# Today's Focus

# Who am I?

- Heather Couture
- MS from Carnegie Mellon University - autonomous science
- PhD from University of North Carolina - computational pathology
- Computer vision consultant

OWKIN    DECIPHEX    ULTIVUE    WattTime    Leica MICROSYSTEMS    QRITIVE

- Keynote speaker at MICCAI workshop on computational pathology
- Contributor to Scientific American, The Pathologist, IEEE Spectrum
- Newsletter and podcast

**Computer Vision Insights** by Pixel Scientia Labs
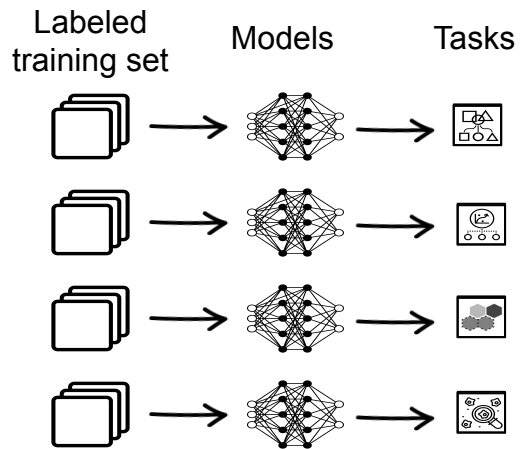
IMPACT AI
WITH HEATHER COUTURE
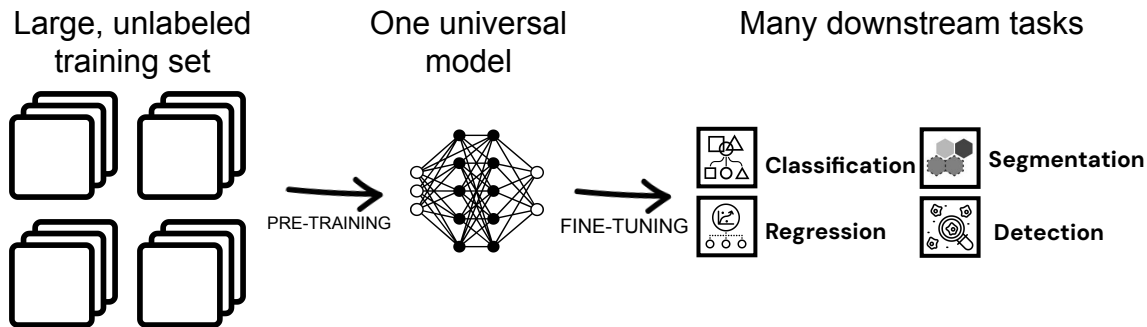
Task-Specific

Foundation Models

Multimodal

Agents

- New domains
- Domain-specific adaptations
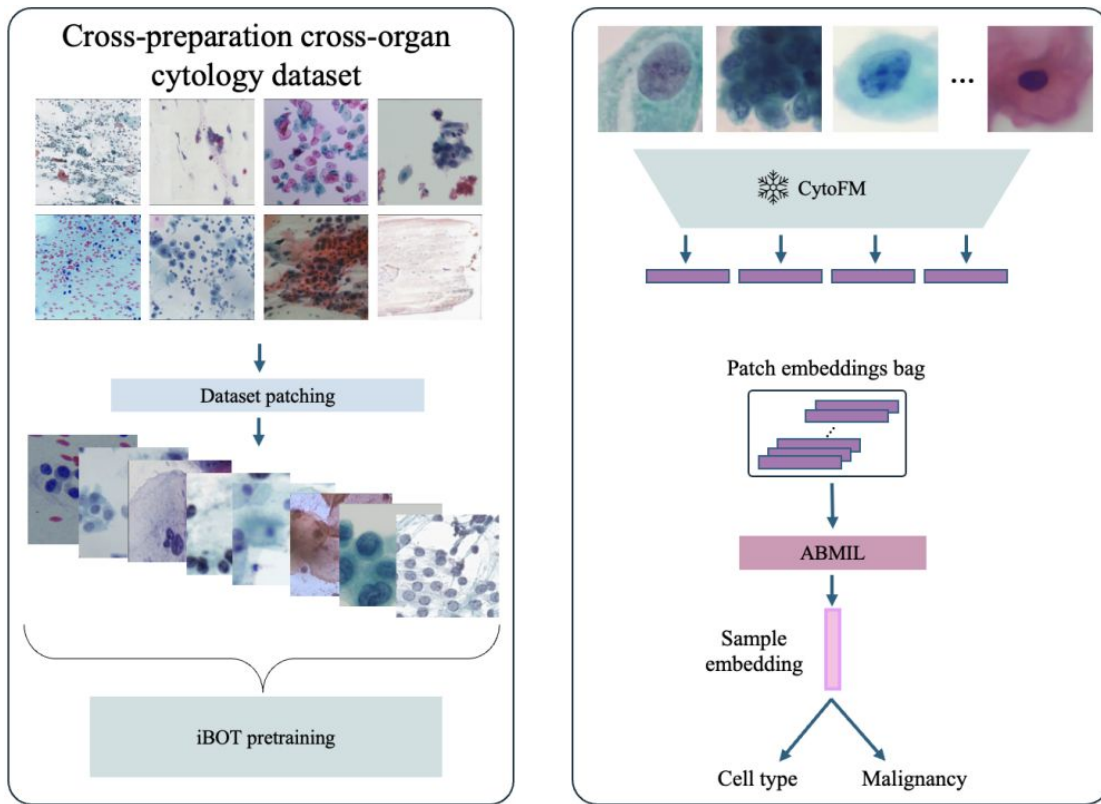- New datasets
- New benchmarks

