

Foundation Models for Earth Observation

EO for Sustainable Development Applications
21.01.2025, Rohan Sawahn

Lecture Overview

Part 1 (50min)

1. Introduction to Foundation Models
2. Foundations Models in EO
3. Understanding Embeddings
4. Training Paradigms for Foundation Models for EO
5. Using EO Foundation Models in Practice

Part 2 (30min)

Hands-On Session: Applying Foundation Models for Croptype classification & Deforestation Detection

Part 3 (20min)

Continue to Work on your Presentations

What is a Foundation Model?

What is a Foundation Model?

*“Foundation Models serve as a **foundation** that has some **general understanding** of the domain and can be applied to **broad range of specialized tasks** without requiring to retrain the model from scratch.”*

What is a Foundation Model?

*“Foundation Models serve as a **foundation** that has some **general understanding** of the domain and can be applied to **broad range of specialized tasks** without requiring to retrain the model from scratch.” ~ a foundation model*

Foundation Models Characteristics

Large-Scale Pretraining

Pretrained on a vast corpus of (unlabelled) data. Typically millions of samples requiring 100s of GPU hours.

Unified Representation

Shared, general purpose representation feature space of the domain (ie. earth) that can be reused for different tasks.

Task Agnostic

Not designed for one specific task, but can easily be adopted to a variety of downstream tasks by the end users.

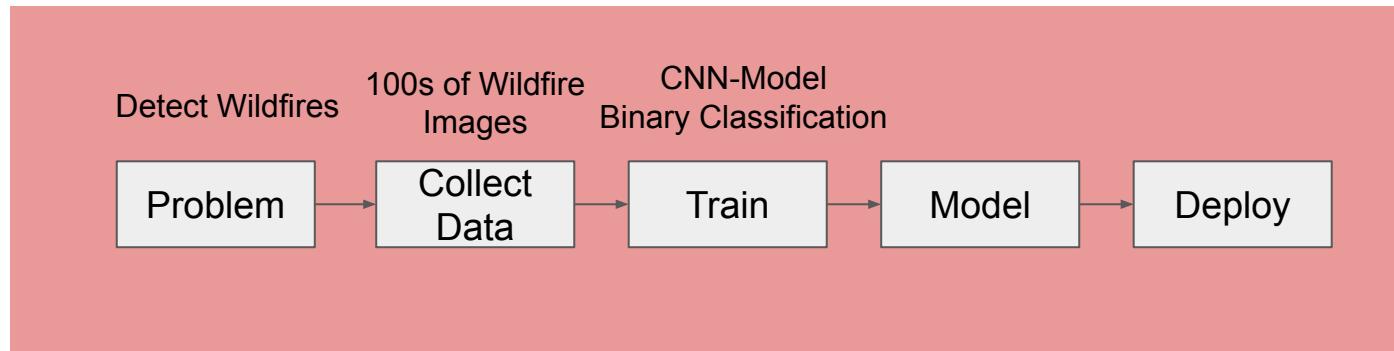
Transfer Learning

Can be fine-tuned for specific downstream tasks using limited data, thus enabling domain adoption.

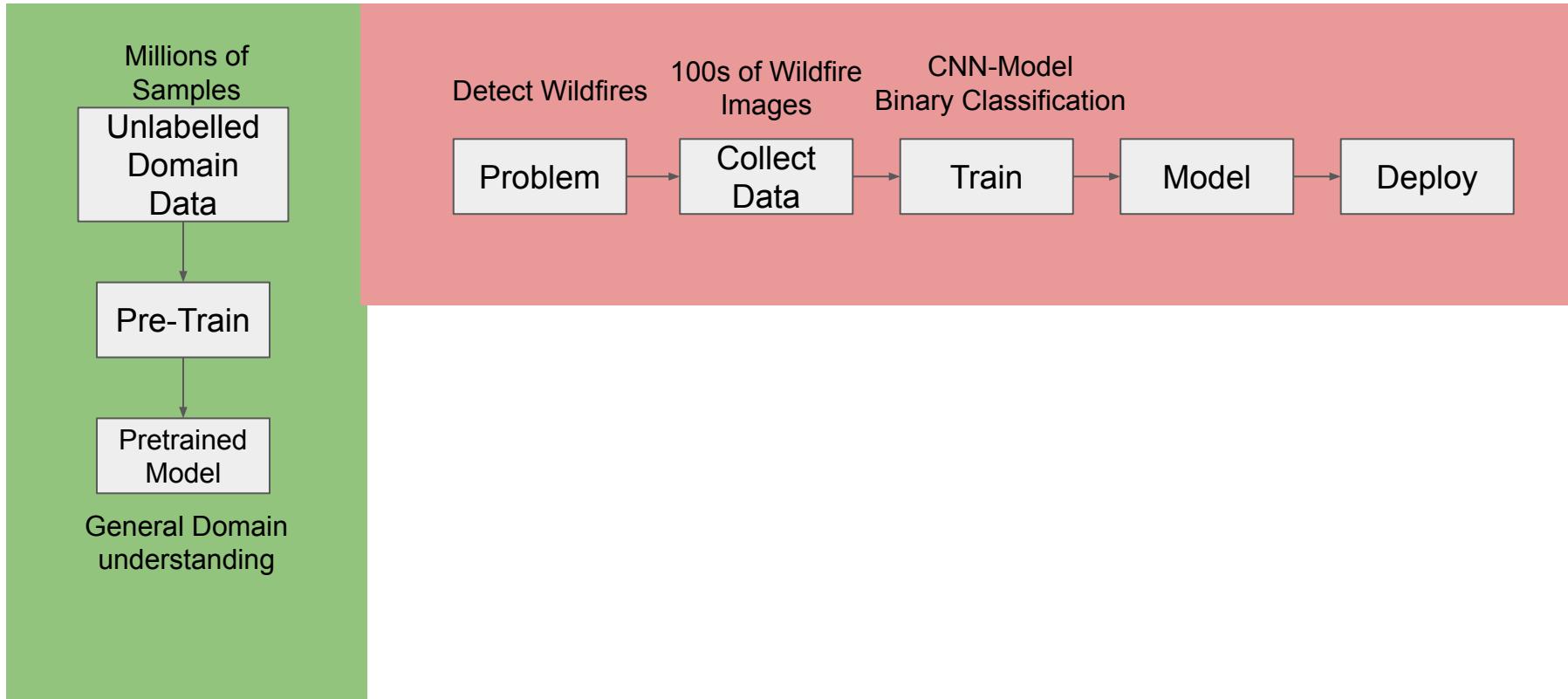
Zero-Shot Capabilities

Can be used without labelled instances with unsupervised methods like clustering etc.

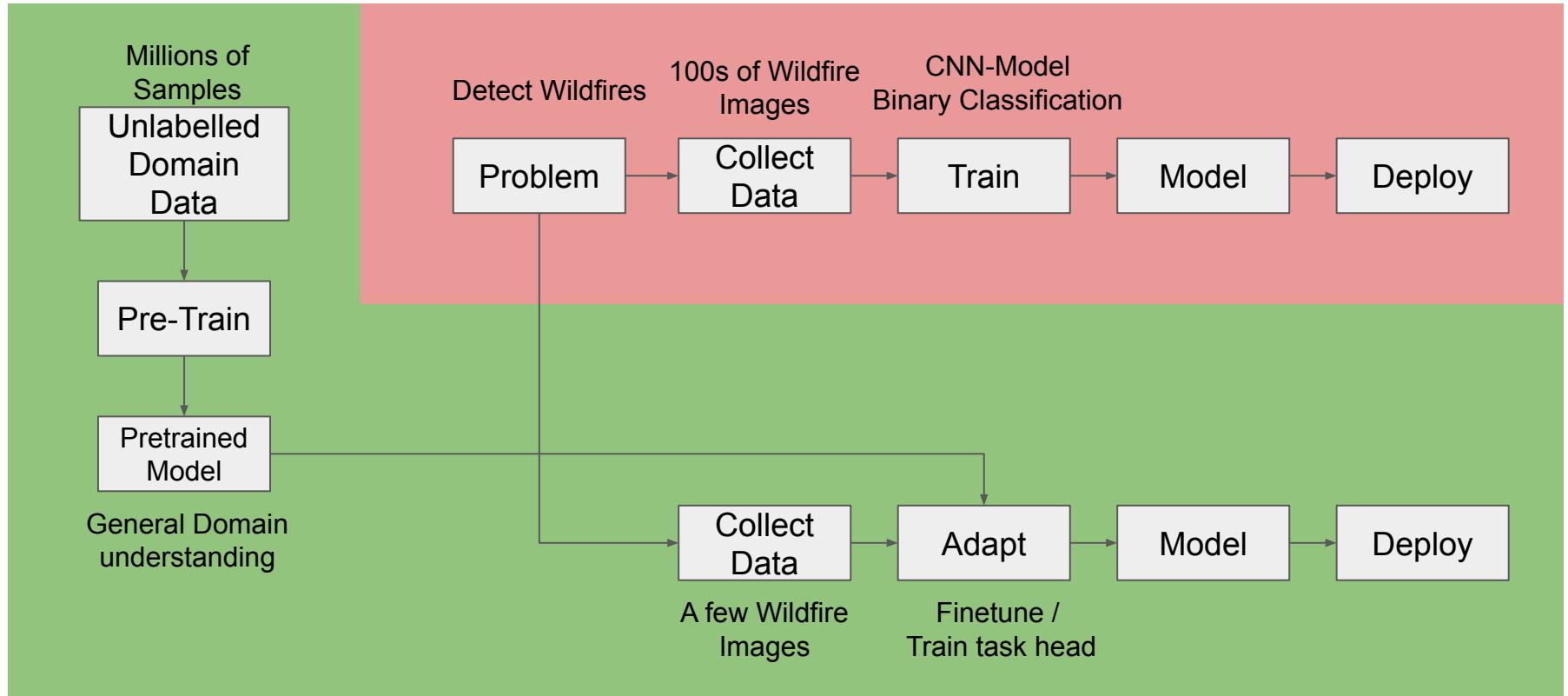
Traditional ML vs Foundation Models



Traditional ML vs Foundation Models



Traditional ML vs Foundation Models



Why Foundation Models in EO?

