

# Advanced AI applications in EO: fine-tuning foundation models

Dr. Lorenzo Papa, Dr. Ruben Cartuyvels, Dr. Valerio Marsocci  
 Internal Research Fellows, ESA ESRIN Frascati, Rome, Italy  
 Email: [lorenzo.papa, ruben.cartuyvels, valerio.marsocci]@esa.int

- What are Geospatial Foundation Models?
- How do we obtain a GFM?
- Benchmarking GFMs
- Conclusion

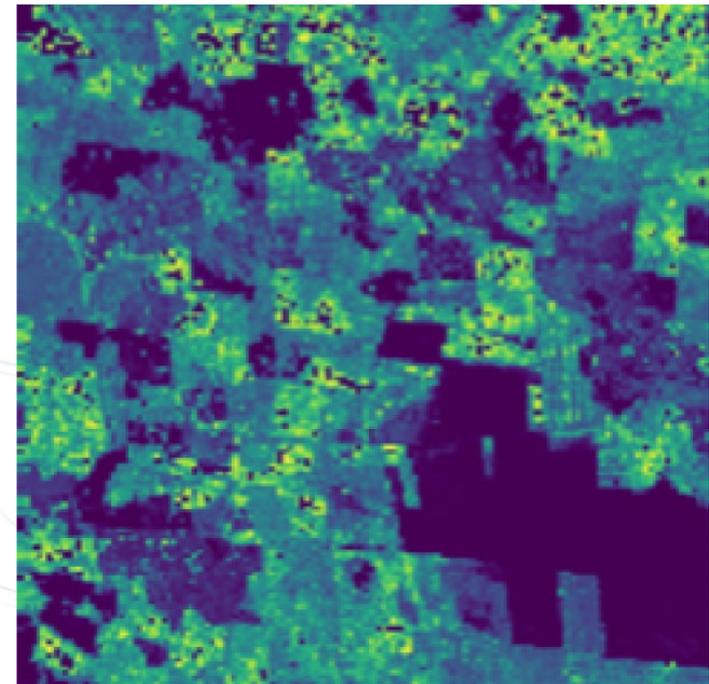
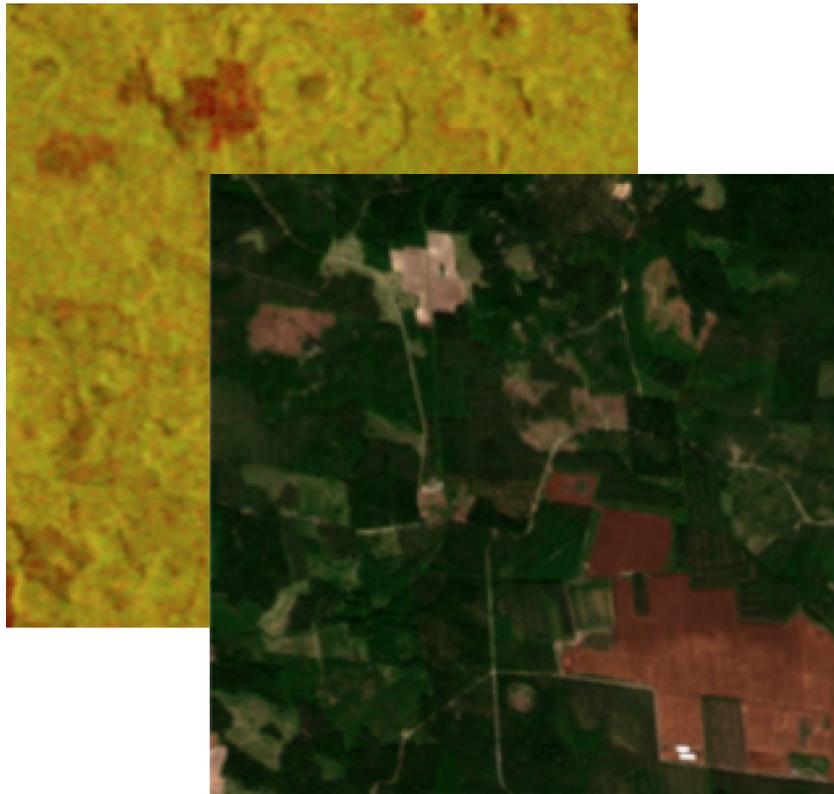
[https://github.com/lorenzopapa5/BiDs\\_training\\_course\\_2025/tree/main](https://github.com/lorenzopapa5/BiDs_training_course_2025/tree/main)