# The Paradigm Shift

## Traditional ML

Labeled
training set

Models

Tasks

## Foundation Models

Large, unlabeled
training set

One universal
model

Many downstream tasks

PRE-TRAINING

FINE-TUNING

**Classification**

**Segmentation**

**Regression**

**Detection**

# A Foundation Model for Cytology



CytoFM

1.4 million image patches
8 datasets
7 institutions

Ivezic, CytoFM: The first cytology foundation model, 2025

# A Foundation Model for Spatial Proteomics



KRONOS

47 million image patches
175 protein markers
16 tissue types
8 imaging platforms

Shaban, A Foundation Model for Spatial Proteomics,2025

# A Foundation Model for Agriculture

Agri-FM+



Nahian, Agri-FM+: A Self-Supervised Foundation Model for Agricultural Vision, 2025
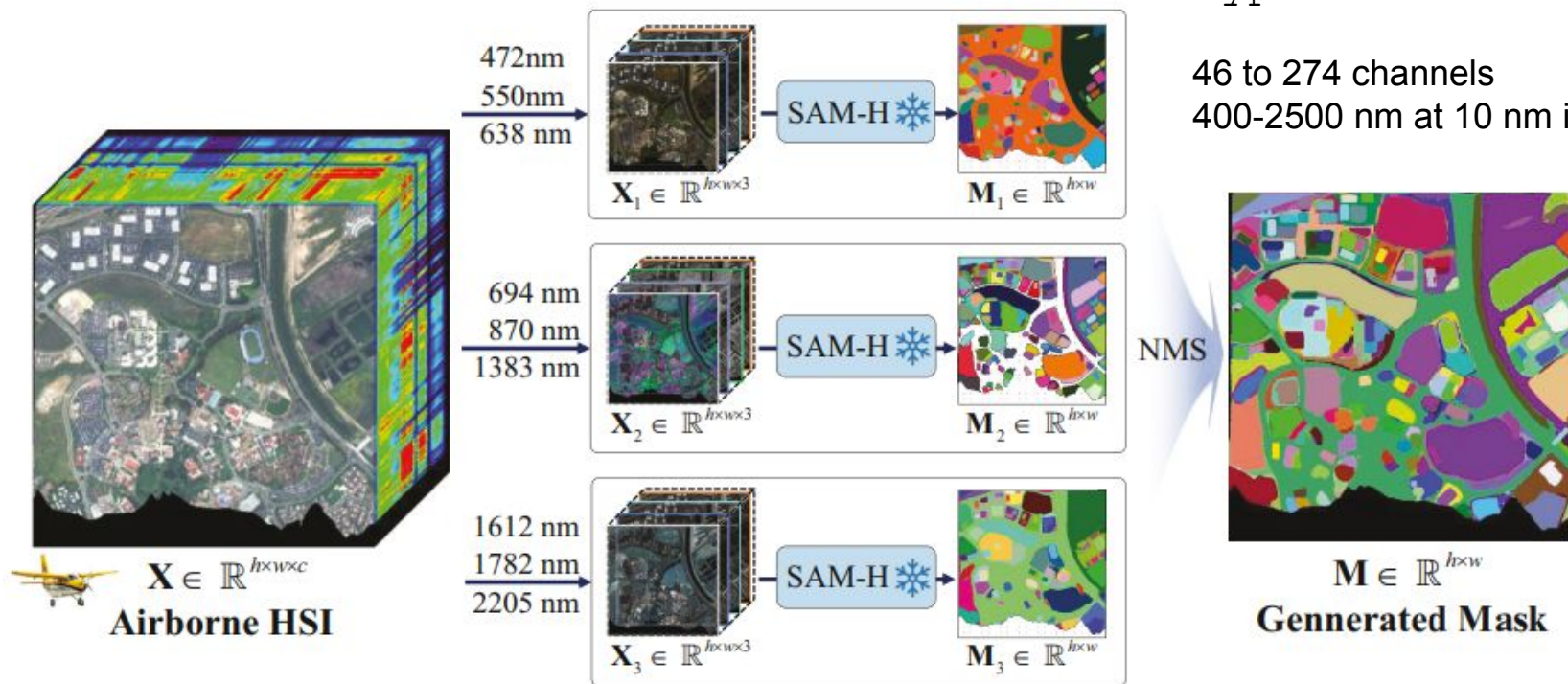
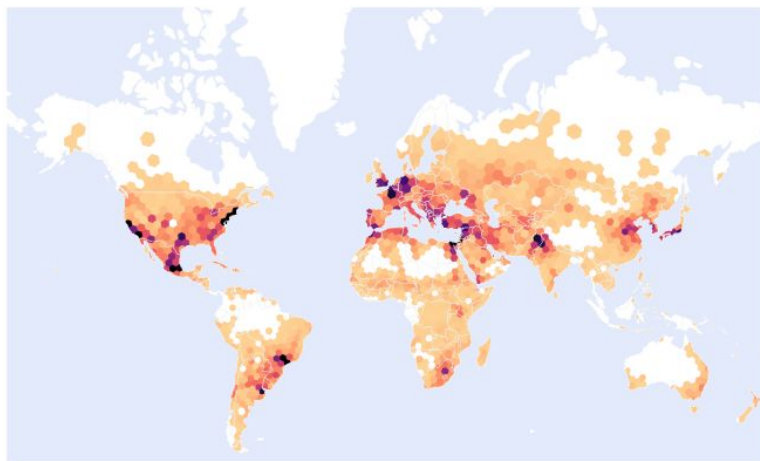# A Foundation Model for Hyperspectral Imagery