Data Availability

- Petabytes of unlabelled data
- Collection task-specific data is expensive & requires expertise

Common Visual Patterns

- Similar visual patterns for different processes - exploit shared understanding
- Similar environments often exhibit similar patterns

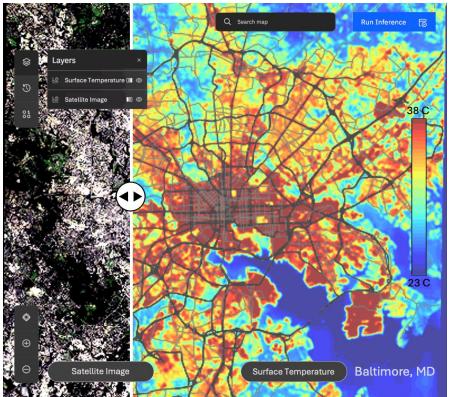
Multimodal Nature of Data

- Optical, SAR, Hyperspectral...
- Climate, Location, Elevation, Descriptions
- Variety of spatial and temporal resolution

Democratization of EO & ML

- Reduces requirements for Compute & Data
- Faster development of solutions for e.g disaster response

Foundation Models in EO Use Cases



<https://research.ibm.com/blog/prithvi2-geospatial>



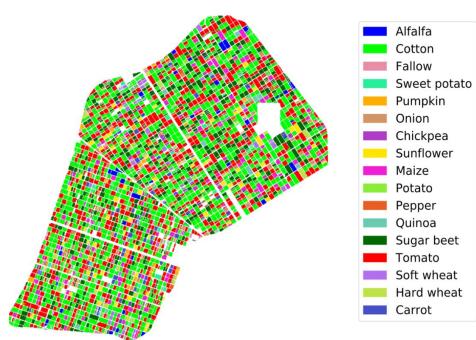
<https://www.nytimes.com/2024/09/15/world/africa/floods-africa.html>



<https://earth.org/amazon-rainforest-deforestation-facts/>



<https://www.bbc.com/news/science-environment-64920236>

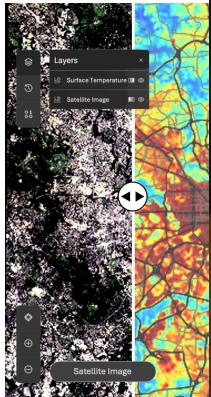


https://www.esa.int/ESA_Multimedia/Images/2021/02/Crop-type_map



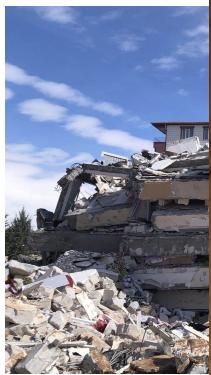
<https://airqoon.com/resources/articles/urban-air-pollution-sources-and-pollutants/>

Foundation Models in EO Use Cases



<https://research.ibm.com>

ONE MODEL TO RULE THEM ALL



<https://www.bbc.com/imgflip.com>



<https://forestation-facts/>



<https://urban-air-pollution-sources.com>

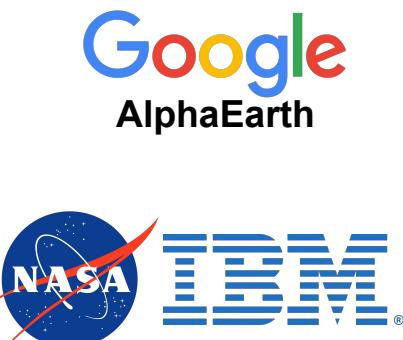
Overview of Foundation Models in EO

Model Name	Architecture	Pre-training Dataset	Resolution (m)	Geographic Coverage	Image Analysis Levels	Pretrain methods	# of Params
CMC-SSSR [14]	ResNet-50	NWV, DOAS [1-1], BigBarbara [1%], BigBarbara [1%]	0.2 to 60	Global	Image-level, Multires Coding	Contrastive, Mutual Coding	23M
SeCo [6]	ResNet-50	Sentinel-2 imagery	10, 20, 60	20K Locations Worldwide	Image-level, Spectral-temporal	CL	23.9M
GeoKE [15]	ResNet-50	Levi-CR [16]	0.8 to 16	-	Image-level, Pixel-level, Region-level, Feature-level, SSL	Geographical Knowledge Supervision	23.9M/23.8M
MATTER [17]	ResNet-54	Sentinel-2 Imagery	-	Read and Reuse Regions with Little Changes	Pixel-level, Region-level	SSL	21.3M
GASSL [18]	ResNet-50	IM2W [1], GoogleSat [19]	-	7 Continents	Image-level, Region-level	CL	23.9M
RSP [20]	ViT-H2-S2	MillenniAD [20, 21]	0.5 to 155	Global	Image-level, Region-level, Spectral-temporal	Supervised Learning	24.8M/23.5M/29M
DINO-MC [21]	ViT-S8	BigBarbara-SM [20]	10	Global	Image-level, Region-level, Pixel-level	SSL	23M
Schubert et al. [22]	Swin Transformer	S2N2M2O [22]	10	-	Image-level, Region-level, Pixel-level	CL	-
Ragle et al. [23]	ViT-based Transformer	2 million RS images	0.3 to 30	6 Continents	Image-level, Region-level, Spectral-temporal	MIM	-
GAoT [24]	ResNet-50	Levi-CR [16]	0.8 to 16	-	Image-level, Region-level, Spectral-temporal	SSL	23.9M
RS-BYOL [25]	BYOL	Sent2MSD [25]	10 to 20	Global	Image-level, Region-level, Pixel-level	SSL	23.9M
CSPF [26]	VIT-B	ImageNet-1K [26]	-	-	Image-level, Region-level	SSL	88M
IVSA [10]	VIT	MillenniAD [20, 21]	0.5 to 155	Global	Image-level, Pixel-level, Region-level	MAE	100M
SiMaF [27]	MLP-Mixer	(IM2W/Sentinel-2) [21]	10, 20, 60	-	Image-level, Region-level, Pixel-level	MAE	307M
TOW [28]	TOV	TOV [28]	-	Global	Image-level, Region-level	SSL	-
CM3D [29]	Teacher-student Model	MillenniAD [20, 21]	Varied	-	Image-level, Pixel-level, Region-level, Spectral-temporal	SSL	25.6M/87.8M
CACo [30]	ViT-B	Sentinel-2 Imagery	10	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	SSL	11.7M/23.5M
Id-SwinLR [27]	ResNet-18	SEN1MS [27]	-	Global	Image-level, Pixel-level, Region-level	CL	11.7M

Model Name	Architecture	Pre-training Dataset	Resolution (m)	Geographic Coverage	Image Analysis Levels	Pretrain methods	# of Params
GPM [31]	Teacher-Student	GeoPhi [31]	-	Global	Image-level, Pixel-level	Contrastive Pretraining	-
SatellitePro [32]	Multi-Branch	SatelliteSet [32]	1, 10	Global	Image-level, Pixel-level	Multi-task Learning	8M
Region-Sense [17]	Multi-Branch	RS Spatiotemporal Dataset	-	Global	Pixel-level	SSL	-
Scale-MAE [33]	VIT-Large	IM2W [1]	10, 20, 60	Global	Image-level, Pixel-level, Region-level	MAE	322.9M
RegMo-Me [34]	CNN-Transformer	AID [34]	0.5 to 30	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	FB-MIM	60%
DeCVR [35]	Multi-modal SSL	SEN1MS [27]	Varied	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	SSL	23.5M
Feng et al. [36]	MSFE-MMHF	Multi-modal Dataset	Varied	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	SSL	-
PG-MAE [37]	VIT	SSLR-GD-512 [37]	10	Global	Image-level, Pixel-level	MAE	-
Pritvi [38]	VIT	Harmosted Landset Sentinel 2	30	Contiguous U.S.	Image-level, Pixel-level	MAE	100M
CRISP [39]	Multi-modal Decoder	SSL4EO [39]	10	Human Settlements	Image-level, Pixel-level	CL, MAE	8M
Uso-1 [40]	VIT	RS [39]	Varied	Global	Image-level, Pixel-level	MAE	8M
Cross-Scale-MAE [30]	VIT-B	DOW [1]	-	Global	Image-level, Pixel-level	MAE	8M
U-BARN [30]	Unet Transformer	Sentinel-2 Imagery	Varied	France	Image-level, Pixel-level	SSL	-
EarthPT [32]	Transformer	Sentinel-2 Imagery	10	UK	Image-level, Pixel-level	Autoregressive SSL	700M
GrASP [41]	Teacher-Student Network	ImageNet [1], MillenniAD [20, 21]	0.5 to 155	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	SSL, SL	-
Sat2MSD [42]	SeNet2MSD [42]	-	Varied	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	SSL	11.7M
SMFLR [43]	Generative ConvNet	GeoPhi [31]	0.05 to 150	Multiple Continents	Image-level, Pixel-level, Region-level, Spectral-temporal	SSL	88M/97M
SpectroGPT [44]	3D GPT	Sentinel-2 Imagery	Varied	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	MAE	100M/300M/600M
Promo [45]	MAS-based Framework	Promo-21-IM [45]	10	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	MAE	40K
SatM4x [46]	SatM4x	DOw [1]	Varied	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	Multi-scale pre-training	-

