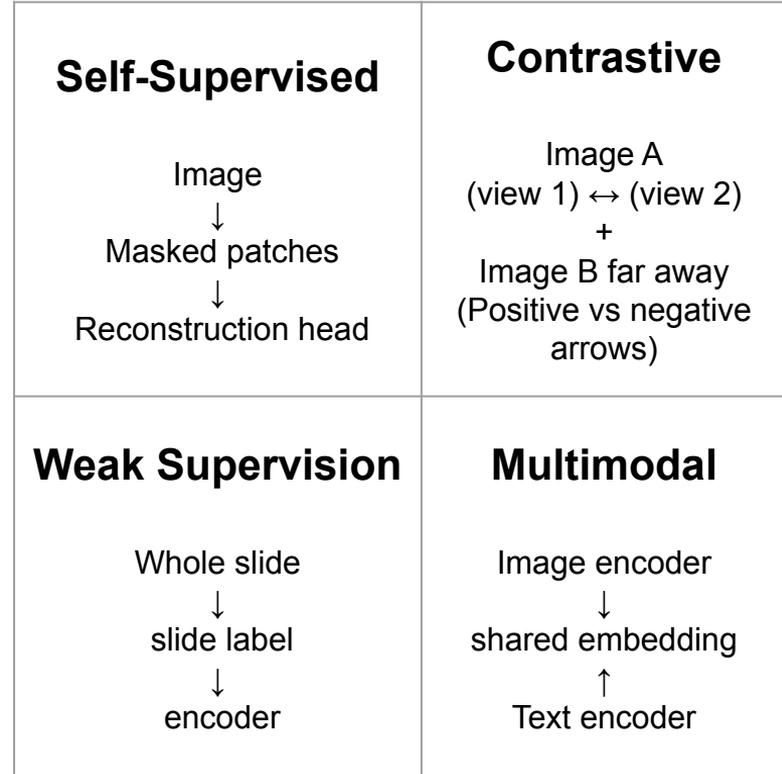


How Foundation Models are Built

Pretraining Signals

- Self-supervised
- Contrastive
- Weak supervision
- Multimodal alignment

Different signals → different representations



Self-Supervised Learning

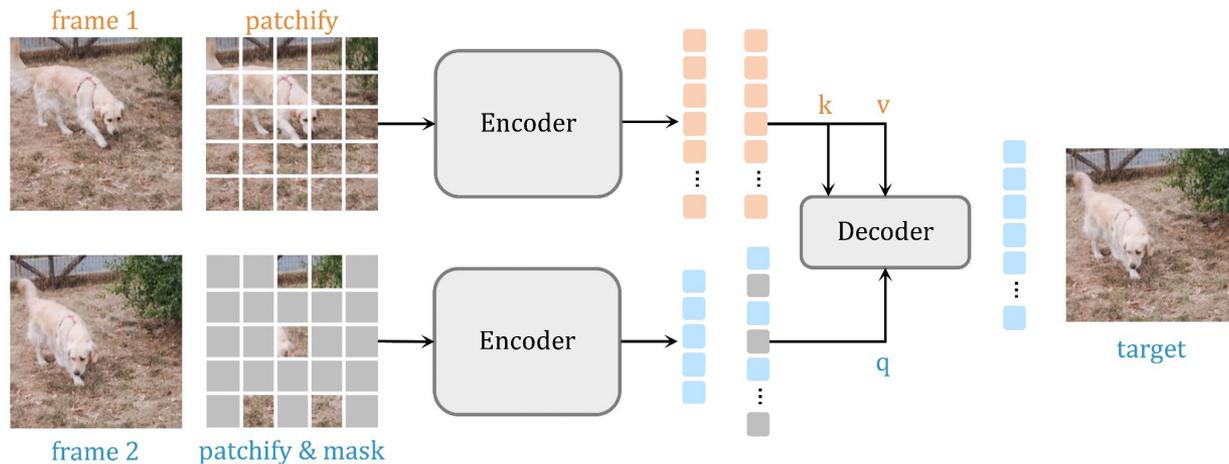
Core Idea

Predict missing information from context

Examples

- MAE (Masked Autoencoders)
- BEiT
- Masked image modeling

Captures structure without labels



Contrastive Learning

Core Idea

Pull positives together

Push negatives apart

Frameworks

- SimCLR
- MoCo
- DINO (self-distillation)