HyperFree

46 to 274 channels
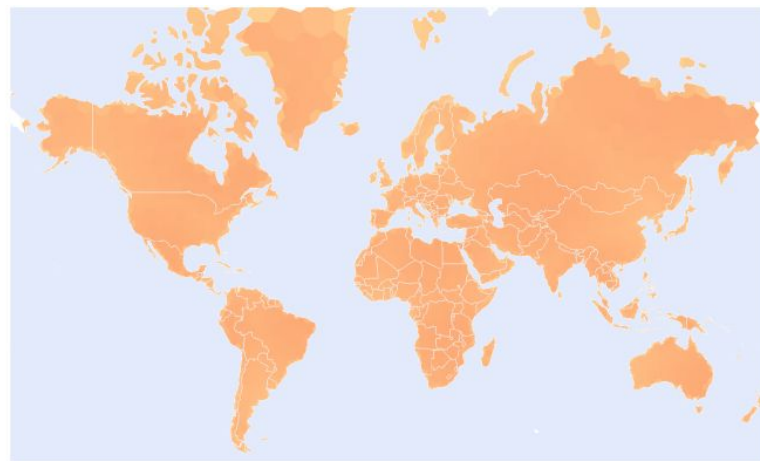400-2500 nm at 10 nm intervals

Li, HyperFree: A Channel-adaptive and Tuning-free Foundation
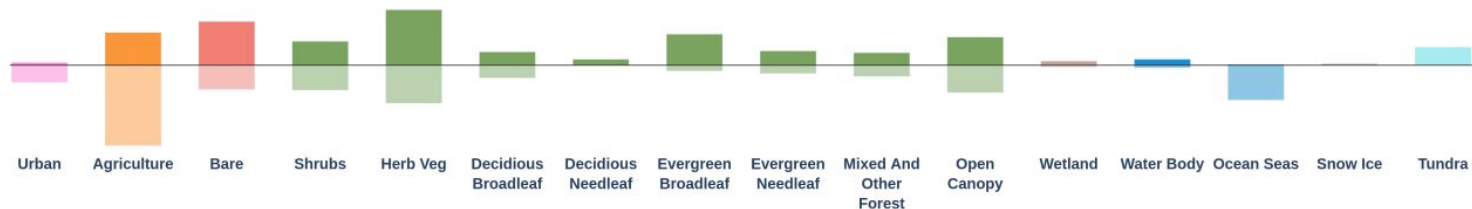Model for Hyperspectral Remote Sensing Imagery, 2025

# Remote Sensing Dataset: Global Distribution



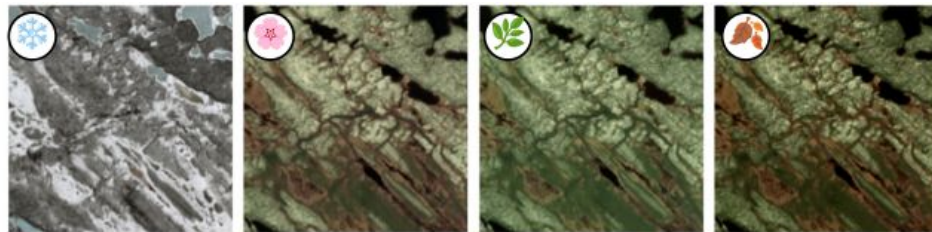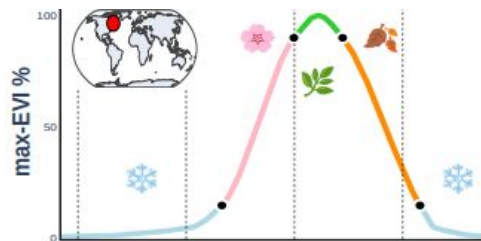(a) Spatial distribution of SSL4EO-S12 [89]