Model Name	Architecture	Pre-training Dataset	Resolution (m)	Geographic Coverage	Image Analysis Levels	Pretrain methods	# of Params
Folio-Bench [47]	VIT	Multiple	Varied	Global	Pixel-level, Region-level, Image-level, Spectral-temporal	MAE	10M/100M
SkySense [48]	Factored Multi-Modal Spatio-temporal Encoder	Multiple	Varied	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	CL	2.0B
LIN [49]	Multi-Modality	GeoPhi [31]	-	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	SSL	0.65M
dat4D [50]	dat4D	GeoPhi [31]	0.1 to 155	Global	Pixel-level, Region-level, Image-level, Spectral-temporal	SSL	10M
DINO-MC [21]	DINO	Set-1000 [51]	10 to 60	Image-level, Pixel-level, Region-level, Spectral-temporal	SSL	-	
GfA-Net [52]	OFANet	Multi-modal Dataset	Varied	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	MIM	~10M
MTP [53]	Task-Specific Decoders	SAMRS [49]	Varied	Global	Region-level, Spectral-temporal	Pretaining	over 300M
BFM [11]	VIT	MillenniAD [20, 21]	0.5 to 155	Global	Image-level, Pixel-level, Region-level	MAE	1.9M/0.42B
MMEarth [54]	MP-MAE	Multi-modal Geospatial Data	-	Global	Image-level, Pixel-level	MP-MAE	3.7M to 650M
CiSiM [55]	VIT	WorldView-3 Imagery	Varied	Asia	Image-level, Pixel-level, Region-level	MIM	88M
SAROT-X [56]	HVIT	SAR Datasets	0.1 to 3	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	MIM	-
SoftCon [57]	SVIT	SSL4EO-S12-M1 [37]	-	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	Cross-Attention	23M, 23M, 300M
LaMeViT [58]	Hierarchical VIT	MillenniAD [20, 21]	-	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	Dual Cross-Attention with Learnable Meta Token	33M
SIMAR [59]	3D Transformer-based	IM2W/Sentinel-2 [59]	-	Global	Image-level, Spectral-temporal	MIM	30M
RS-DIM [49]	Multi-resolution Inference Framework	AI-Co-MultiTask [59]	-	-	3D Region-level, Pixel-level	Generalized Feature Matching with Relative Depth	over 1B
A2-MAE [57]	VIT-Large	STND [59]	0.8 - 30s	Global	Image-level, Pixel-level, Region-level, Spectral-temporal	Matching-PAT and Geographic Encoding Module	-
HyperSIGMA [60]	VIT-based	Hypersig-SIGK [59]	30m	Global	Image-level, Region-level, Spectral-temporal	MAE	over 1B