# WHAT ARE GEOSPATIAL FOUNDATION MODELS?

---

# Think of a problem connected to an SDG

15 LIFE ON LAND

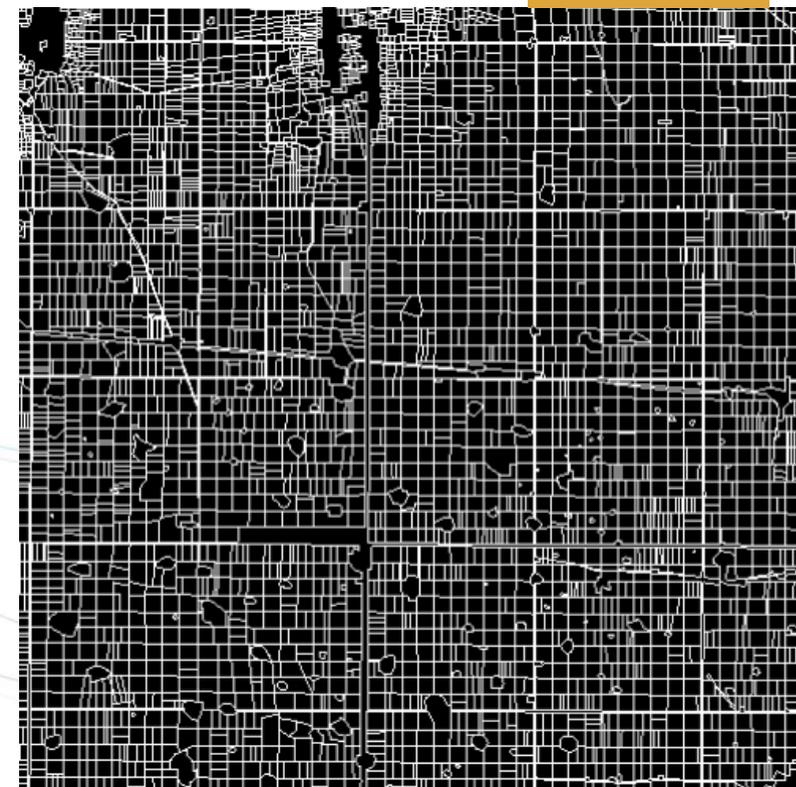
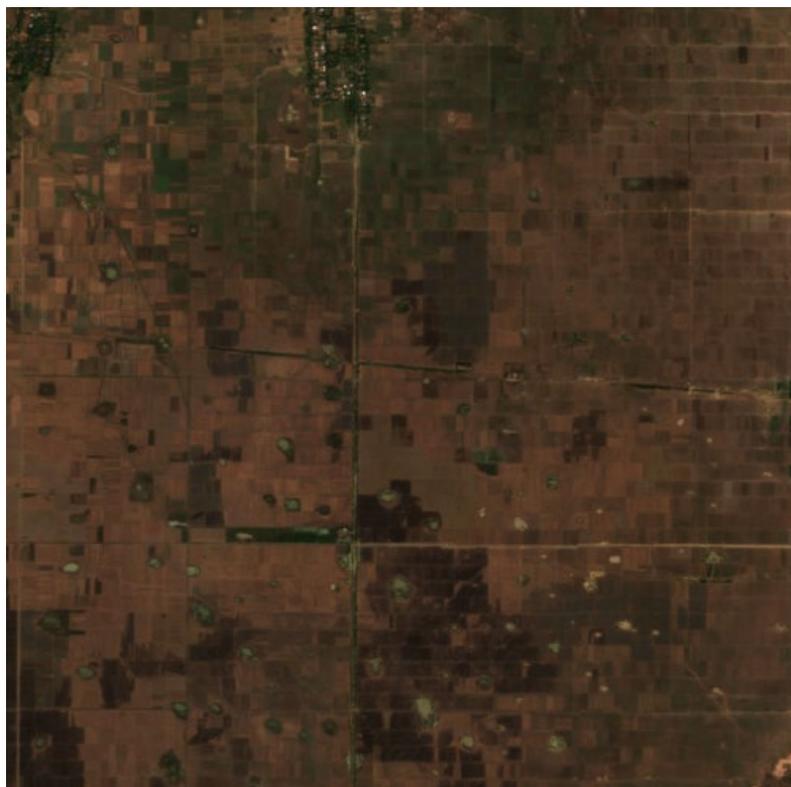