Strong transfer performance

Weak Supervision at Scale

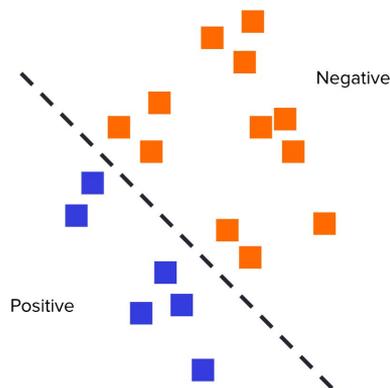
Label Sources

- Slide-level diagnosis
- Radiology reports
- Billing codes

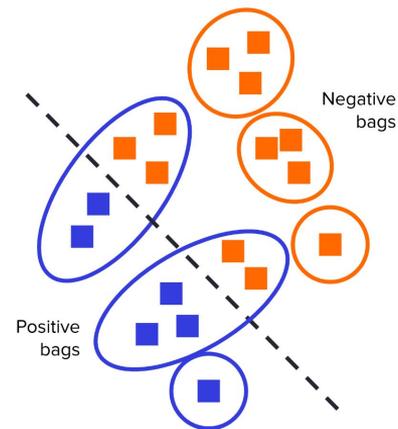
Tradeoff

Scale \uparrow
Label quality \downarrow

Traditional Supervised Learning



Multiple Instance Learning

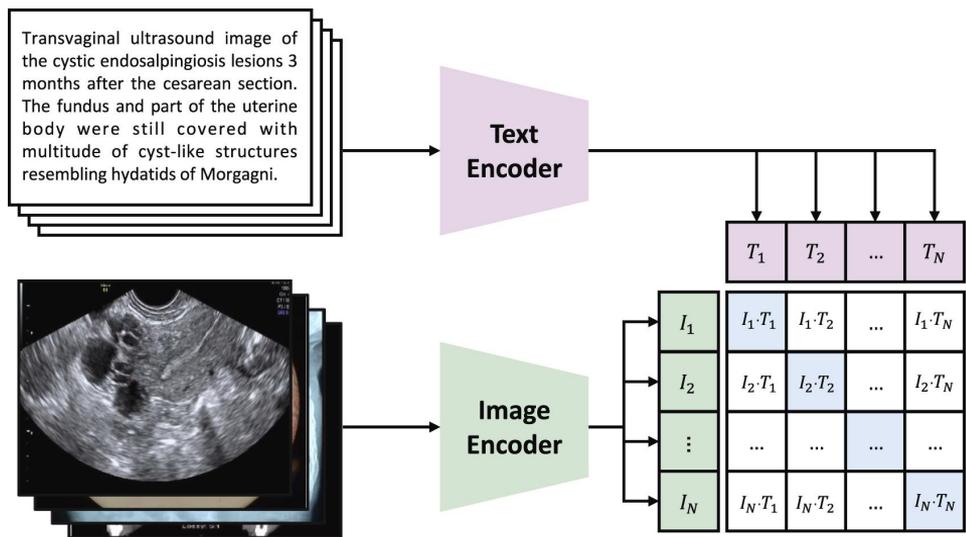


Multimodal Alignment

Align images to language or omics

Image encoder \leftrightarrow Text/omics encoder

Toward semantic representations



Lecture Outline

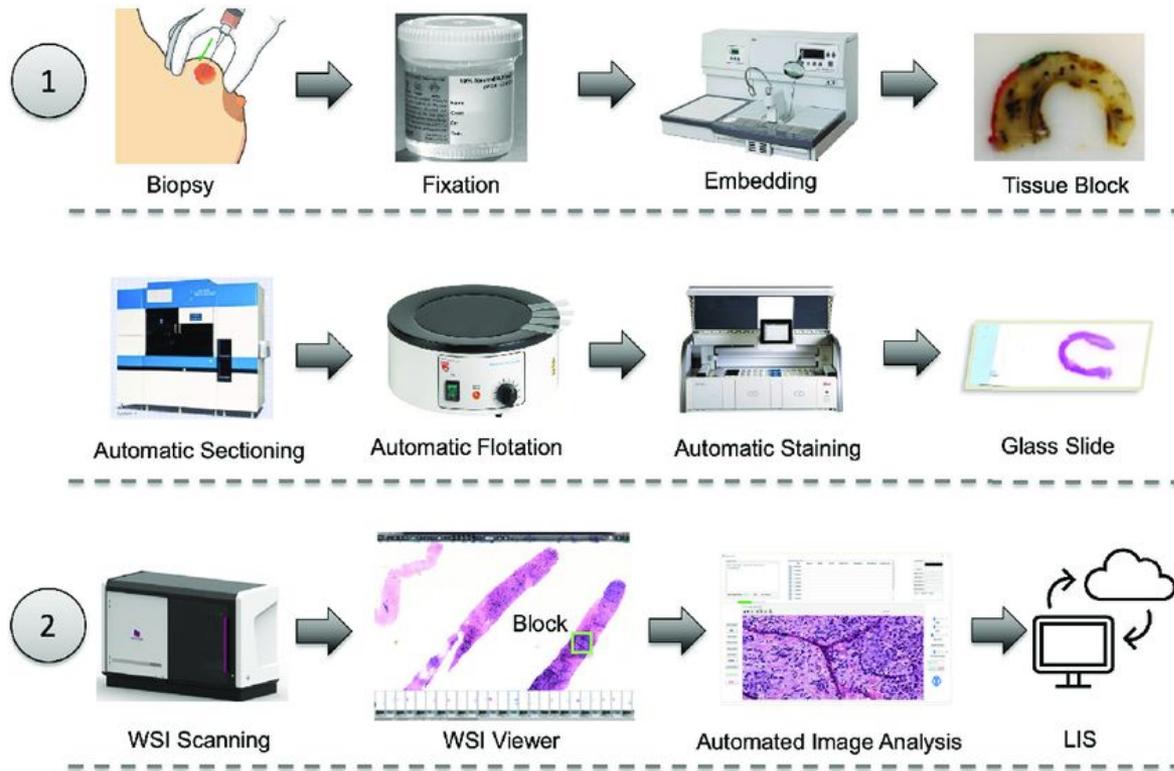
- ~~1. Foundation Models~~
2. Pathology Foundational Models
3. Model Interpretability
4. Model Generalizability

Why Pathology Became the Scaling Laboratory

Digital Pathology Enables Scale

- Gigapixel whole-slide images
- Routine digitization → massive archives
- Slide-level labels available at scale