(b) Spatial distribution of SSL4Eco

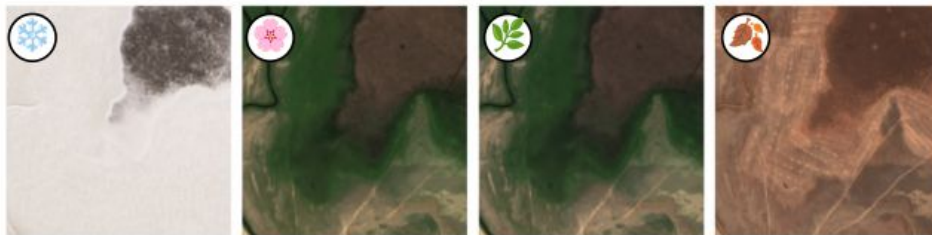Urban    Agriculture    Bare    Shrubs    Herb Veg    Decidious Broadleaf    Decidious Needleaf    Evergreen Broadleaf    Evergreen Needleaf    Mixed And Other Forest    Open Canopy    Wetland    Water Body    Ocean Seas    Snow Ice    Tundra

(c) Copernicus land cover [55] distribution for SSL4Eco (upwards) and SSL4EO-S12 [89] (downwards)

Plekhanova, SSL4Eco: A Global Seasonal Dataset for Geospatial Foundation Models in Ecology, 2025

# Remote Sensing Dataset: Seasonal Distribution
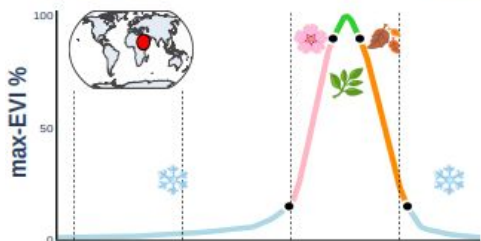


(a) EVI-based seasons

(b) Seasonal images

`SSL4Eco`

250k locations across entire landmass

Enhanced Vegetation Index-based seasonal sampling

`SeCo-Eco`

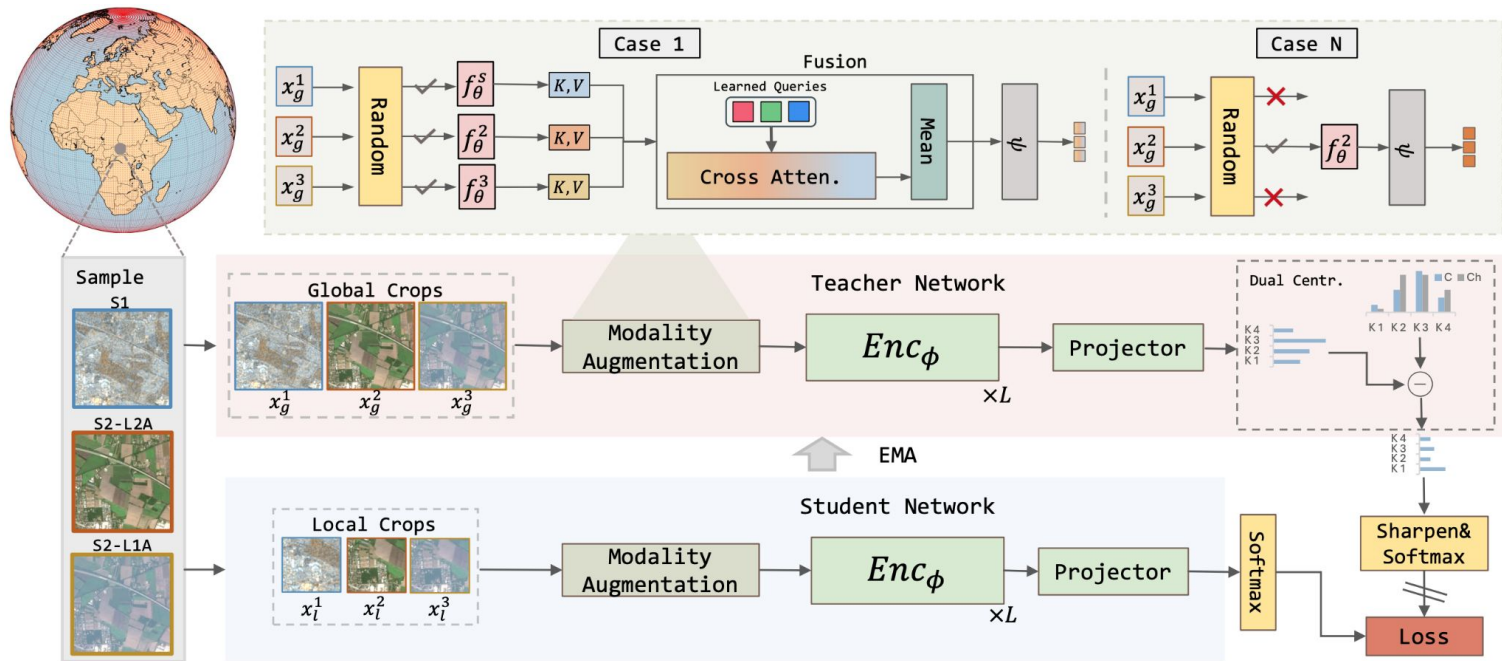Seasonal contrastive learning

Plekhanova, SSL4Eco: A Global Seasonal Dataset for Geospatial Foundation Models in Ecology, 2025