# Think of another one



# One more

2 ZERO HUNGER



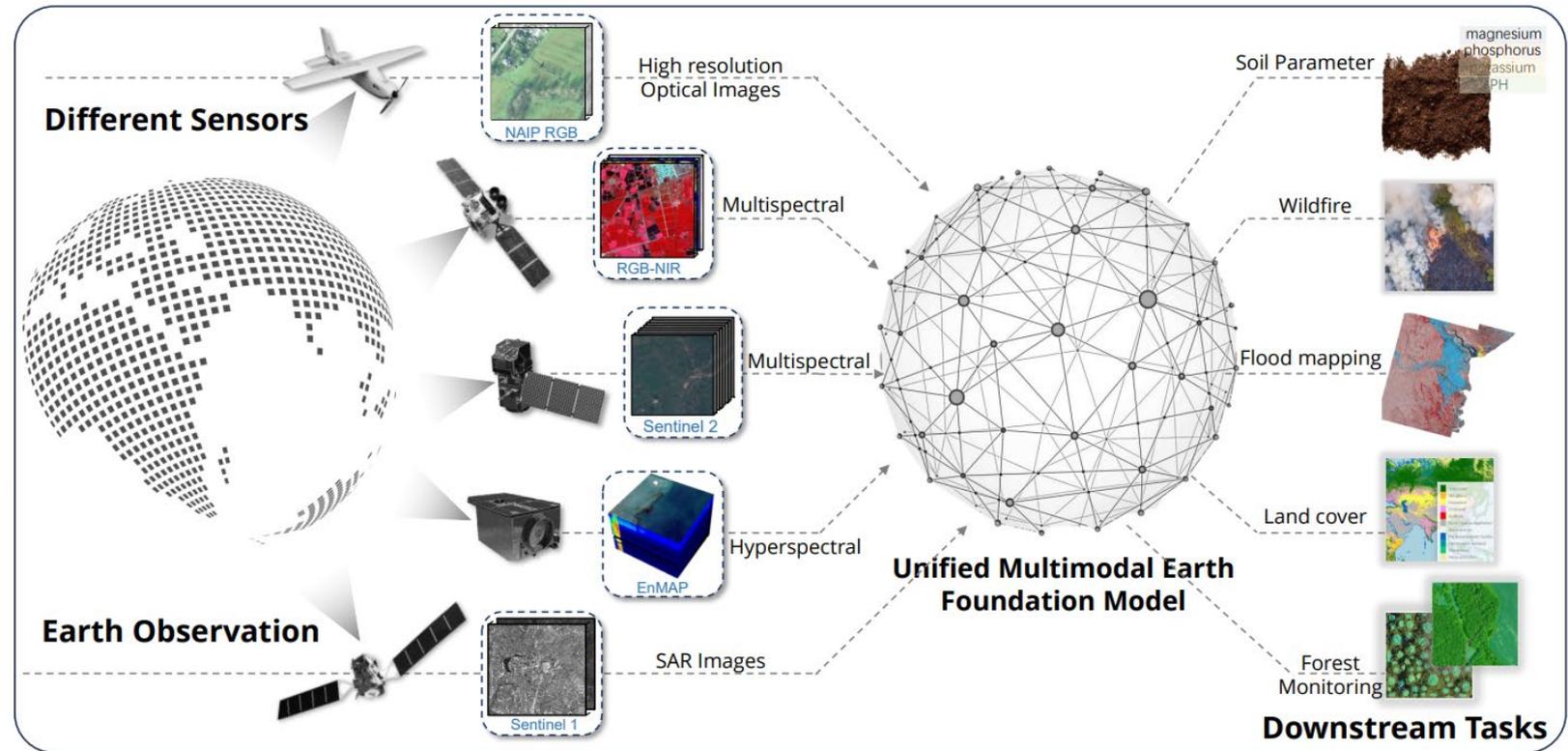
# How to solve all these tasks with the same model?



Φ-lab

## Let's use a Geospatial Foundation Model

- task agnosticism
- spatio-temporal awareness
- sensor agnosticism
- multimodality
- adaptability



Source: DOFA - [2403.15356](#)

- Geospatial foundation models are **large-scale** AI models trained on vast amounts of diverse **geospatial data**
- GFMs learn **temporal** and **spatial** patterns, relationships, and **features** across different geographic scales

- **Not limited to a single task**  
Can generalize across multiple geospatial applications
- **Enhances scalability & efficiency**  
Reduces need for retraining and Supports diverse use cases

Task  
Agnosticism

- Model understands both spatial (**geographic**) and temporal (**time-related**) patterns
- Essential for **dynamic applications** like climate monitoring, disaster response, and urban development

Spatio-temporal  
awareness

- Model can process data from **various sensors** (e.g., satellite, UAV, LiDAR)
- Increases **versatility**  
Allows **integration** of **diverse** data sources

Sensor  
agnosticism

- Model integrates and processes **multiple types** of data (e.g., imagery, text, GIS, data point)
- Provides richer insights by leveraging **diverse information**

## Multimodality

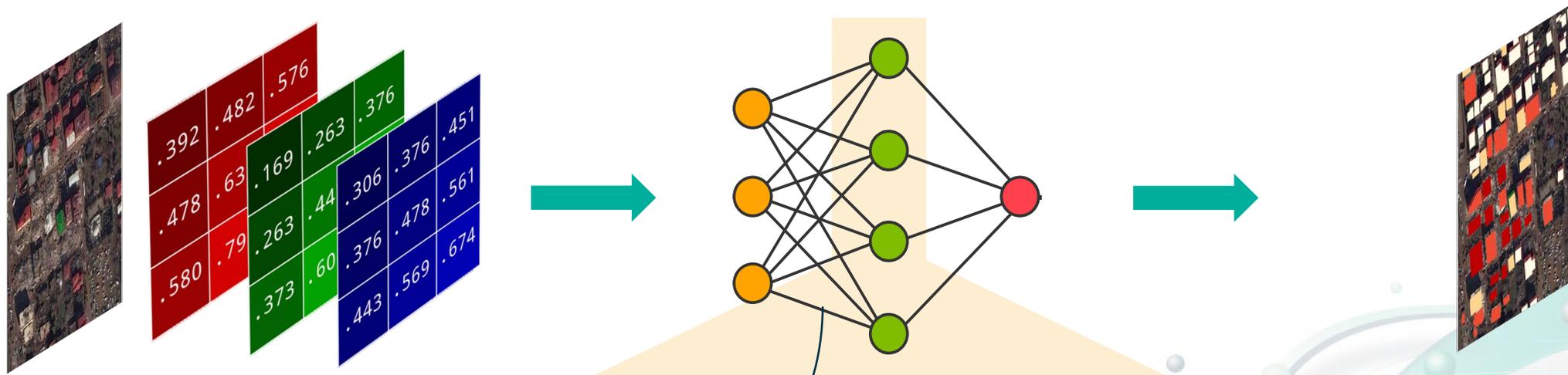
- Model can be **easily fine-tuned** or adjusted for different tasks or datasets
- **Reduces the need for retraining from scratch**, saving time and resources