Pathology became the foundation-model sandbox



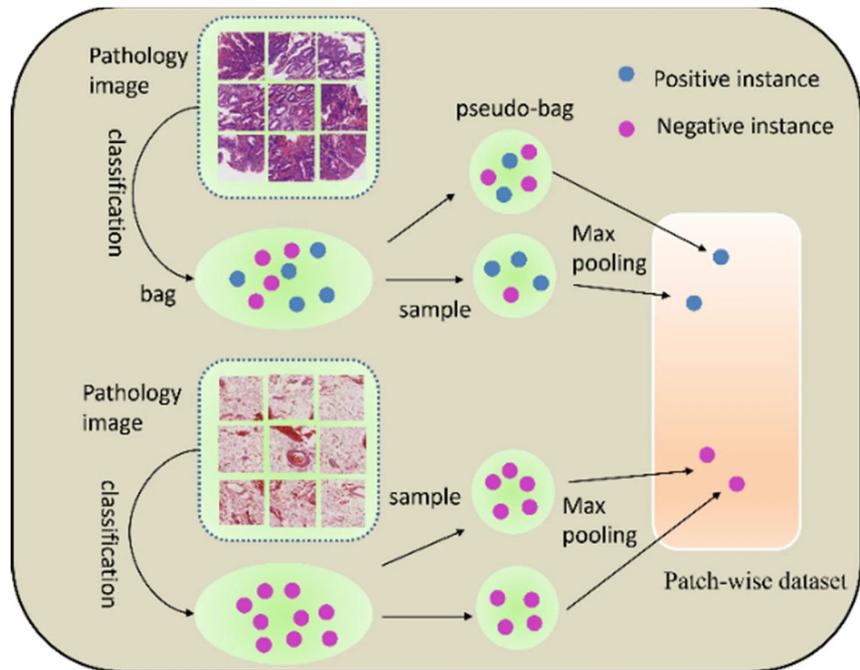
1) Tissue processing, 2) WSI acquisition and automated image analysis using WSI

Whole-Slide Images as “Bags of Patches”

WSI → Many Patches → One Label

- Supervision at slide level
- Signal is sparse
- Aggregation is the challenge

Multiple-instance learning framing

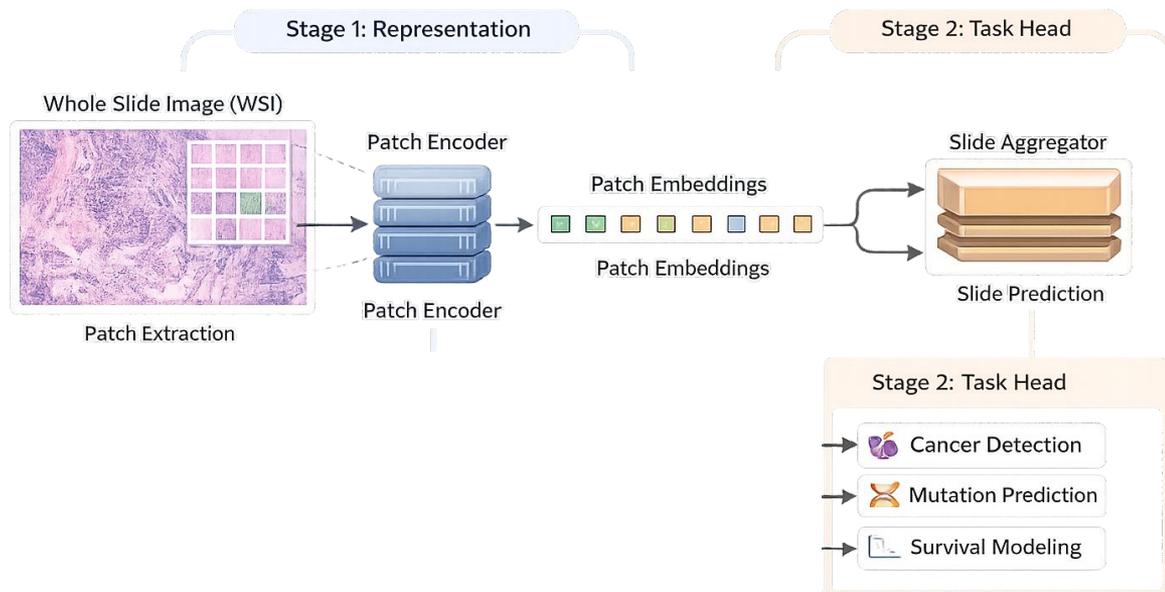


Patch Encoder + Slide Aggregator

Two-Stage System

1. Patch encoder (representation learning)
2. Slide aggregator (task-specific pooling)

Foundation move: make encoder reusable



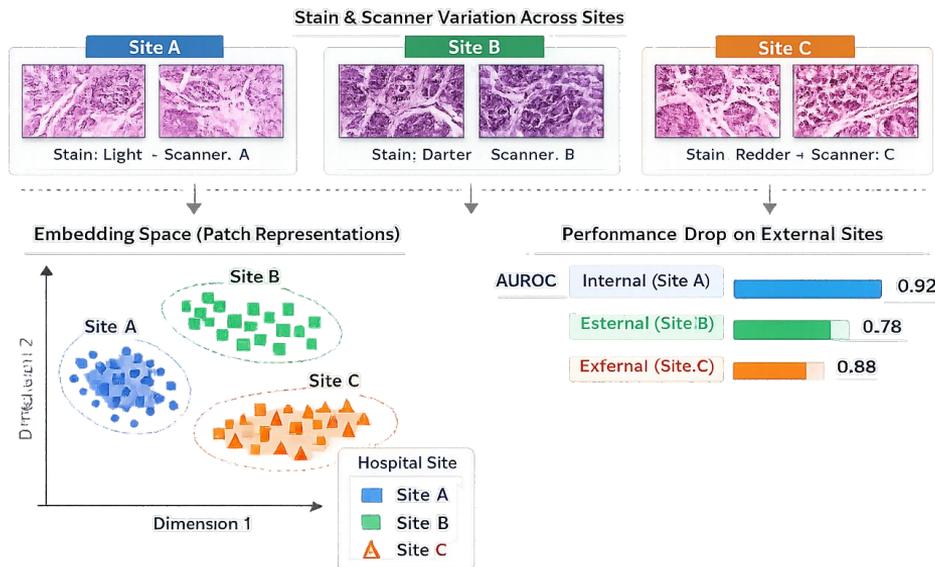
Why Generalization is Difficult in Pathology