Task-Specific → Foundation Models → **Multimodal** → Agents

- Open weights and open data
- Grounding
- VLMs - language is the "glue"
- Multi-stage pretraining
- Larger inputs
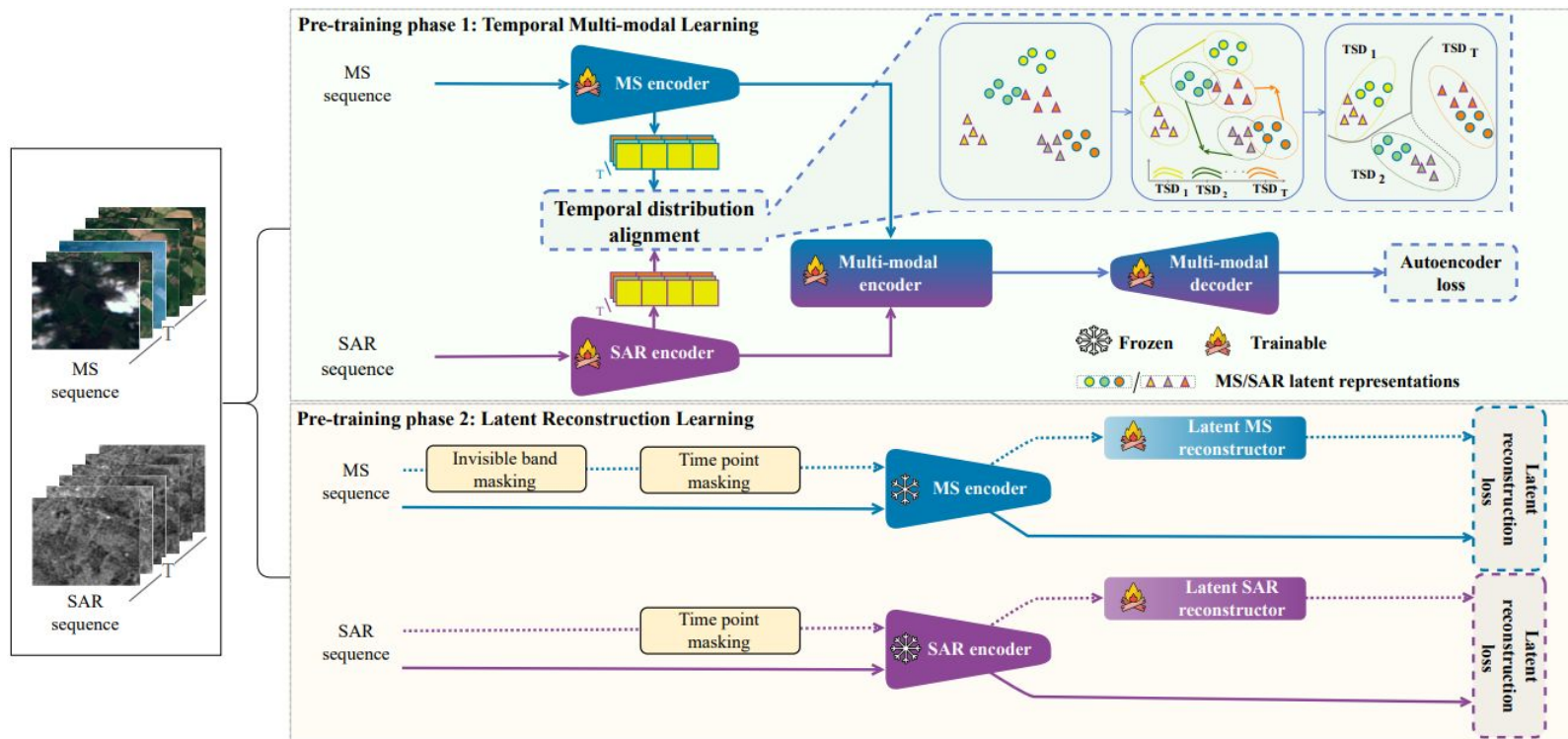- Benchmarks