## Adaptability

# DEVELOPING GEOSPATIAL FOUNDATION MODELS

---

# Neural networks (short recap)

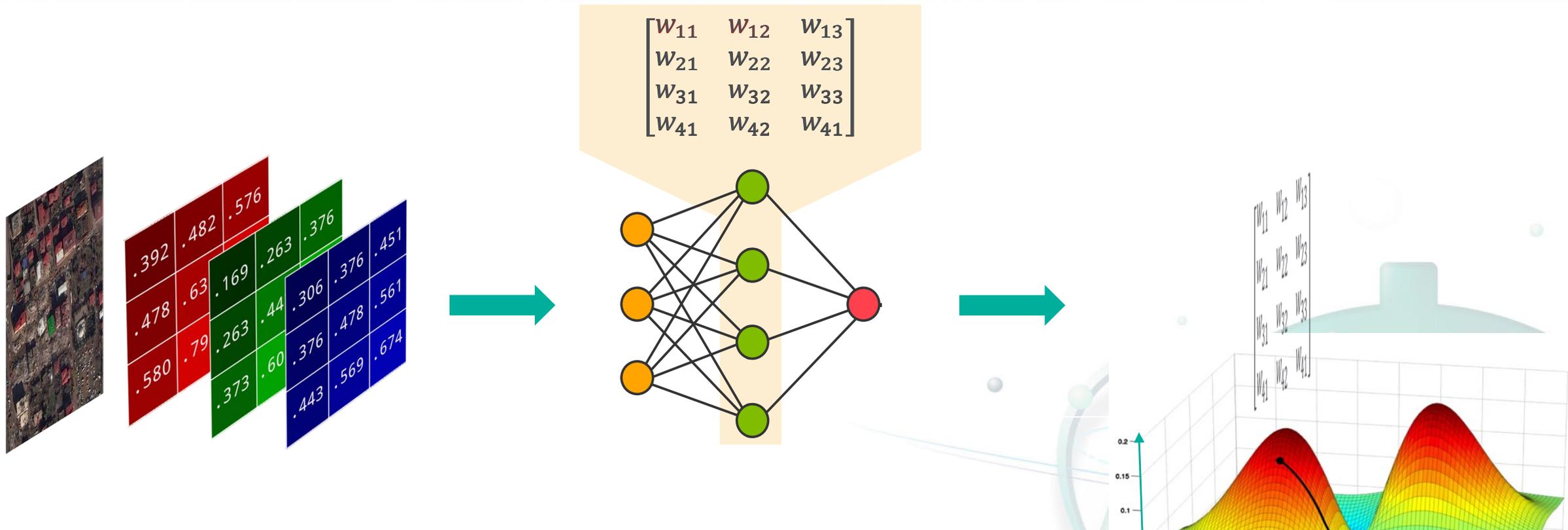


$$\begin{bmatrix} y_1 \\ y_2 \\ y_3 \\ y_4 \end{bmatrix} = \sigma \left( \begin{bmatrix} W_{11} & W_{12} & W_{13} \\ W_{21} & W_{22} & W_{23} \\ W_{31} & W_{32} & W_{33} \\ W_{41} & W_{42} & W_{43} \end{bmatrix} * \begin{bmatrix} .392 & .482 & .576 \\ .478 & .63 & .169 & .263 & .376 \\ .580 & .79 & .263 & .44 & .306 & .376 & .451 \\ .373 & .60 & .376 & .478 & .561 \\ .443 & .569 & .674 \end{bmatrix} + \begin{bmatrix} b_1 \\ b_2 \\ b_3 \\ b_4 \end{bmatrix} \right)$$