Sources of Hidden Shift

- Stain variation
- Scanner variation
- Tissue processing differences
- Site-specific protocols

Shortcut risk: predicting hospital instead of disease

Hidden Shift in Digital Pathology



Three Papers, One Story

Article | Published: 04 September 2024

A pathology foundation model for cancer diagnosis and prognosis prediction

[Xiyue Wang](#), [Junhan Zhao](#), [Eliana Marostica](#), [Wei Yuan](#), [Jietian Jin](#), [Jiayu Zhang](#), [Ruijiang Li](#), [Hongping Tang](#), [Kanran Wang](#), [Yu Li](#), [Fang Wang](#), [Yulong Peng](#), [Junyou Zhu](#), [Jing Zhang](#), [Christopher R. Jackson](#), [Jun Zhang](#), [Deborah Dillon](#), [Nancy U. Lin](#), [Lynette Sholl](#), [Thomas Denize](#), [David Meredith](#), [Keith L. Ligon](#), [Sabina Signoretti](#), [Shuji Ogino](#), ... [Kun-Hsing Yu](#)  [+ Show authors](#)

[Nature](#) **634**, 970–978 (2024) | [Cite this article](#)

70k Accesses | 385 Citations | 388 Altmetric | [Metrics](#)

Scale → strong transfer (diagnosis, prognosis)

“General-purpose”
=

task breadth + external validation

Article | Published: 19 March 2024

Towards a general-purpose foundation model for computational pathology

[Richard J. Chen](#), [Tong Ding](#), [Ming Y. Lu](#), [Drew F. K. Williamson](#), [Guillaume Jaume](#), [Andrew H. Song](#), [Bowen Chen](#), [Andrew Zhang](#), [Daniel Shao](#), [Muhammad Shaban](#), [Mane Williams](#), [Lukas Oldenburg](#), [Luca L. Weishaupt](#), [Judy J. Wang](#), [Anurag Vaidya](#), [Long Phi Le](#), [Georg Gerber](#), [Sharifa Sahaj](#), [Walt Williams](#) & [Faisal Mahmood](#) 

[Nature Medicine](#) **30**, 850–862 (2024) | [Cite this article](#)

89k Accesses | 990 Citations | 248 Altmetric | [Metrics](#)

Perspective | Published: 12 April 2023

Foundation models for generalist medical artificial intelligence

[Michael Moor](#), [Oishi Banerjee](#), [Zahra Shakeri Hossein Abad](#), [Harlan M. Krumholz](#), [Jure Leskovec](#), [Eric J. Topol](#)  & [Pranav Rajpurkar](#) 

[Nature](#) **616**, 259–265 (2023) | [Cite this article](#)

254k Accesses | 1454 Citations | 751 Altmetric | [Metrics](#)

Generalist medical AI across modalities

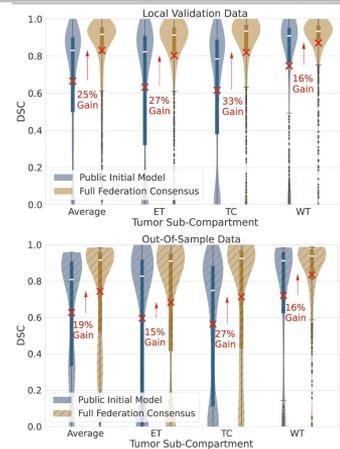
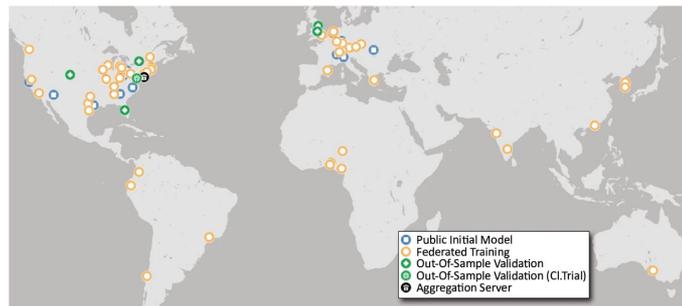
How to Read a “Foundation Model” Paper

When evaluating claims, check:

- Data scale & diversity
- Pretraining objective
- Transfer protocol
- External validation design

“General-purpose” is an empirical claim

- ✓ Data Scale
- ✓ Objective
- ✓ Transfer
- ✓ External Validation



Lecture Outline

- ~~1. Foundation Models~~
- ~~2. Pathology Foundational Models~~
3. Model Interpretability
4. Model Generalizability

Why Interpretability Exists

Interpretability

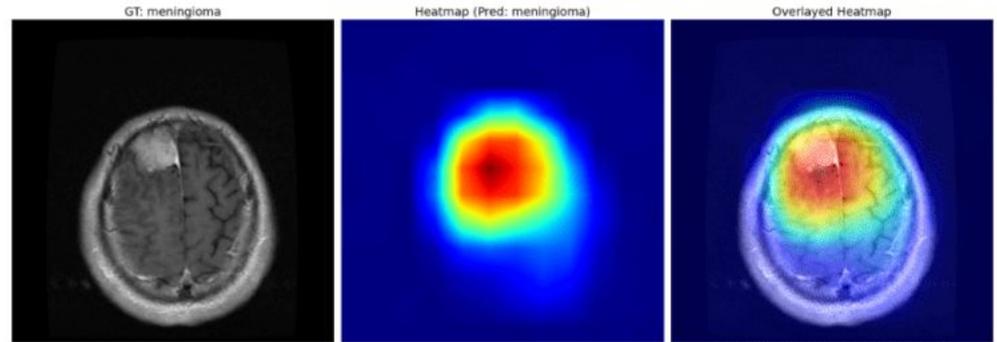
Methods that help humans understand *what a model uses to make predictions*

Why it exists:

- Debugging & failure detection
- Communication & trust

Not proof of causality

Keep it clean. No more than that.