Lu et al., 2025, Vision Foundation Models in Remote Sensing: A survey



UNIVERSITY OF
CAMBRIDGE
Tessera

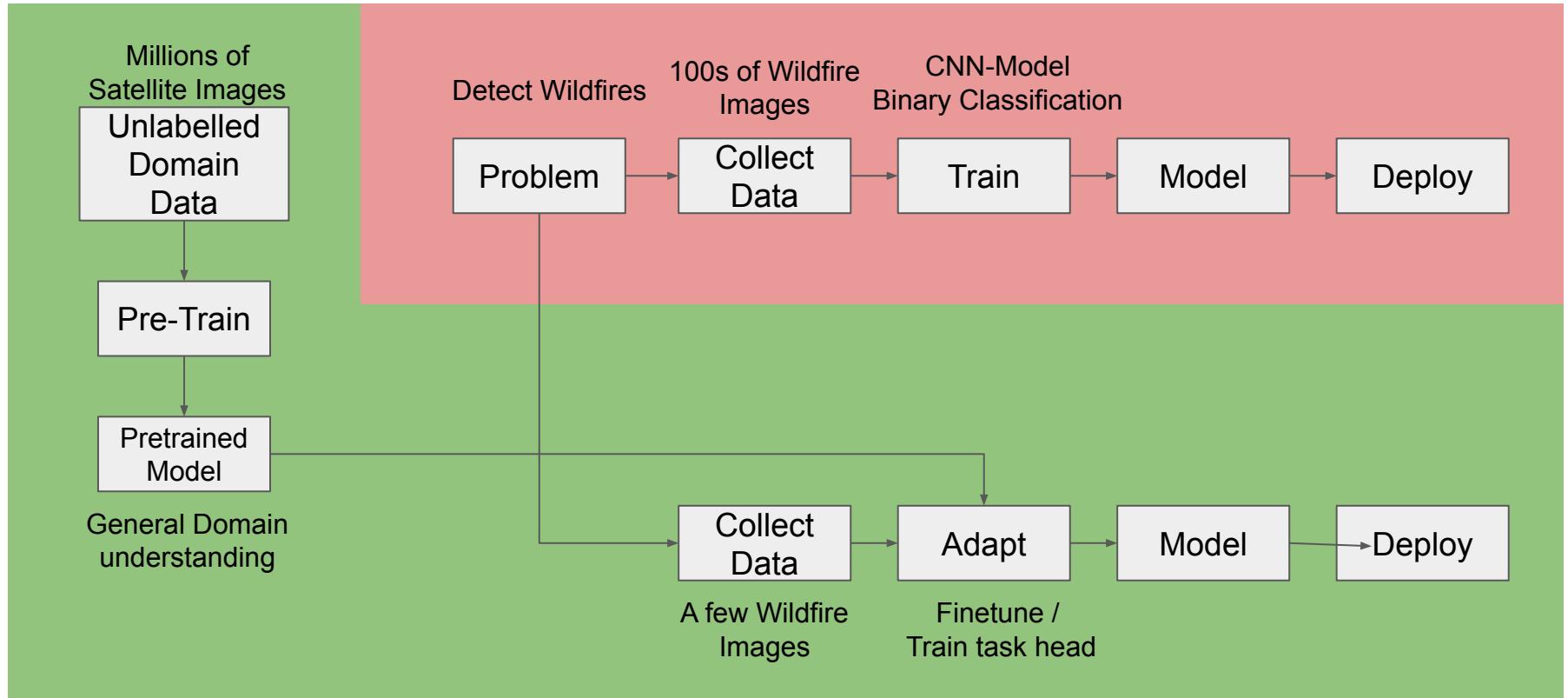

Presto/Galileo


Φ-lab

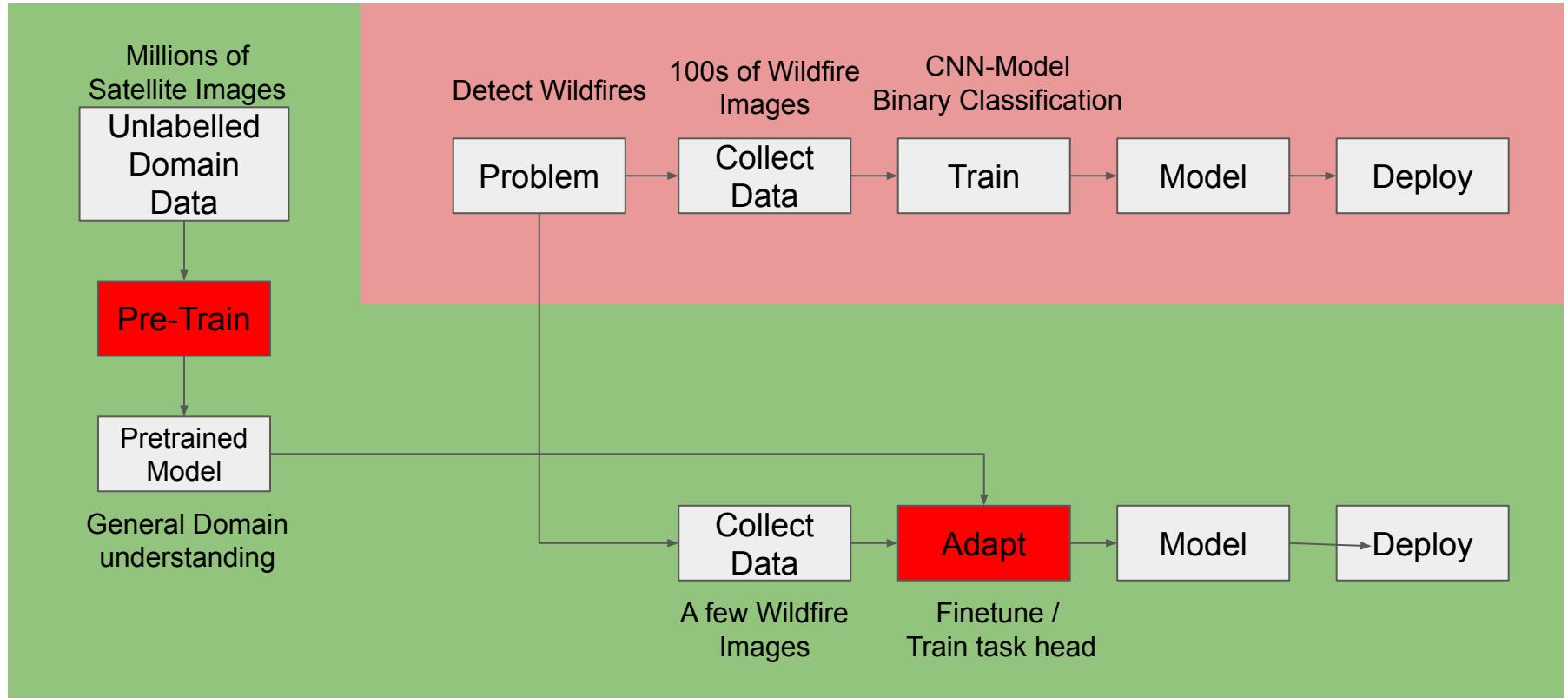
MajorTom


OlmoEarth

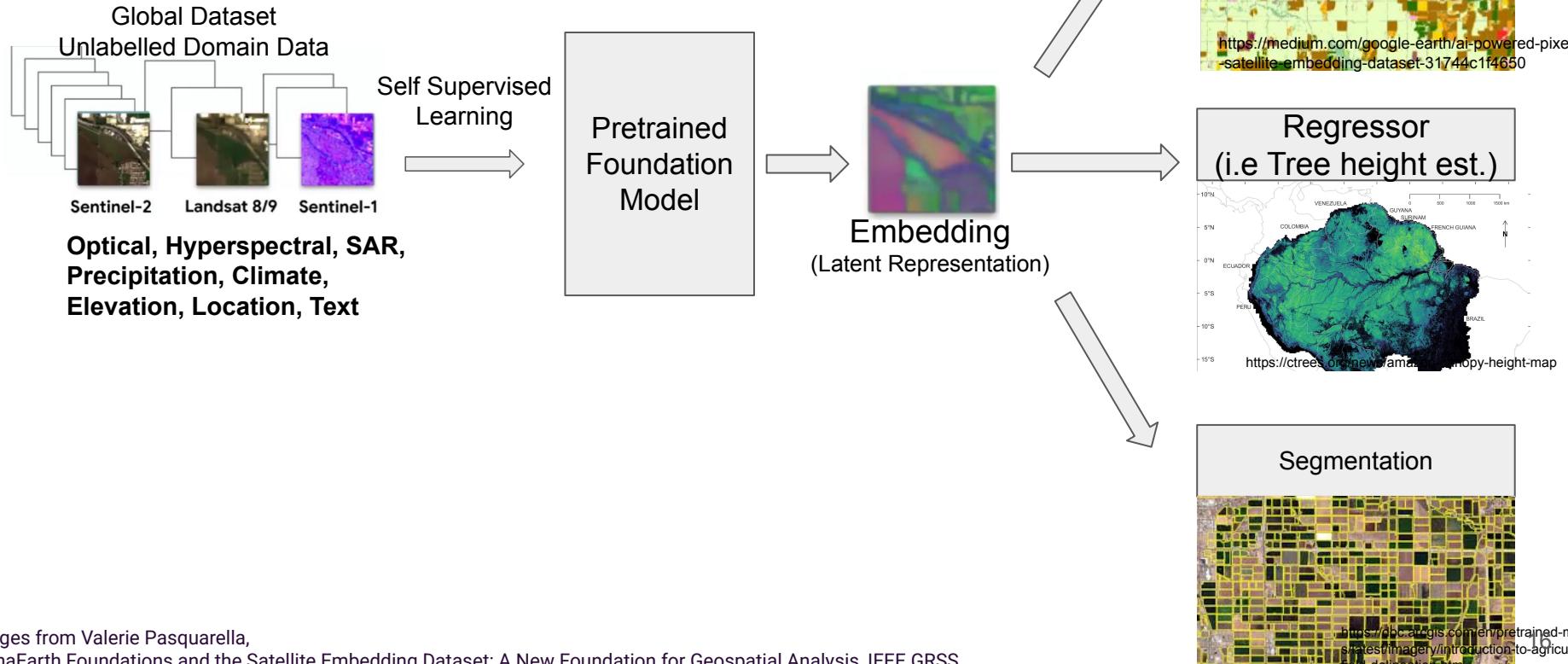
Traditional ML vs Foundation Models



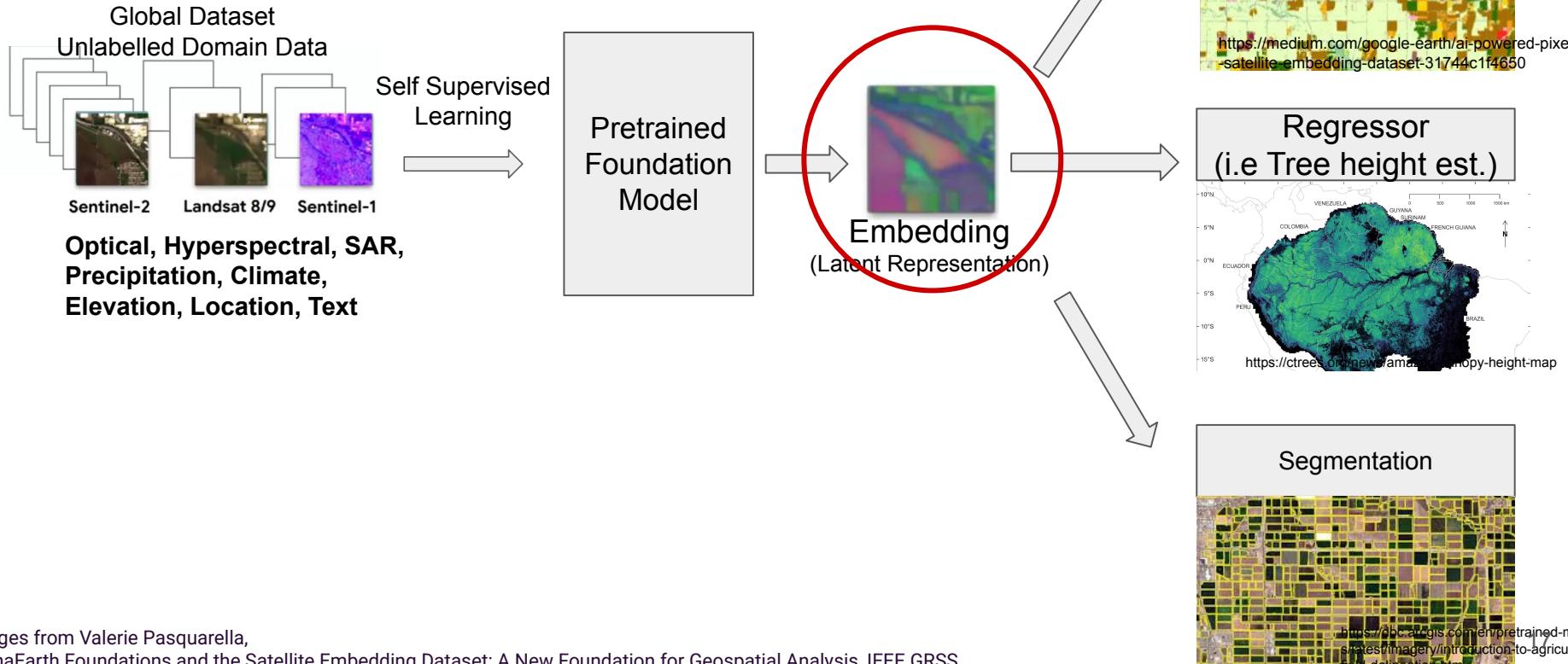
Traditional ML vs Foundation Models



Foundation Models in EO

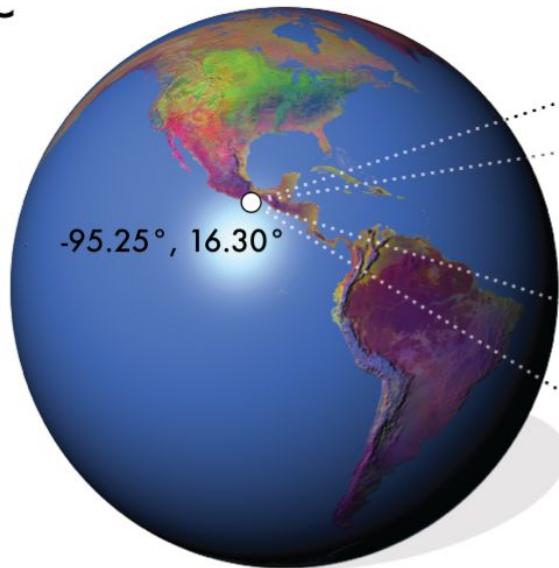


Foundation Models in EO

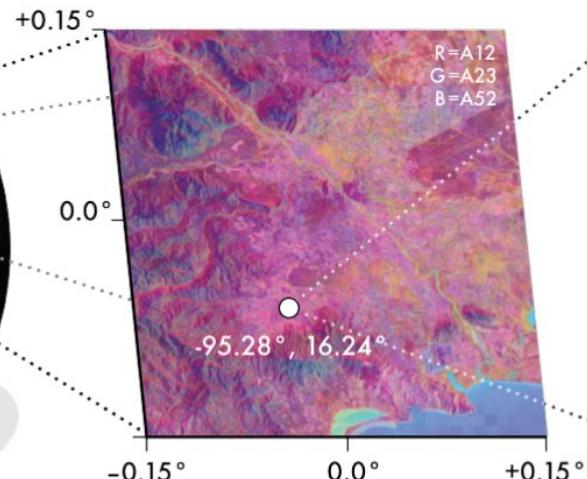


What's in an Embedding?

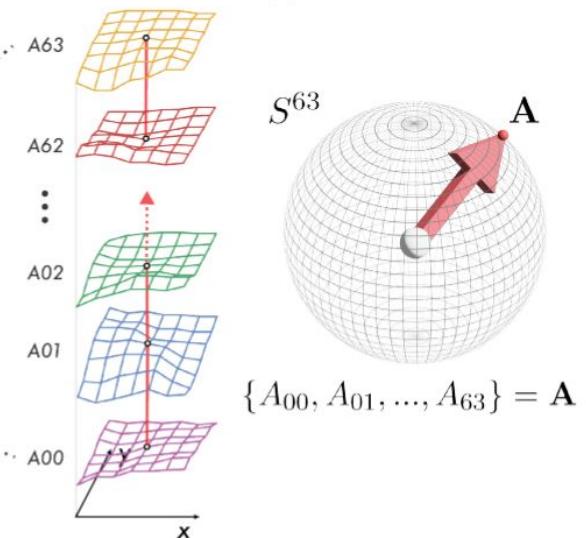
C



D



E



Brown et al, 2025, AlphaEarth Foundations

What's in an Embedding?



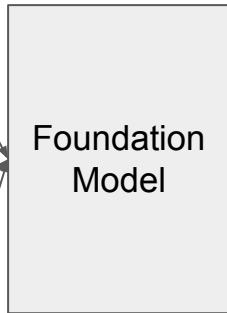
Sentinel-2



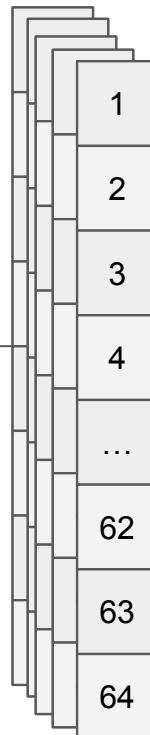
Sentinel-1



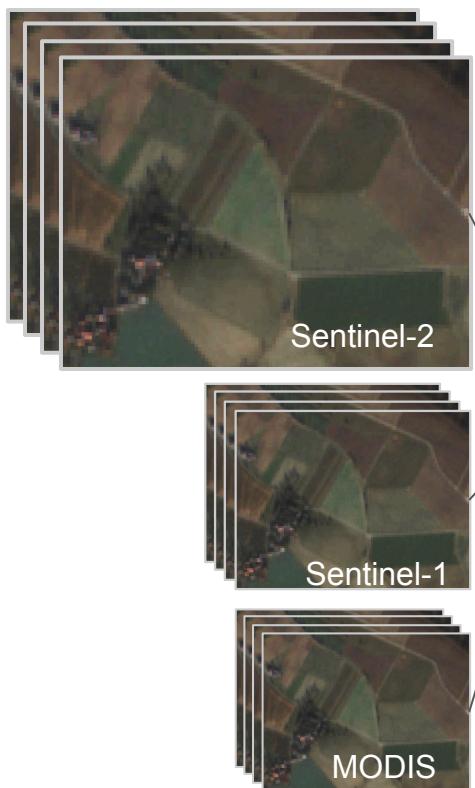
MODIS



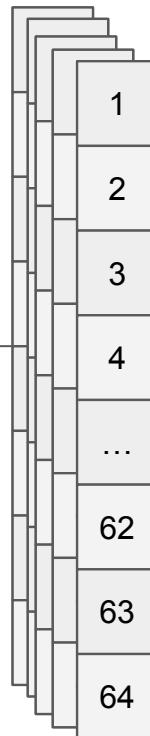
1 embedding vector
per pixel!