ESA UNCLASSIFIED - For ESA Official Use Only



# Neural networks (short recap)



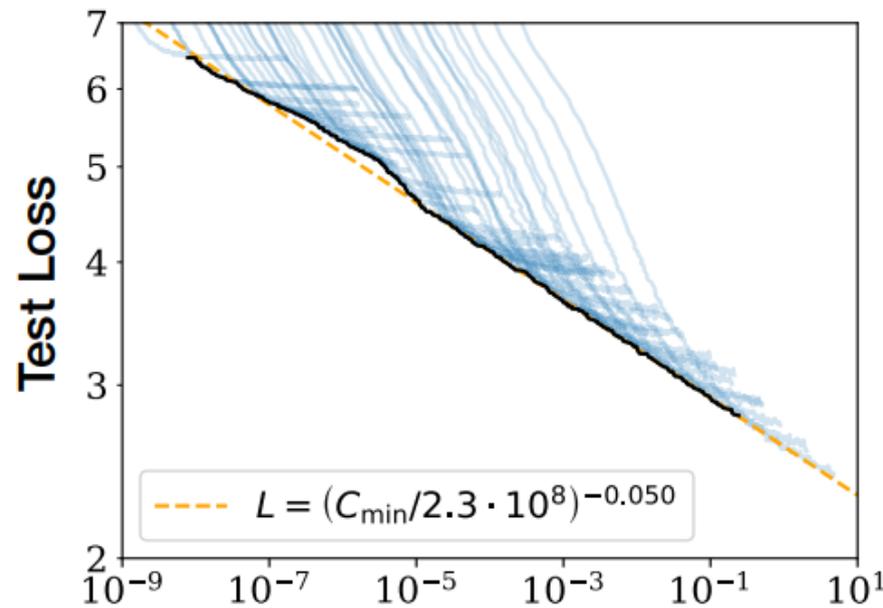
Objective function:

$$L(\text{Image}, \begin{bmatrix} W_{11} & W_{12} & W_{13} \\ W_{21} & W_{22} & W_{23} \\ W_{31} & W_{32} & W_{33} \\ W_{41} & W_{42} & W_{43} \end{bmatrix}) = 2.4 \times$$

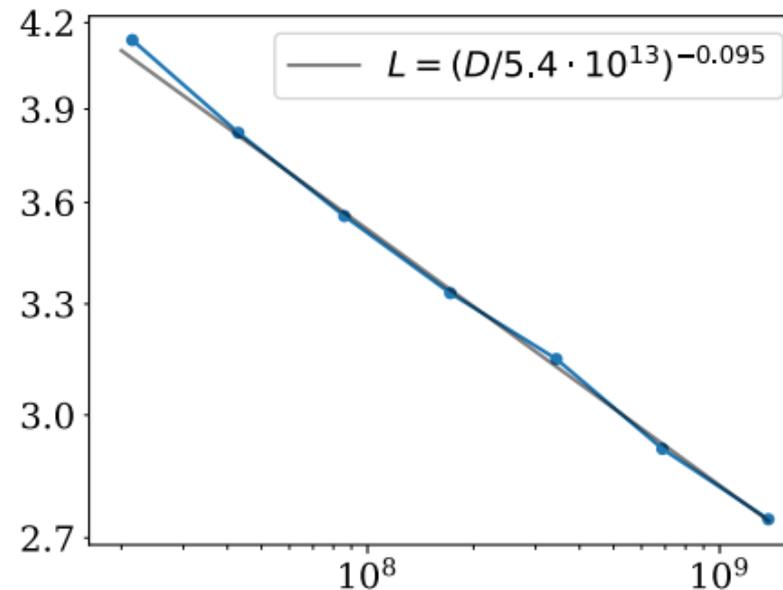
ESA UNCLASSIFIED - For ESA Official Use Only

Current w values  
Better w values

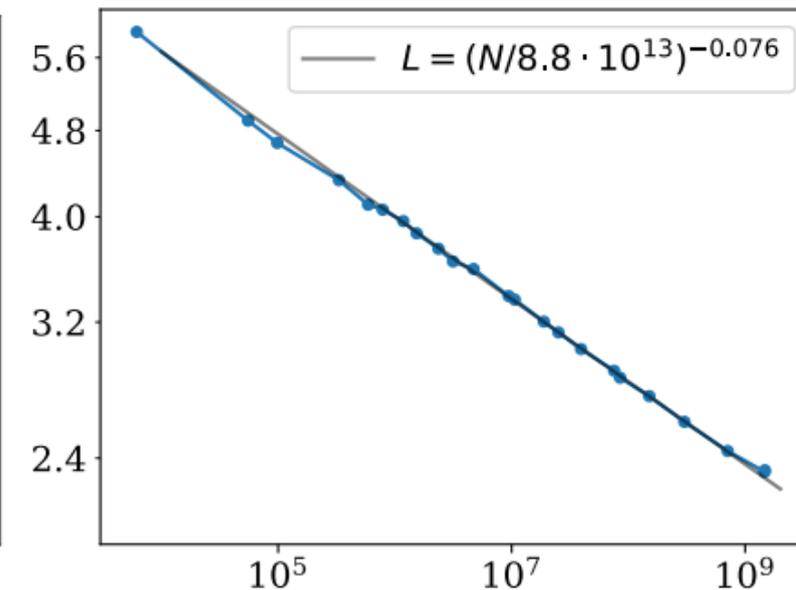
## Scaling Laws for Neural Language Models