# Remote Sensing: Optical + Radar



TerraFM

SAR and optical
18.7 million tiles
534x534 tiles

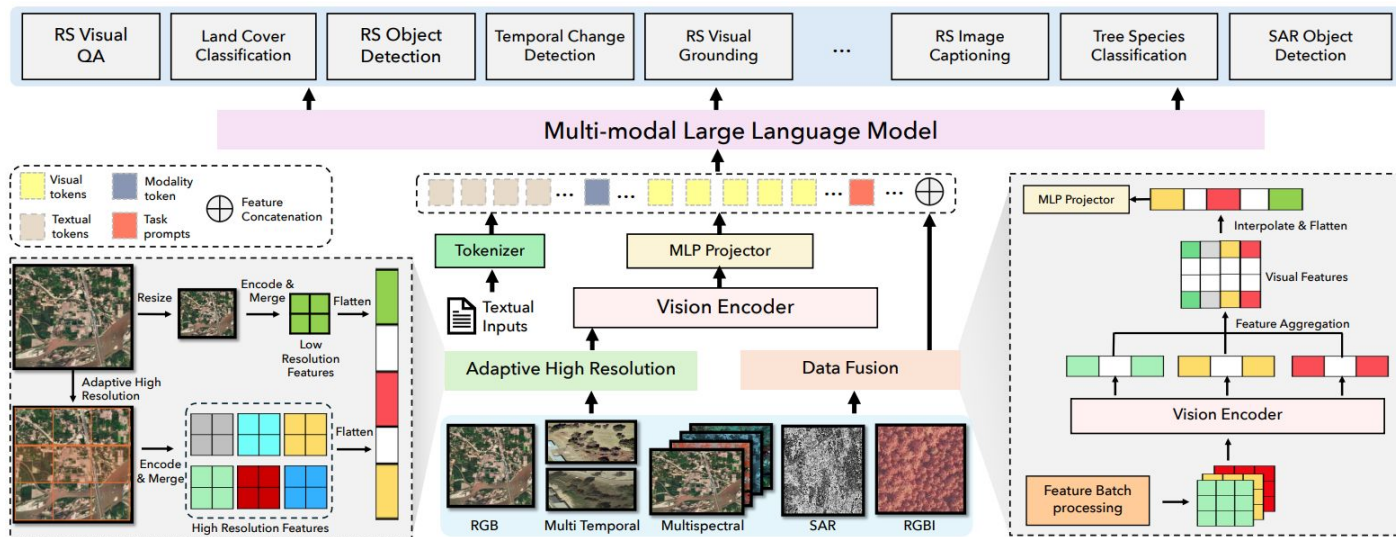Danish, TerraFM: A Scalable Foundation Model for Unified Multisensor Earth Observation, 2025

# Remote Sensing: Robust to Missing Data

RoboSense



Do, RobSense: A Robust Multi-modal Foundation Model for Remote Sensing
with Static, Temporal, and Incomplete Data Adaptability, 2025

# A VLM for Remote Sensing



EarthDial

RGB, multispectral, infrared, SAR
11 million instruction pairs

Soni, EarthDial: Turning Multi-sensory Earth Observations to Interactive Dialogues, 2025

# Benchmarking VLMs for Remote Sensing
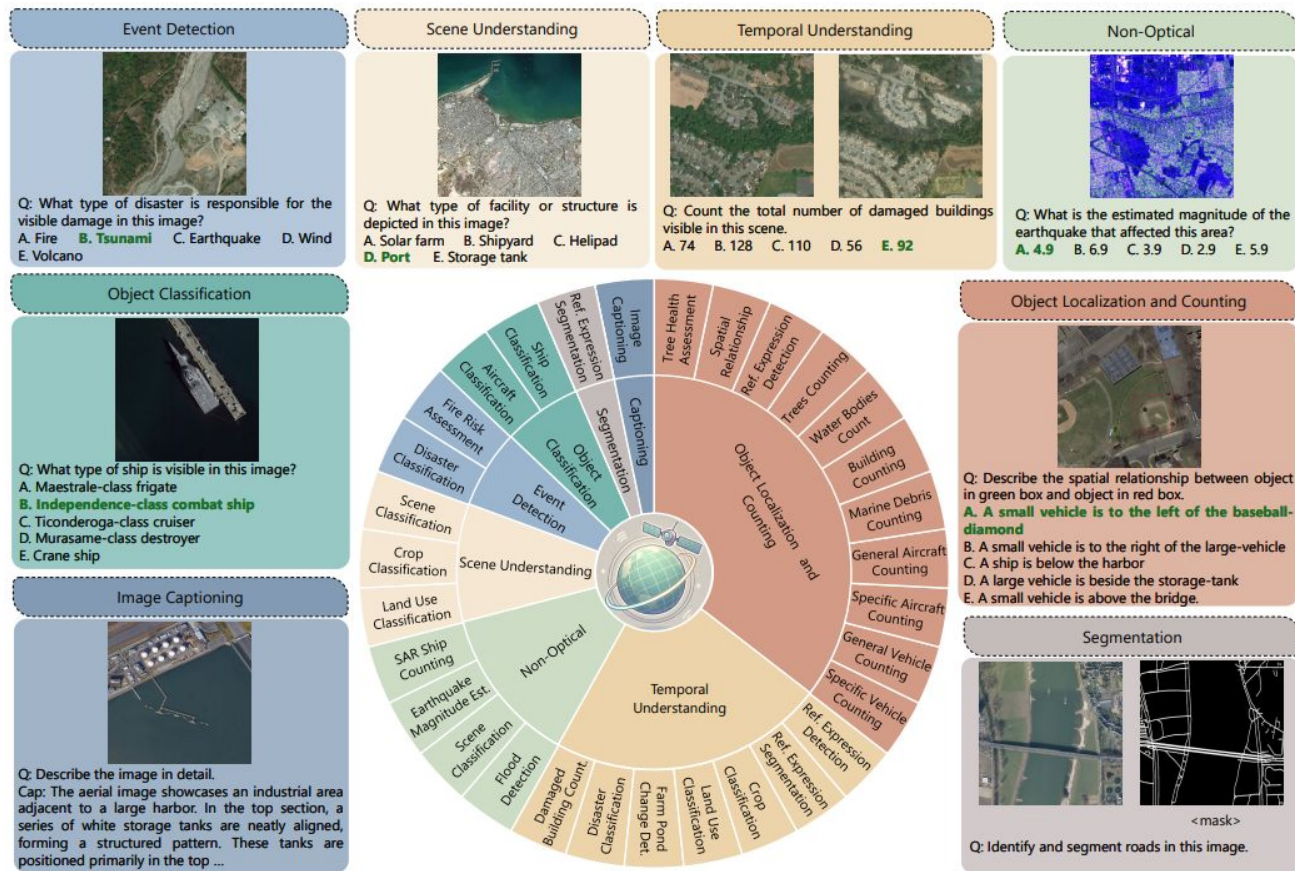
GEOBench-VLM

31 fine-grained tasks
8 categories
Optical, multispectral, SAR,
temporal
10+ manually verified instructions
Multiple-choice format