What's in an Embedding?

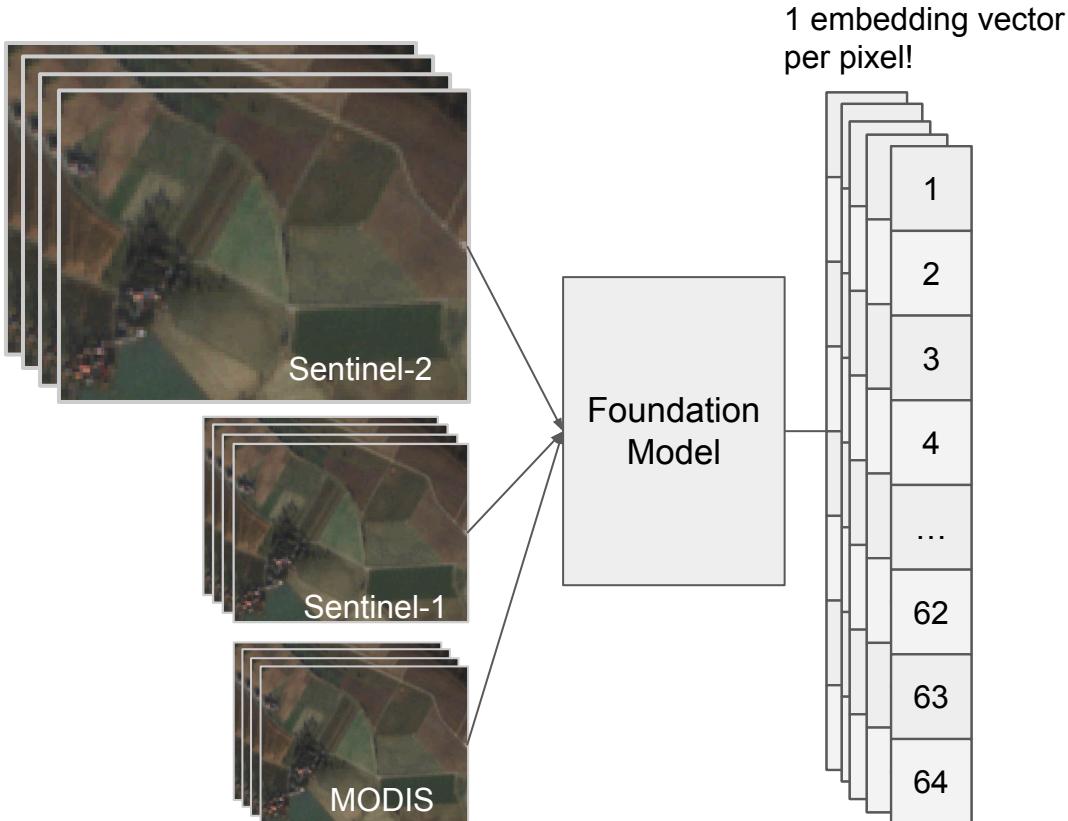


1 embedding vector
per pixel!



Latent Space Representation
Learned coordinate system where “nearby” points mean “similar in the way the model cares about.”

What's in an Embedding?

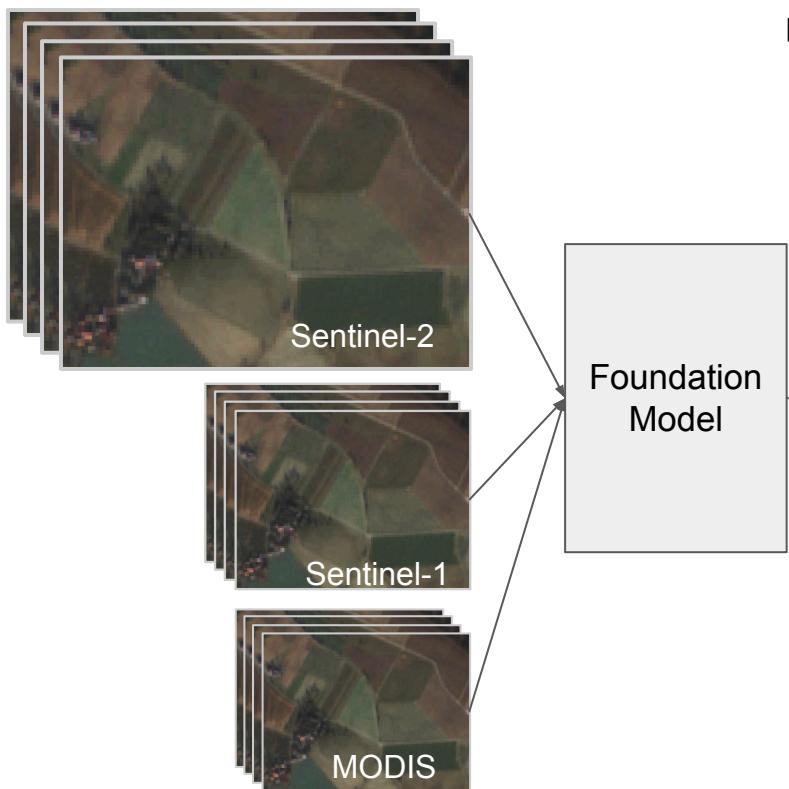


Latent Space Representation
Learned coordinate system where “nearby” points mean “similar in the way the model cares about.”

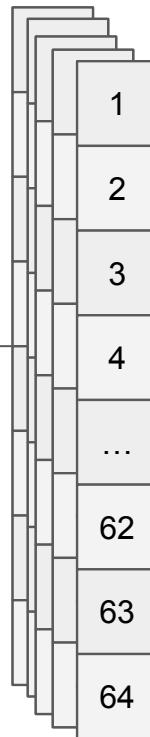
Characteristics of Embeddings

- Values are not directly observed, but are (hidden) features that the model learns
- “Compression” of different data modalities & dates per pixel
- Individual dimensions usually don’t have clean human labels
- Directions can encode attributes

What's in an Embedding?

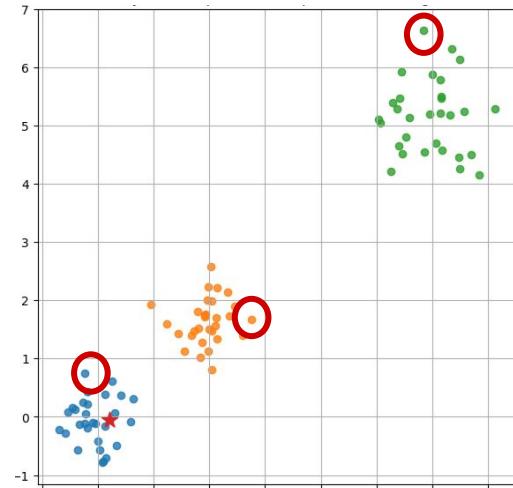


1 embedding vector
per pixel!



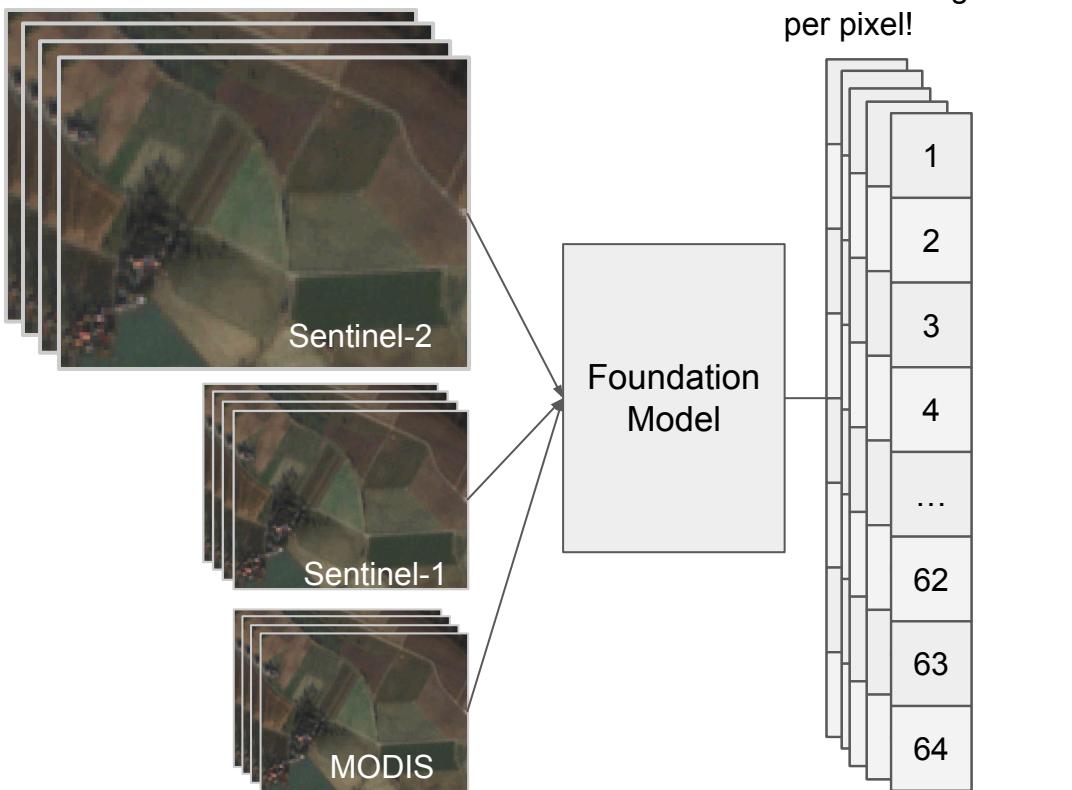
Latent Space Representation
Learned coordinate system where “nearby” points mean “similar in the way the model cares about.”