**Compute**  
PF-days, non-embedding



**Dataset Size**  
tokens

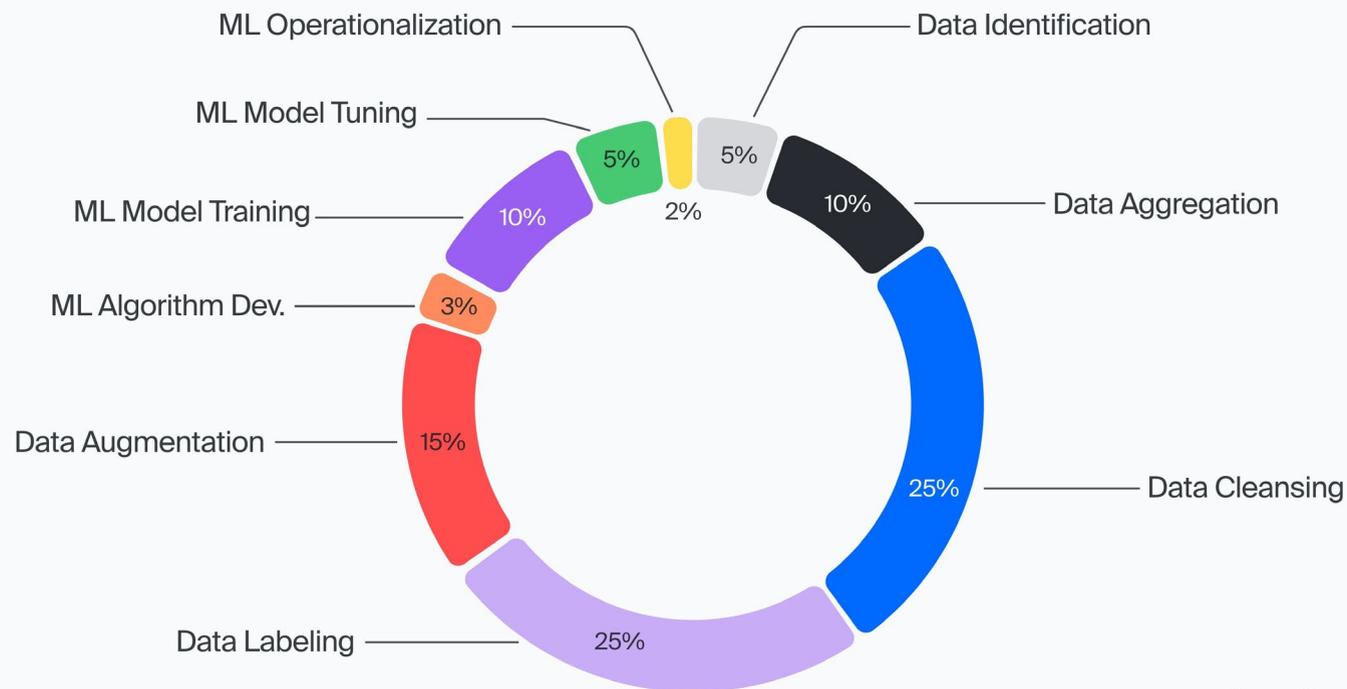


**Parameters**  
non-embedding

Kaplan, Jared, et al. "Scaling laws for neural language models." arXiv preprint arXiv:2001.08361 (2020).