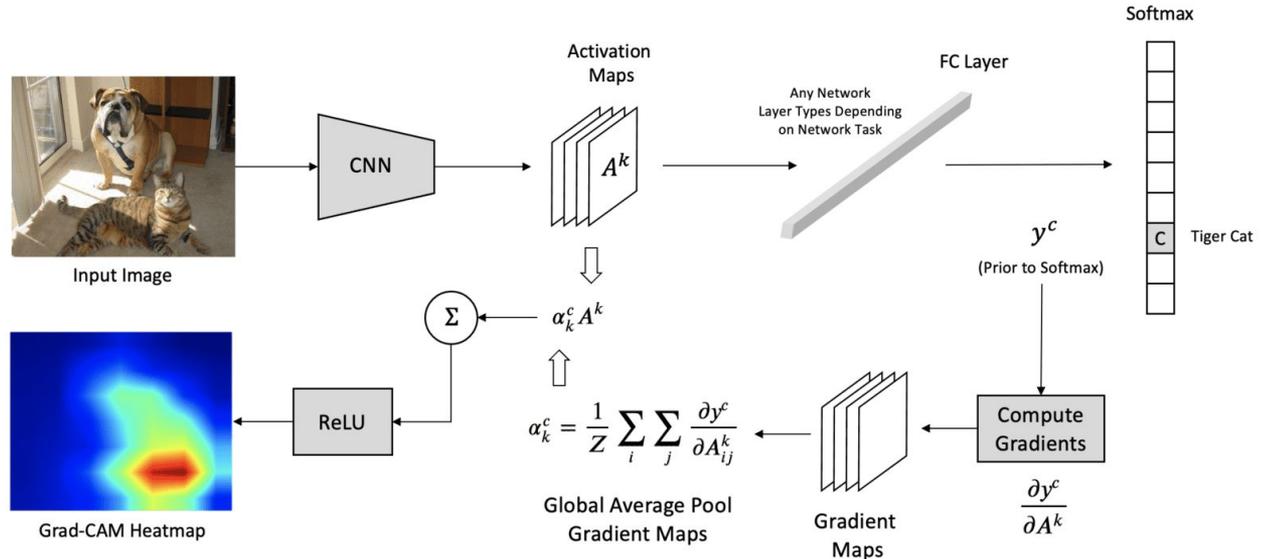
Attribution Methods: What They Show (and What They Don't)

Common attribution tools:

- Saliency / Grad-CAM
- Attention maps
- Occlusion tests

They measure **prediction sensitivity**, not biological reasoning.

Keep that last line visually separated.



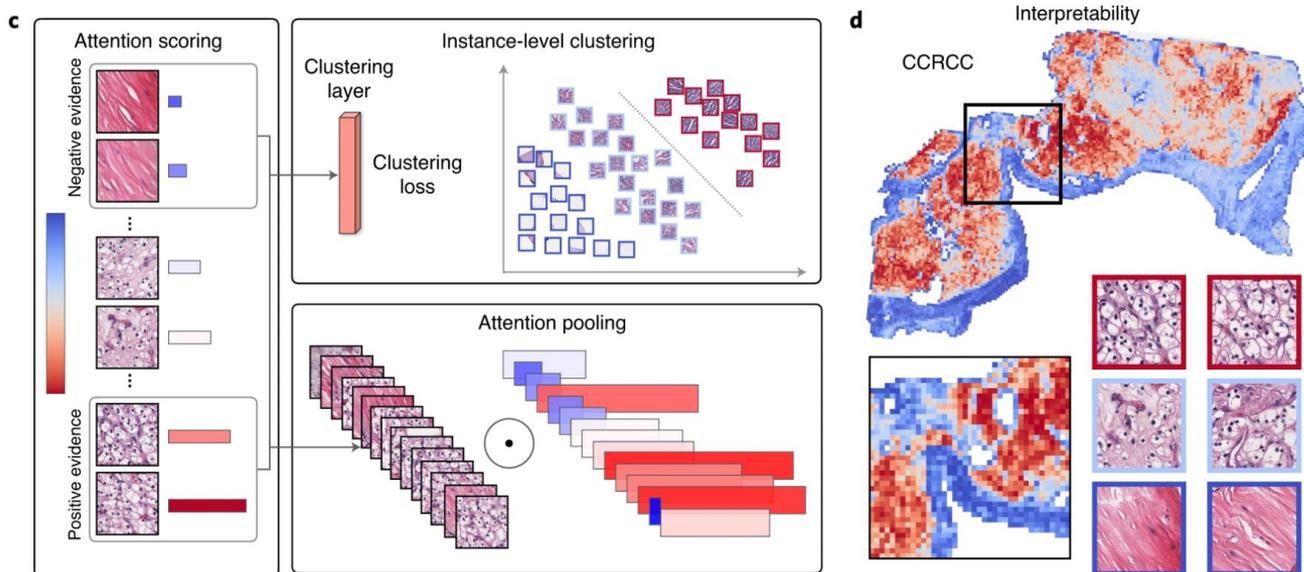
Example — Attention-Based MIL in Pathology (CLAM)

Clustering-constrained Attention MIL (CLAM)

- Attention-based multiple instance learning
- Identifies discriminative regions
- Enables weakly supervised localization

Attention \neq causation

Keep it tight. Let the figure do most of the work.