Similarity of Representations



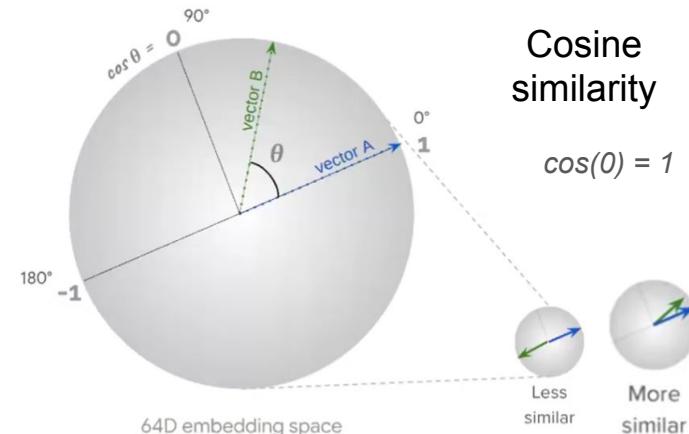
UMAP dimensionality reduction of embeddings.
Star = Wheat field

What's in an Embedding?



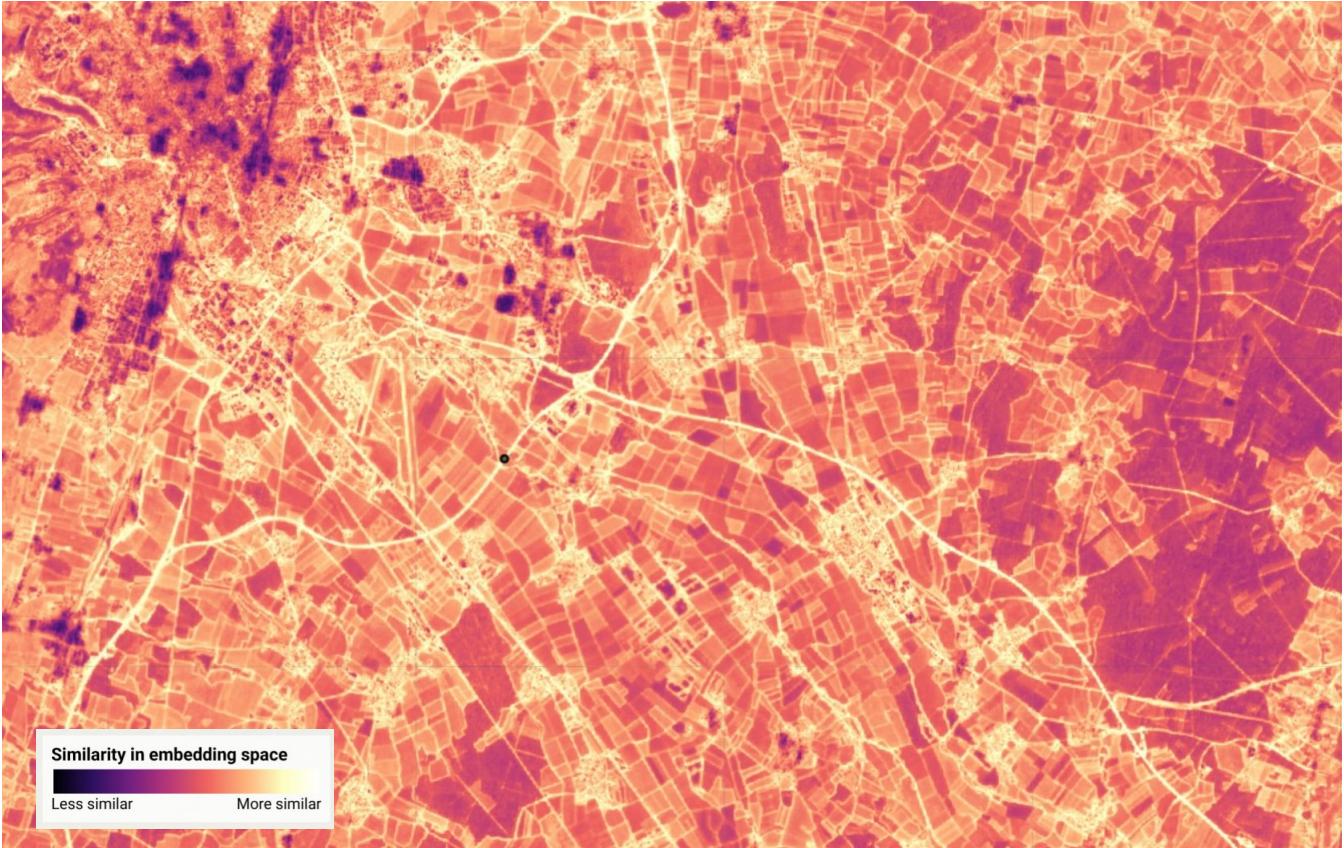
Latent Space Representation
Learned coordinate system where “nearby” points mean “similar in the way the model cares about.”

Similarity of Representations



Cosine similarity
 $\cos(\theta) = 1$
Valerie Pasquarella,
AlphaEarth Foundations and the Satellite Embedding Dataset: A New Foundation for Geospatial Analysis, IEEE GRSS

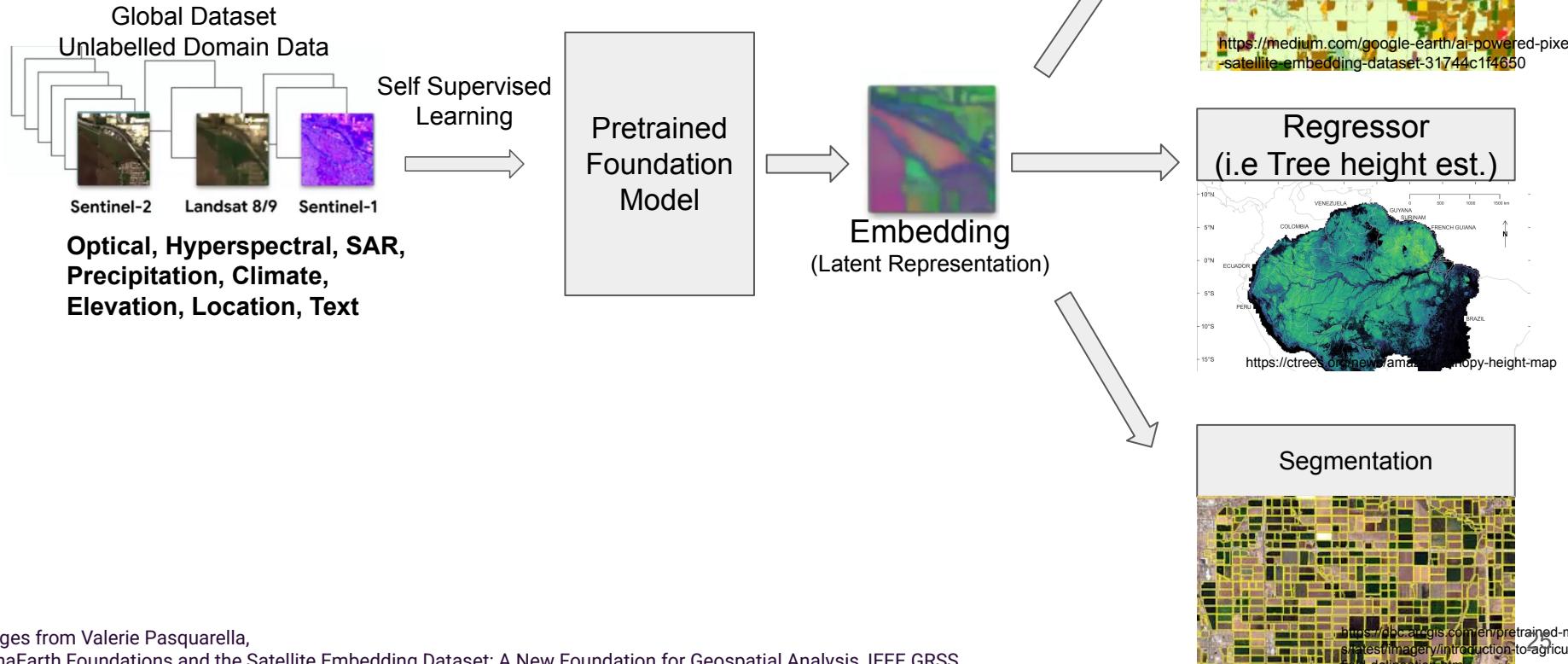
What's in an Embedding?



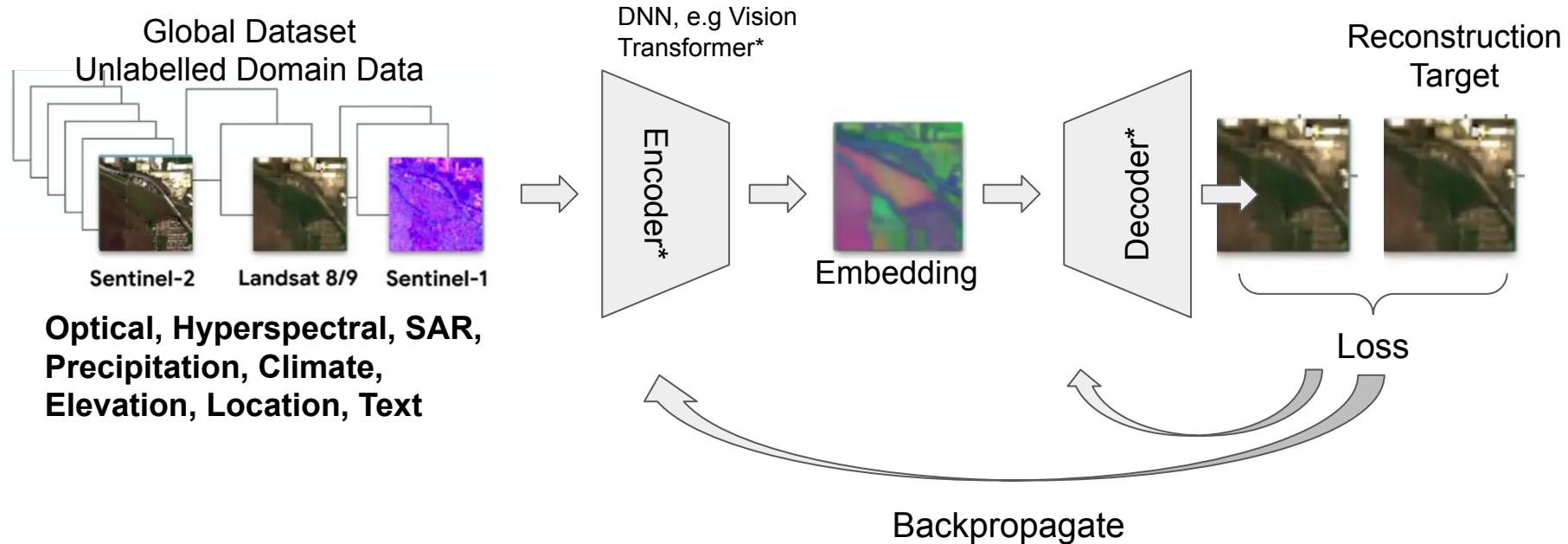
Try it out yourself:

<https://earthengine-ai.projects.earthengine.app/view/embedding-similarity-search#year=2024;zoom=18;lon=5.122757904832156;lat=47.25973734720962;clicked=true;>

Foundation Models in EO



Foundation Models in EO: Training



*Simplified schematics. Actual model architecture varies significantly by model.