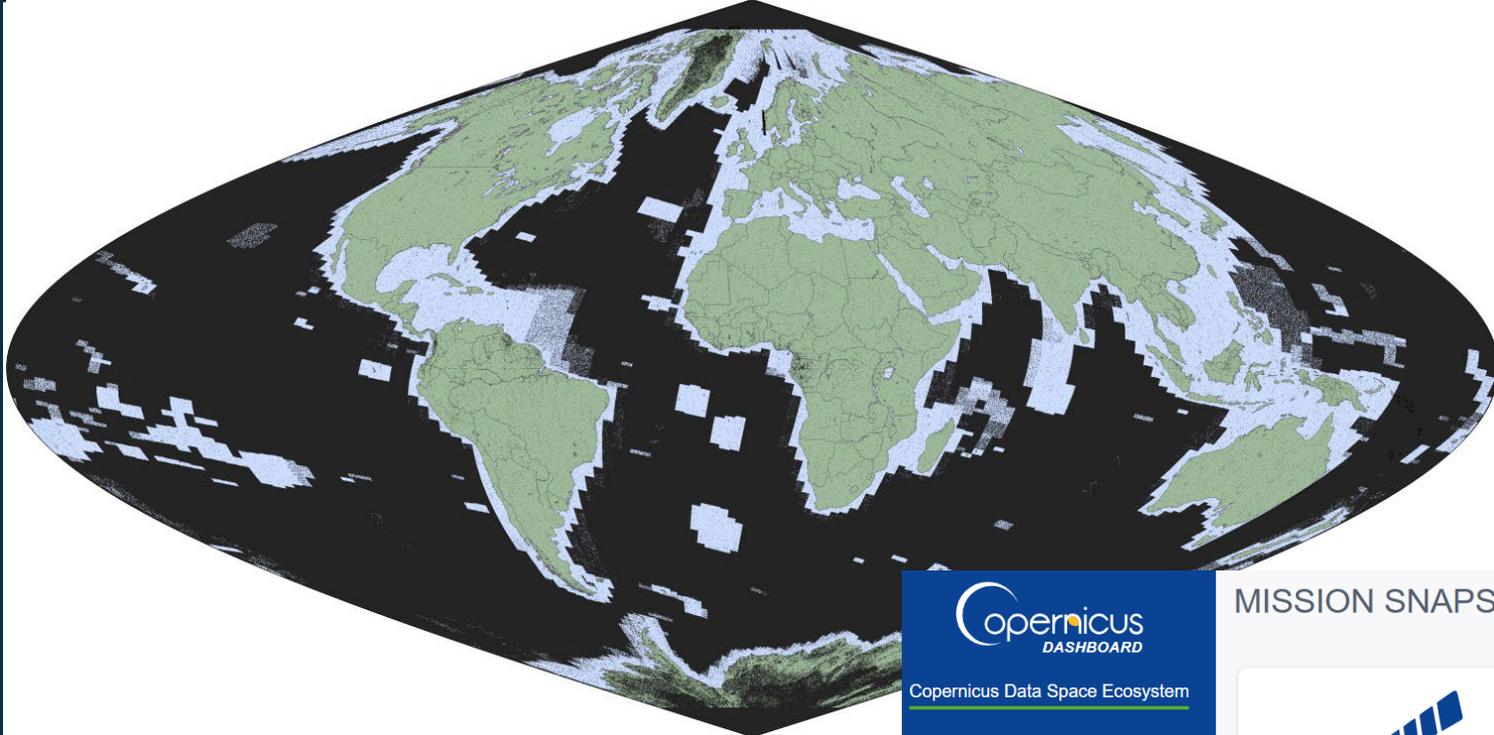
# (Labelled) data is expensive

### Percentage of Time Allocated to Machine Learning Project Tasks



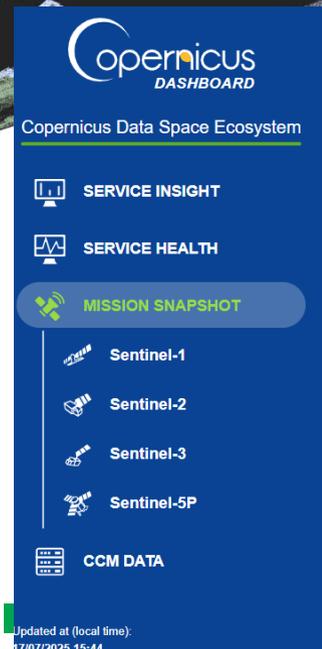
Source: Cognilytica

# There is a lot of EO data



<https://dashboard.dataspace.copernicus.eu/#/mission-snapshot>

<https://philab.esa.int/hello-major-tom-esa-%CF%86-lab-releases-largest-ml-ready-sentinel-2-dataset-ever-published/>



**Copernicus DASHBOARD**  
Copernicus Data Space Ecosystem

- SERVICE INSIGHT
- SERVICE HEALTH
- MISSION SNAPSHOT**
- Sentinel-1
- Sentinel-2
- Sentinel-3
- Sentinel-5P
- CCM DATA