Danish, GEOBench-VLM:
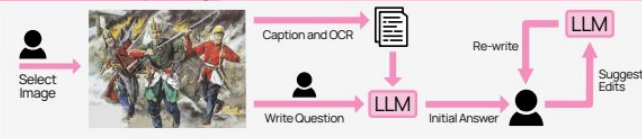Benchmarking Vision-Language
Models for Geospatial Tasks, 2025

# A VLM with Open Weights and Open Data



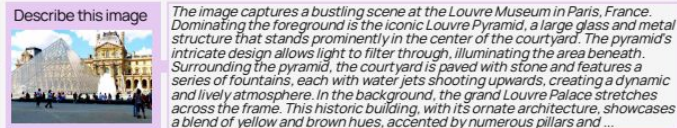Deitke, Molmo and PixMo: Open Weights and Open Data for State-of-the-Art Vision-Language Models, 2025

# A VLM with Open Weights and Open Data + Grounding



Deitke, Molmo and PixMo: Open Weights and Open Data for State-of-the-Art Vision-Language Models, 2025

# Grounding for Agriculture



**Prompt: wheat spike head**

Singh, Few-Shot Adaptation of Grounding DINO for Agricultural Domain, 2025

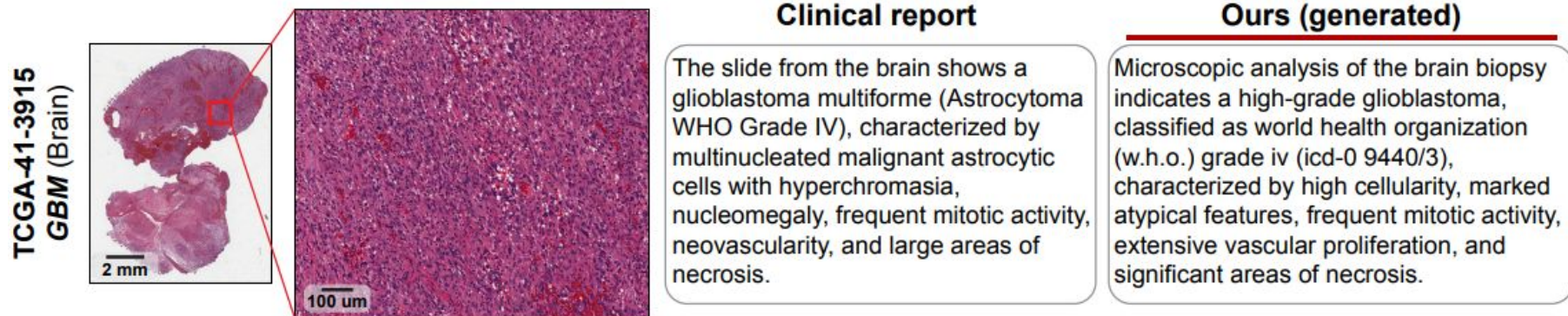# Remote Sensing VLM with Grounding

GeoPixel



Can you give a thorough description of this image, including interleaved segmentation masks to highlight key objects?

The image is an aerial view of a sports complex featuring a running track and a swimming pool. The track has multiple lanes, and the field is marked for various sports. The surrounding area includes roads, buildings, and green spaces. The ground track field at the center is a large athletic venue with a multi lane running track encircling a grassy field marked for various sports. The swimming pool at the top is large, rectangular, and filled with a blue substance, likely water. Numerous small vehicles can be seen at different locations. A soccer field within the sports complex features a rectangular pitch with goals at each end, surrounded by a running track. The image gives a static and undisturbed appearance as no people or moving objects are visible.

Shabbir, GeoPixel: Pixel Grounding Large Multimodal Model in Remote Sensing, 2025

# A VLM for Whole Slide Images - Pathology Reports



**Clinical report**

The slide from the brain shows a glioblastoma multiforme (Astrocytoma WHO Grade IV), characterized by multinucleated malignant astrocytic cells with hyperchromasia, nucleomegaly, frequent mitotic activity, neovascularity, and large areas of necrosis.

**Ours (generated)**

Microscopic analysis of the brain biopsy indicates a high-grade glioblastoma, classified as world health organization (w.h.o.) grade iv (icd-0 9440/3), characterized by high cellularity, marked atypical features, frequent mitotic activity, extensive vascular proliferation, and significant areas of necrosis.