Self Supervised Learning

*“Self-supervised learning is a way to train a model on **unlabeled data** by creating its own **training signal** from the data itself.”*

1. Masked Reconstruction
2. Multi-Modal Alignment
3. Contrastive Learning
4.

Self Supervised Learning

Masked Reconstruction



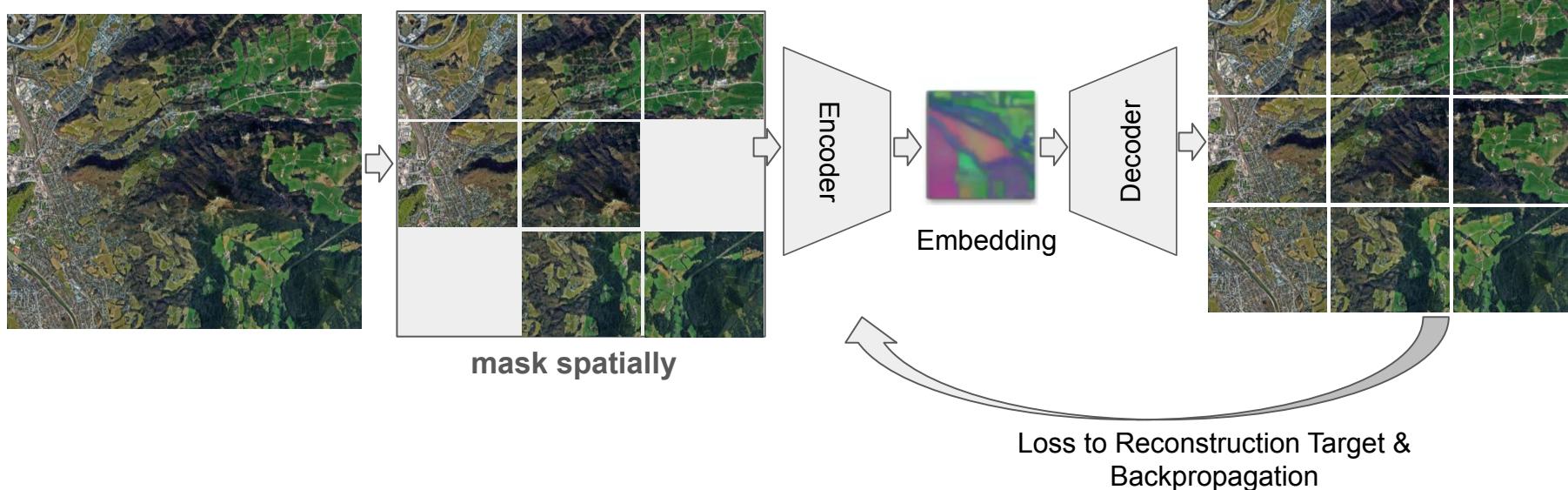
Self Supervised Learning

Masked Reconstruction



Self Supervised Learning

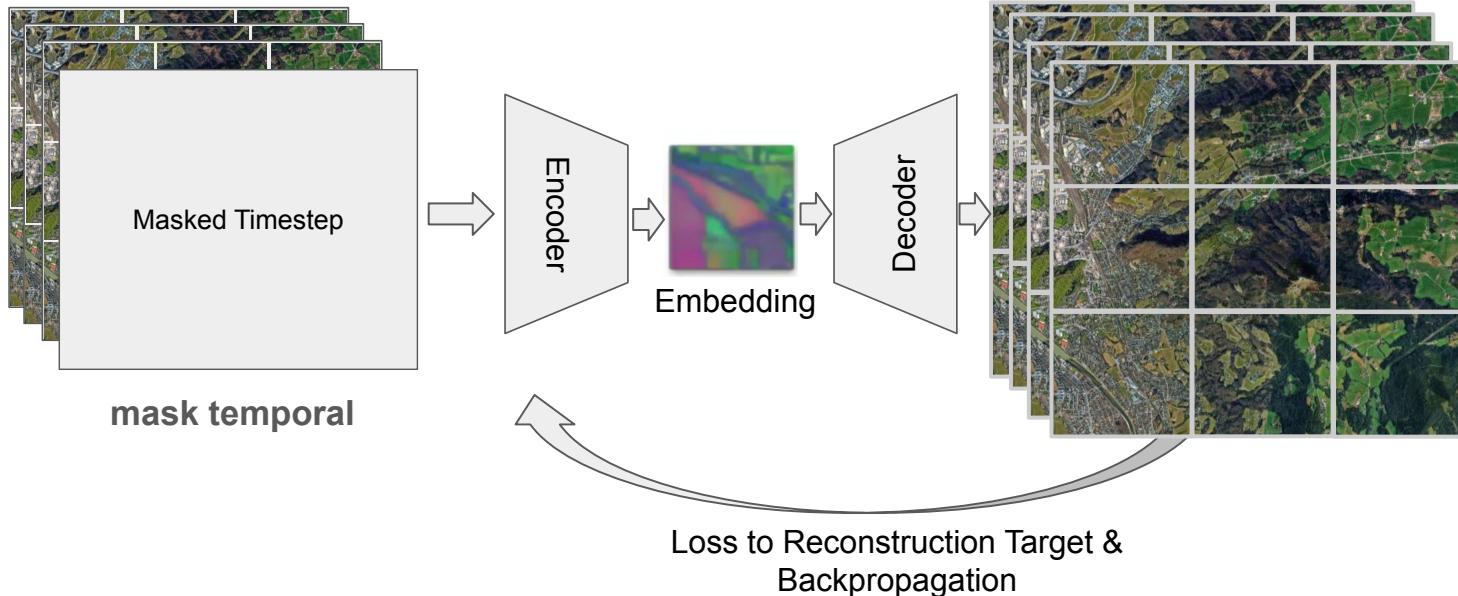
Masked Reconstruction



=> Model learns to reconstruct masked data on the spatial axis.

Self Supervised Learning

Masked Reconstruction



- => Model learns to reconstruct missing data on the temporal axis.
- => Similar process for the spectral bands.

Self Supervised Learning

Multi-Modal Alignment

Data from different sensors available: Multispectral, Hyperspectral, SAR, different resolutions, ...

Goal: Make Representations of the same locations with temporal alignment have similar embeddings.
Important: Need to ensure temporal alignment of different sensor modalities!

Different Methods

- Also use Masked Autoencoding for different modalities!
- Separate encoders (+decoders) per modality + (learnt) fusion into latent space embedding
- Contrastive-Style Losses
-

Self Supervised Learning

Pre-Training for EO Foundation Models differs from model to model!

Some great further reads with details:

- Szwarcman, Daniela, et al. "Prithvi-eo-2.0: A versatile multi-temporal foundation model for earth observation applications." IEEE Transactions on Geoscience and Remote Sensing (2025).
- Brown, Christopher F., et al. "Alphaearth foundations: An embedding field model for accurate and efficient global mapping from sparse label data." arXiv preprint arXiv:2507.22291 (2025).
- Tseng, Gabriel, et al. "Galileo: Learning Global & Local Features of Many Remote Sensing Modalities." Forty-second International Conference on Machine Learning.

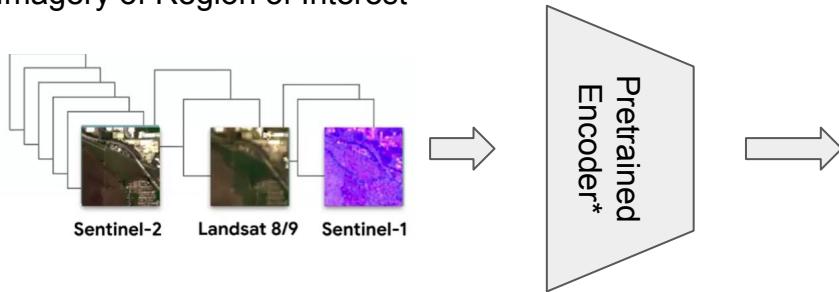
Foundation Models in Practice

Goal: Learn how to apply a pre-trained foundation model for an earth observation task for sustainable development.

Applying EO Foundation Models

Unsupervised - Zero Shot (No Labels required)

Imagery of Region of Interest

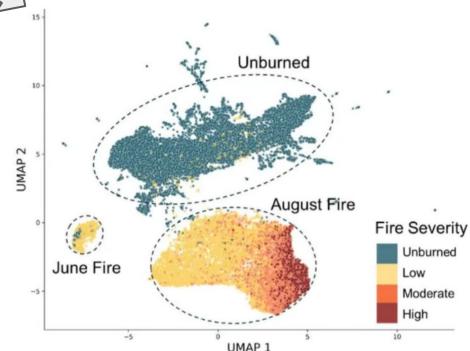


Embeddings for
Region of Interest

Similarity Search



Dimensionality Reduction & Clustering



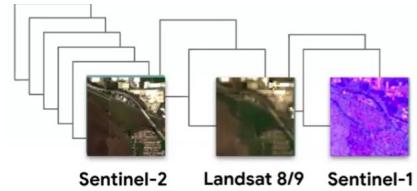
*Decoder dropped! Most Geo-FMs do this. However, some models might have an architecture that uses both or a different architecture altogether..