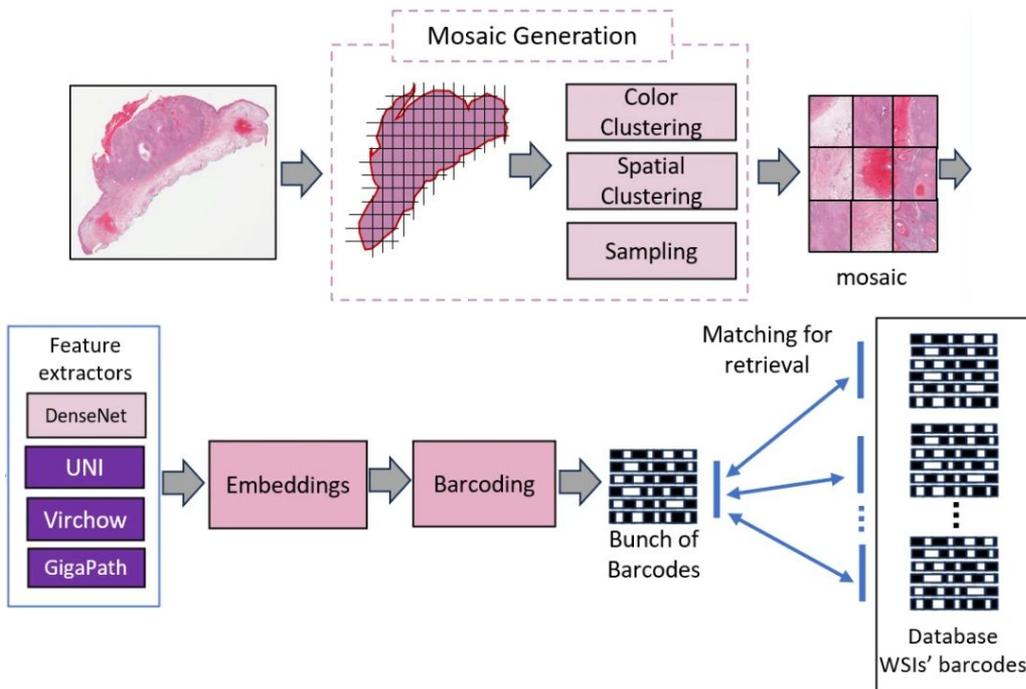


Retrieval-Based Interpretability in Pathology

Zero-shot whole-slide
retrieval from learned
embeddings

Retrieval provides *case
evidence*, but:

Similarity \neq diagnosis



Lecture Outline

- ~~1. Foundation Models~~
- ~~2. Pathology Foundational Models~~
- ~~3. Model Interpretability~~
4. Model Generalizability

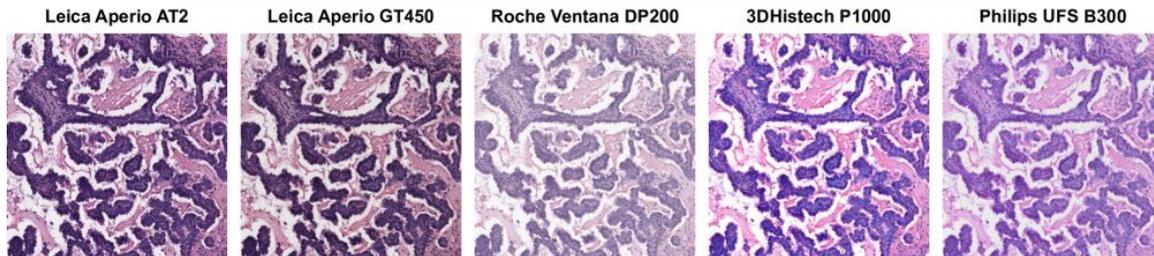
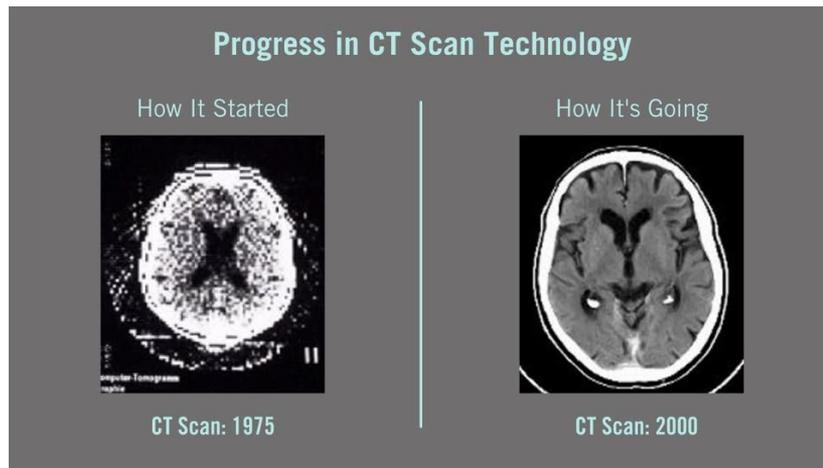
Why Models Fail Outside the Lab

Distribution shift is the norm

Common shifts:

- Site / scanner / protocol
- Temporal
- Population
- Label & prevalence

Shortcut learning amplifies shift



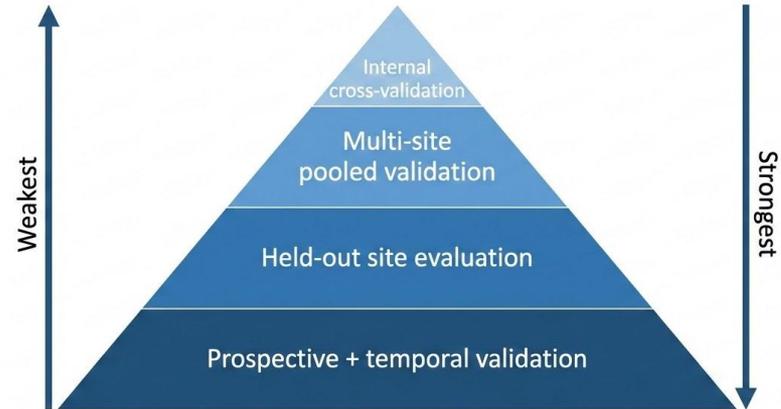
What Robust Evaluation Actually Requires