Ding, Multimodal Whole Slide Foundation Model for Pathology, 2025

TITAN

335k WSIs with 423k synthetic captions
8192 x 8192 tiles
Three-stage training:
     SSL on patches
     Vision-language alignment
     Alignment with pathology reports

# A VLM for Whole Slide Images - VQA



[Caption]
#Q: *Describe the whole slide image.*
#A: The bladder exhibits an invasive poorly differentiated urothelial carcinoma, grade 4/4, extending through the bladder neck wall......

[VQA-Diagnosis]
#Q: From the given whole slide image, which histological subtype of lung adenocarcinoma is most predominantly observed?
#A: It is acinar subtype observed.

[VQA-Microscopy]
#Q: What histological feature observable in a whole slide image differentiates anaplastic astrocytoma from glioblastoma ?
#A: Absence of endothelial proliferation.

[VQA-Clinical]
#Q: Based on the histopathological features, what is the primary recommended treatment?
#A: Simple mastectomy followed by hormone therapy.

Chen, SlideChat: A Large Vision-Language Assistant for Whole-Slide Pathology Image Understanding, 2025

SlideChat

4k slide captions
176k question-answer pairs
Two-stage training:
    Cross-domain alignment
    Visual instruction learning

# A VLM for Whole Slide Images - Multiple Use Cases



CPath-Omni

Virchow2 (DINOv2-based) + CLIP
Four-stage training:
    Vision-language alignment
    Patches: VQA, classification, captioning
    Whole slide pathology reports
    Slide and patch training

Sun, CPath-Omni: A Unified Multimodal Foundation Model for Patch and Whole Slide Image Analysis in Computational Pathology, 2025

Task-Specific → Foundation Models → Multimodal → Agents

- Multi-agent systems
- Specialization
- Orchestration
- Mimic how human experts work

# Agent for Histopathology Diagnosis

PathFinder



Ghezloo, PathFinder: A Multi-Modal Multi-Agent System for Medical Diagnostic Decision-Making Applied to Histopathology, 2025

# Agent for Clinical Decision-Making in Oncology



Ferber, Development and validation of an autonomous artificial intelligence agent for clinical decision-making in oncology, 2025

# Geospatial Agent

`GeoLLM-Squad`



Lee, Multi-Agent Geospatial Copilots for Remote Sensing Workflows, 2025

# Benchmark for Geospatial Agents

ThinkGeo



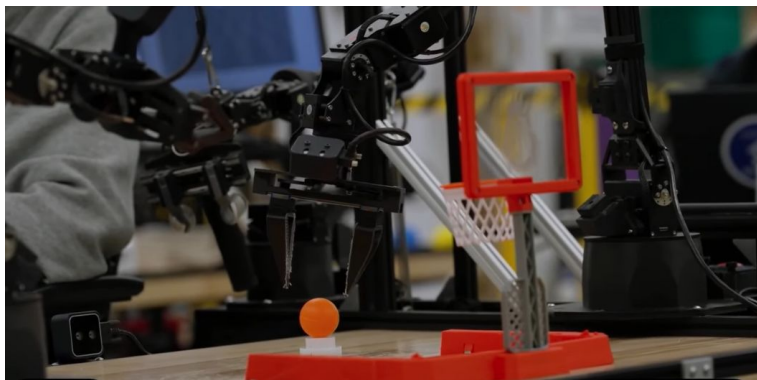Shabbir, ThinkGeo: Evaluating Tool-Augmented Agents for Remote Sensing Tasks, 2025

# Embodied Intelligence

"Can you pack me some trail mix?"

"Put the pen with the other pencils."



"Pick up the basketball and slam dunk it."

"Put the objects in the right container."



Keynote: Carolina Parada, Gemini Robotics, https://www.youtube.com/watch?v=o38k1k7f9Hk

# Trends

**Task-Specific** ➤ **Foundation Models** ➤ **Multimodal** ➤ **Agents**

**Foundation Models**
- New domains
- Domain-specific adaptations
- New datasets
- New benchmarks

**Multimodal**
- Open weights and open data
- Grounding
- VLMs - language is the "glue"
- Multi-stage pretraining
- Larger inputs
- Benchmarks

**Agents**
- Multi-agent systems
- Specialization
- Orchestration
- Mimic how human experts work

# Resources

Recordings of keynotes and workshops:

https://cvpr.thecvf.com/Conferences/2025/Videos

Voxel51's "Best of CVPR" Series:

https://voxel51.com/events



Visual AI in Healthcare

June 27, 12-2 pm EDT

My talk: "Leveraging Foundation Models for Pathology: Progress and Pitfalls"

https://voxel51.com/events/visual-ai-in-healthcare-june-27-2025

# Take the Next Step

**Pixel Clarity Call**

A free 30-minute call where we'll dive into your unique challenges and goals—whether you're seeking sharper models, deeper insights, or a new direction for your AI projects.

- Gain expert perspective on your current approach
- Discover high-leverage opportunities tailored to your mission

Book now: https://calendly.com/hdcouture/pixel-clarity-call

# Q&A

Let's flip the conversation:

1) If you attended CVPR or another recent conference, what is one thing you learned?

2) What imaging domain do you work with?

Which "phase" is your work in?

| Task-Specific | Foundation Models | Multimodal | Agents |

Is the current "phase" sufficient to solve a real world problem?

   OR

What is blocking the path to the next "phase"?