Applying EO Foundation Models: Supervised

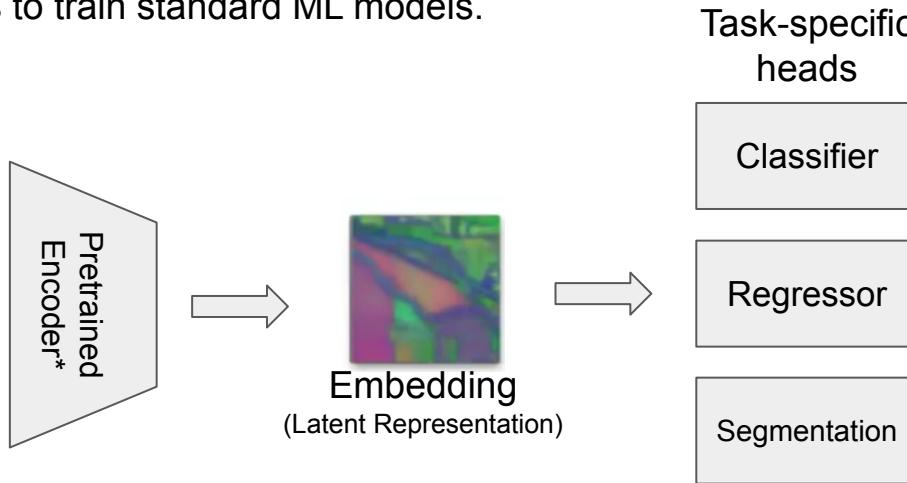
Supervised - Few Shot

Embeddings are used as features to train standard ML models.

Few samples of locations that are relevant for the task.



+ Labels



Option A: End-to-End finetuning (Train the task-head + finetune the FM)

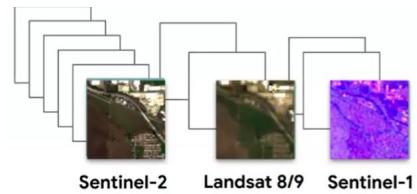
Option B: Only train the task-head and freeze the foundation model (backbone).

Applying EO Foundation Models: Supervised

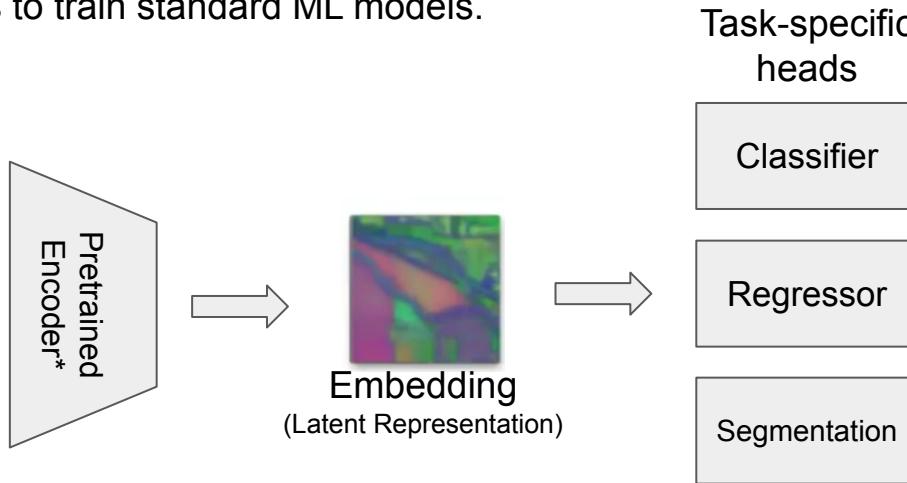
Supervised - Few Shot

Embeddings are used as features to train standard ML models.

Few samples of locations that are relevant for the task.



+ Labels

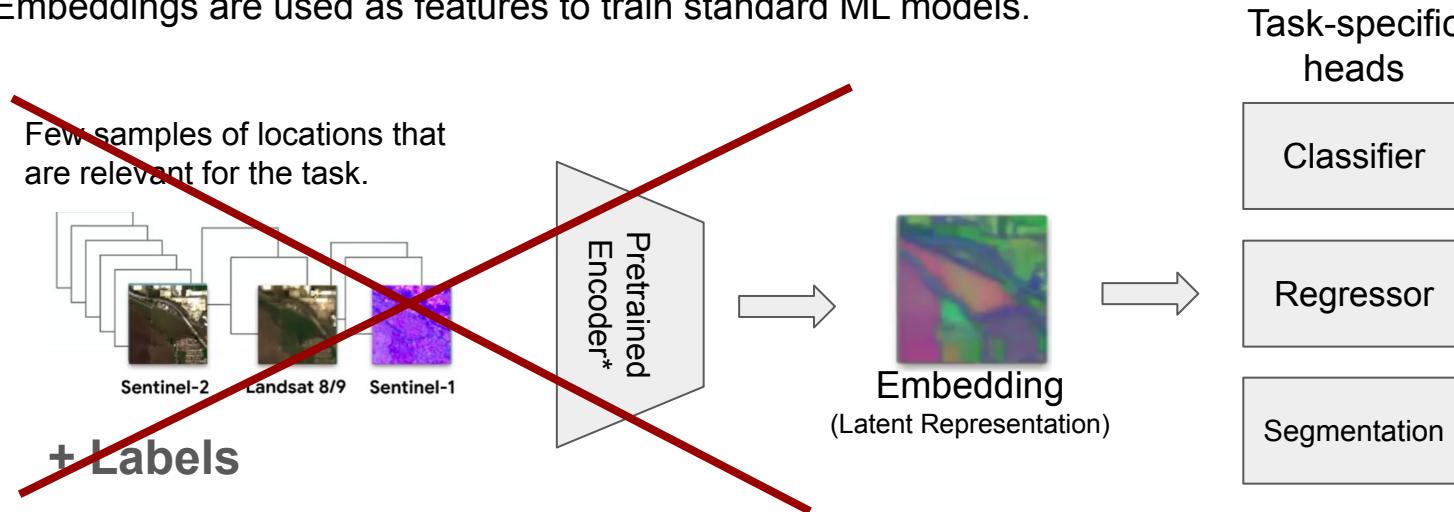


Problem: Still requires access to all the satellite imagery + need all the different modalities to infer the embeddings. Imagery can be multiple TBs of data + Compute on GPUs is expensive!

Applying EO Foundation Models: Supervised

Supervised - Few Shot

Embeddings are used as features to train standard ML models.



Problem: Still requires access to all the satellite imagery + need all the different modalities to infer the embeddings. Imagery can be multiple TBs of data + Compute on GPUs is expensive!

=> Available Stores (Datasets) with free access to **precomputed** embeddings. Democratizes access to EO data. E.g Google Earth Engine, Source Cooperative

Hands-On Session Instructions

Goal: Learn how to use pre-computed embeddings for your work!

Train your own Model for

- 1: Crop Type Mapping: <https://shorturl.at/wEY4V>
- 2: Advanced / Exercise: Deforestation Detection <https://shorturl.at/auXW3>

Step 1: Open the link and sign into Colab

Step 2: Run the code cell by cell and try to understand what is happening

Step 3: Once you have the model working, try to think of ways how to improve it and play around with the parameters. What changes when you change the number of training samples?

Pitfalls & Limitations of GeoFMs

Fast Changing Space with many Unanswered Questions:

- “Black-Box Models”
- High price- and energy costs for pre-training
- Timely availability of pre-computed embeddings
- Seasonality and temporal resolution in inputs
- Hard to quantify uncertainty
- Standard RS+ML Methods often achieve competitive performance
 - Limited understanding of when they work well

Crop	Scheme	AEF				RS			
		R ²	RMSE	R ²	RMSE	R ²	RMSE	R ²	RMSE
Corn	County-Year	0.43	2.28	0.44	2.27	0.45	2.24	0.46	2.24
	Yearly	0.31	2.54	0.30	2.55	0.33	2.50	0.33	2.49
	Scale-Transfer	0.24	2.67	0.25	2.64	0.37	2.43	0.38	2.40
Soybean	County-Year	0.36	0.87	0.37	0.87	0.35	0.90	0.35	0.90
	Yearly	0.29	0.92	0.31	0.91	0.28	0.93	0.25	0.95
	Scale-Transfer	0.13	1.02	0.11	1.03	0.22	0.97	0.24	0.95
Winter Wheat	County-Year	0.43	1.32	0.43	1.31	0.42	1.30	0.41	1.31
	Yearly	0.38	1.34	0.40	1.32	0.36	1.35	0.36	1.36
	Scale-Transfer	-0.18	1.85	-0.04	1.73	0.20	1.52	0.22	1.50

Task	Scale	Crop	Scheme	AEF-based RF		RS-based RF	
				R ²	RMSE	R ²	RMSE
Yield Prediction	County	Corn	East → West	0.02	2.71	0.66	1.61
		Corn	West → East	0.51	1.23	0.60	1.12
	Field	Soybean	East → West	-0.28	0.94	0.53	0.56
		Soybean	West → East	0.39	0.44	0.59	0.37
Field	Corn	Corn	East → West	0.25	2.65	0.38	2.41
		Corn	West → East	0.33	2.48	0.34	2.45
	Soybean	Soybean	East → West	0.08	0.99	0.19	0.93
		Soybean	West → East	0.23	1.00	0.28	0.97

Future of Foundation Models in EO

- Dealing with Seasonality
- Pipelines for easier inference & finetuning
- LoRa & Adapter based finetuning
- Physics constrained FMs
- Explainability of Predictions & Quantifying Uncertainty
- More pre-computed embeddings

Vision-Language Models for EO

- ChatGPT-style interaction with EO data

...

**Time to work on your Paper Presentation.
(~20min, groups)**

References & Further Material

Literature Recommendations

- Brown et Al, 2025, Alphaearth foundations: An embedding field model for accurate and efficient global mapping from sparse label data
- Tseng et Al, 2023, Lightweight, pre-trained transformers for remote sensing timeseries
- Xiao et Al, 2025, Foundation models for remote sensing and earth observation: A survey
- Vaswani et Al, 2017, Attention is all you need
- Herzog et Al, 2025, OlmoEarth: Stable Latent Image Modeling for Multimodal Earth Observation
- Lu et al, 2025, Vision Foundation Models in Remote Sensing: A survey
- Ma et al, 2025, Harvesting AlphaEarth: Benchmarking the Geospatial Foundation Model for Agricultural Downstream Tasks

References and Sources

- <https://www.youtube.com/watch?v=yg5I9ED05-s&t=2387s>
- <https://medium.com/google-earth/ai-powered-pixels-introducing-googles-satellite-embedding-dataset-31744c1f4650>
- <https://deepmind.google/blog/alphaearth-foundations-helps-map-our-planet-in-unprecedented-detail/>

**Interested in doing an Internship with NASA
Harvest @ Unistra?**