Robust evaluation includes:

- External site holdout
- Temporal split
- Subgroup performance
- Calibration under shift

Stress tests > average AUROC

Validation Patterns: Hierarchy of Evaluation Strength



Federated Learning

Motivation

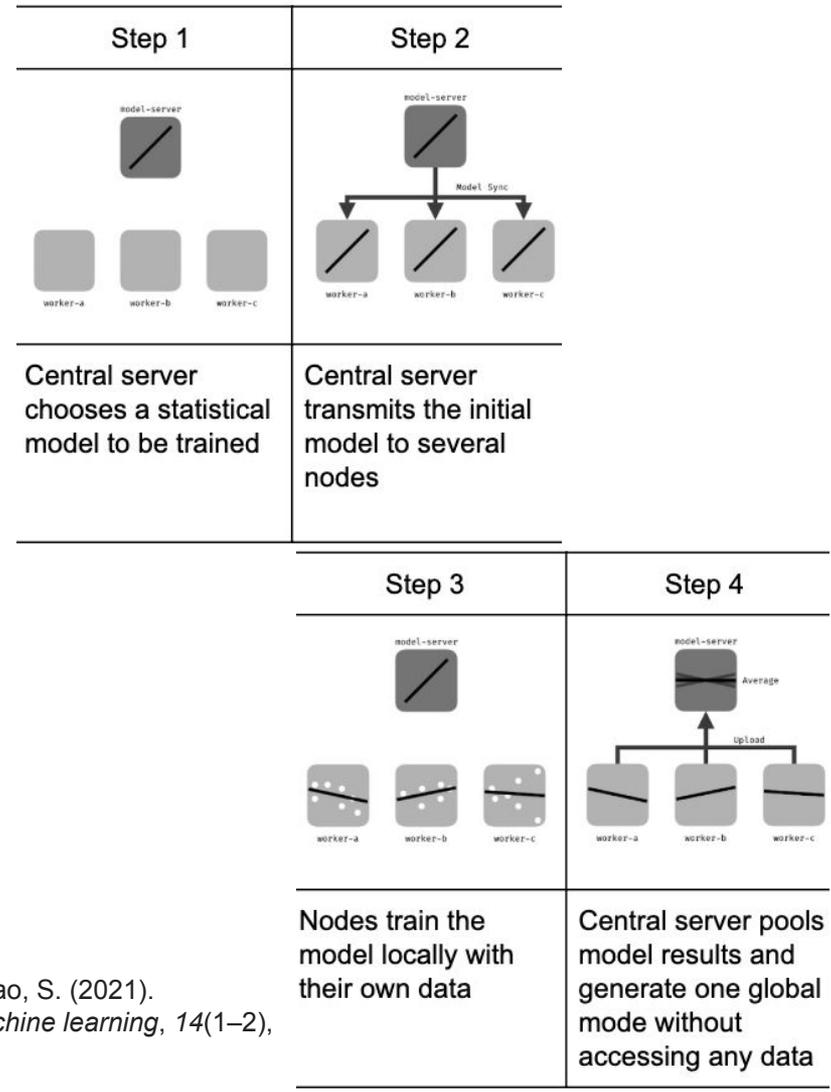
- Train collaboratively across hospitals without sharing raw data
- Larger effective dataset, privacy preserved

Challenges

- Communication overhead, model synchronization
- Data heterogeneity (different scanners, protocols)

Potential Impact

- Improved model generalizability
- Regulatory compliance (HIPAA, GDPR)



Federated Learning

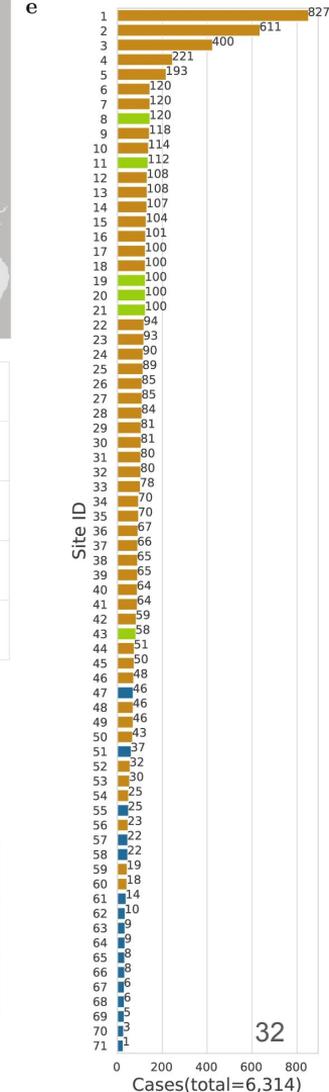
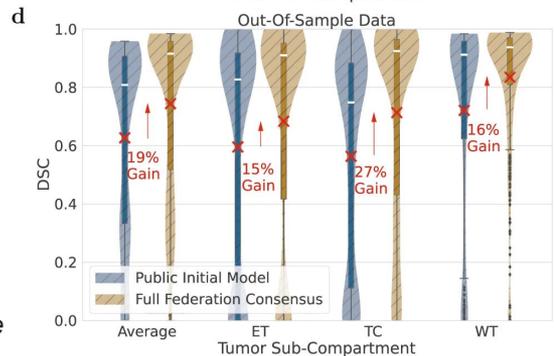
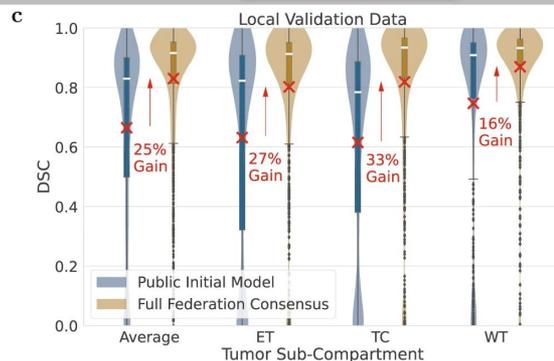
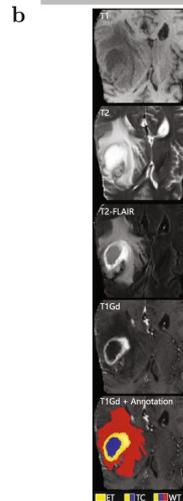
- **Largest FL study in medical imaging** → 71 sites, 6 continents
- **6,314 cases** → Largest glioblastoma dataset
- **No data sharing** → Privacy-preserving model training

Key Results

- **+33% improvement** in surgically targetable tumor segmentation
- **+23% improvement** in complete tumor segmentation
- **Validated on:**
 - **Local site data (n = 1,043 cases)**
 - **Out-of-sample data (n = 518 cases)**

Impact

- **More diverse, generalizable AI models**
- **Public release of the consensus model**
- **New standard for multi-site AI training**



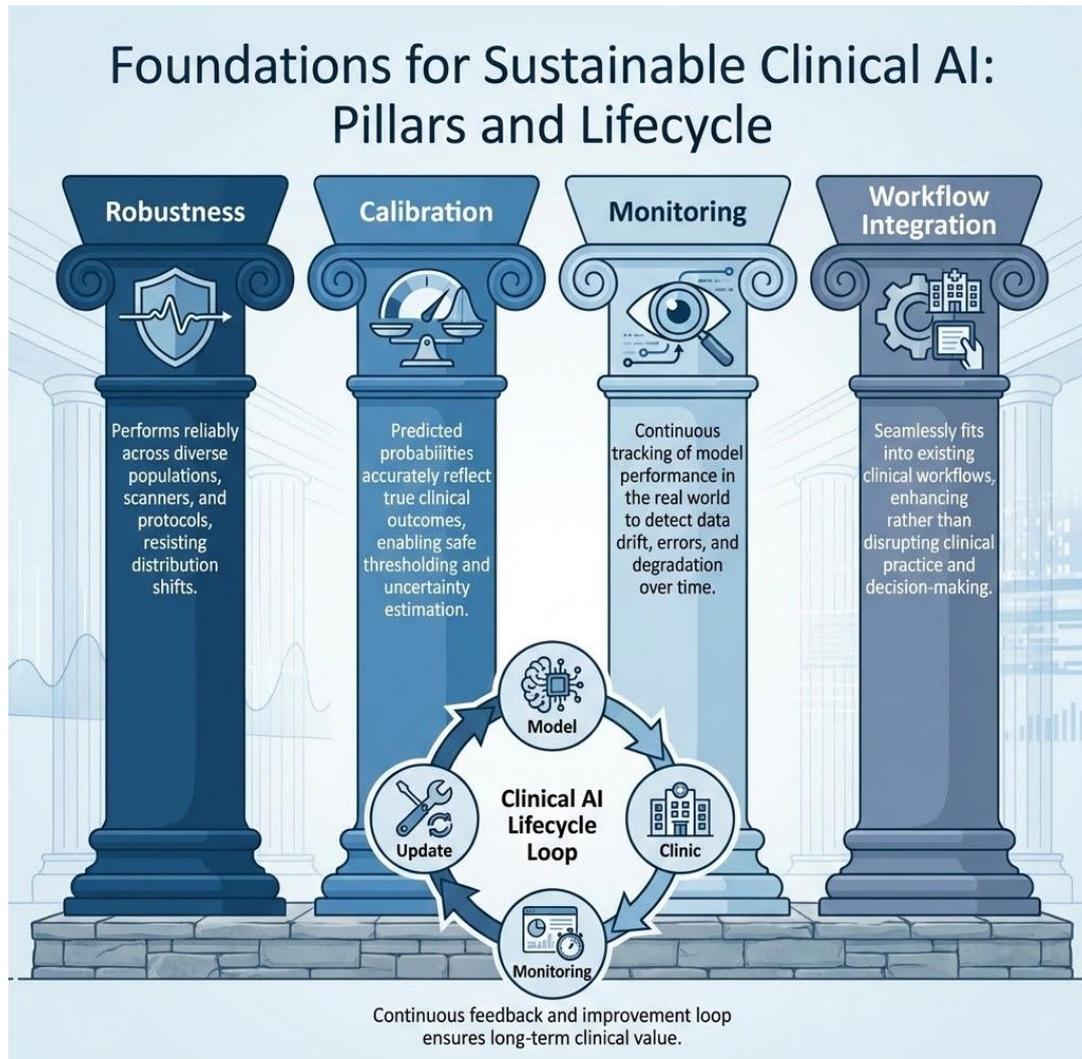
Pati, S., Baid, U., Edwards, B., Sheller, M., Wang, S. H., Reina, G. A., ... & Poisson, L. (2022). Federated learning enables big data for rare cancer boundary detection. *Nature communications*, 13(1), 7346.

The Real Bar for “Generalist” AI

A generalist model must be:

- Stable across sites
- Calibrated under shift
- Transparent about failures
- Continuously monitored

Representation quality \neq clinical safety



The End