Updated at (local time): 17/07/2025 15:44

## MISSION SNAPSHOT

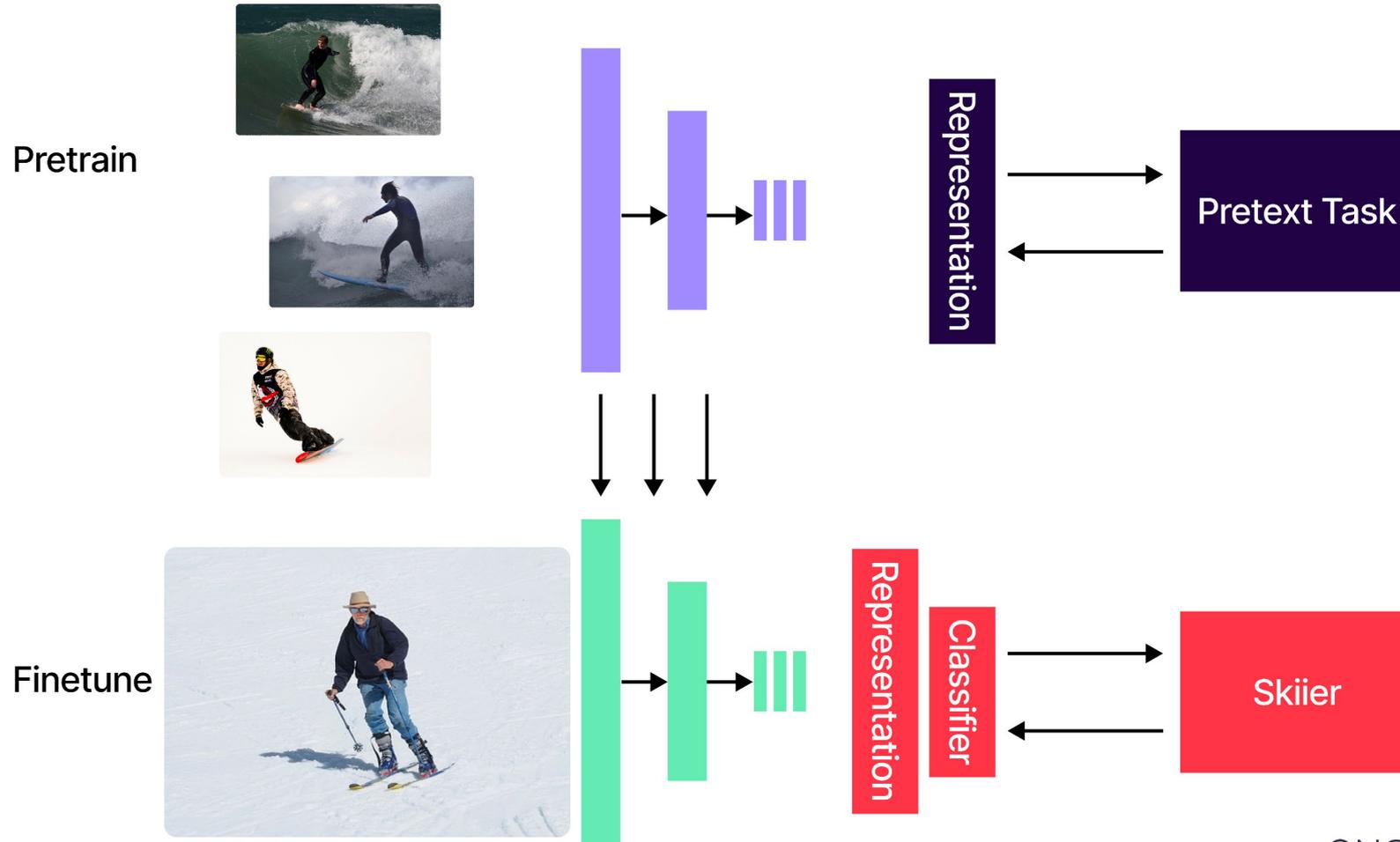
7d 30d

Sentinel-1	Sentinel-2	Sentinel-3	Sentinel-5P
			
<b>Sentinel-1</b>	<b>Sentinel-2</b>	<b>Sentinel-3</b>	<b>Sentinel-5P</b>
Total volume of published products			
<b>34.14 PB</b>	<b>52.17 PB</b>	<b>6.22 PB</b>	<b>1.75 PB</b>
Total number of published products			
<b>16.4M</b>	<b>108.3M</b>	<b>21.06M</b>	<b>1.59M</b>
Sentinel-1 Default Timeliness	Sentinel-2 Default Timeliness	Sentinel-3 OLCI, SLSTR, SRAL, NRT timelin...	Sentinel-5P NRT timeliness
<b>24h from sensing</b>	<b>24h from sensing</b>	<b>3h from sensing</b>	<b>3h from sensing</b>

ESA UNCLASSIFIED - For ESA Official Use Only



# Self-supervised learning

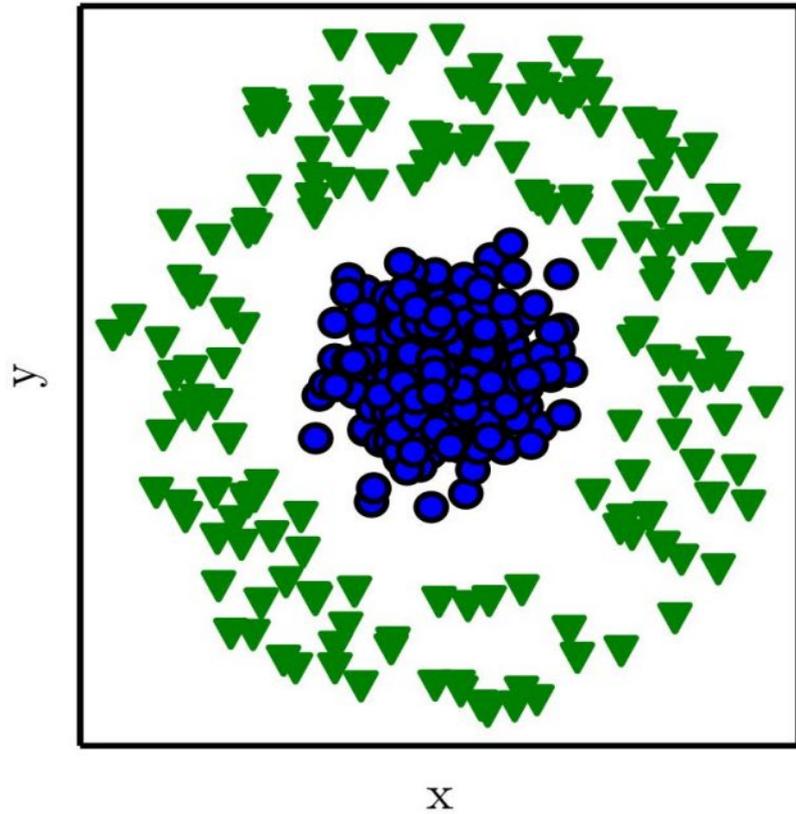


ENCORD

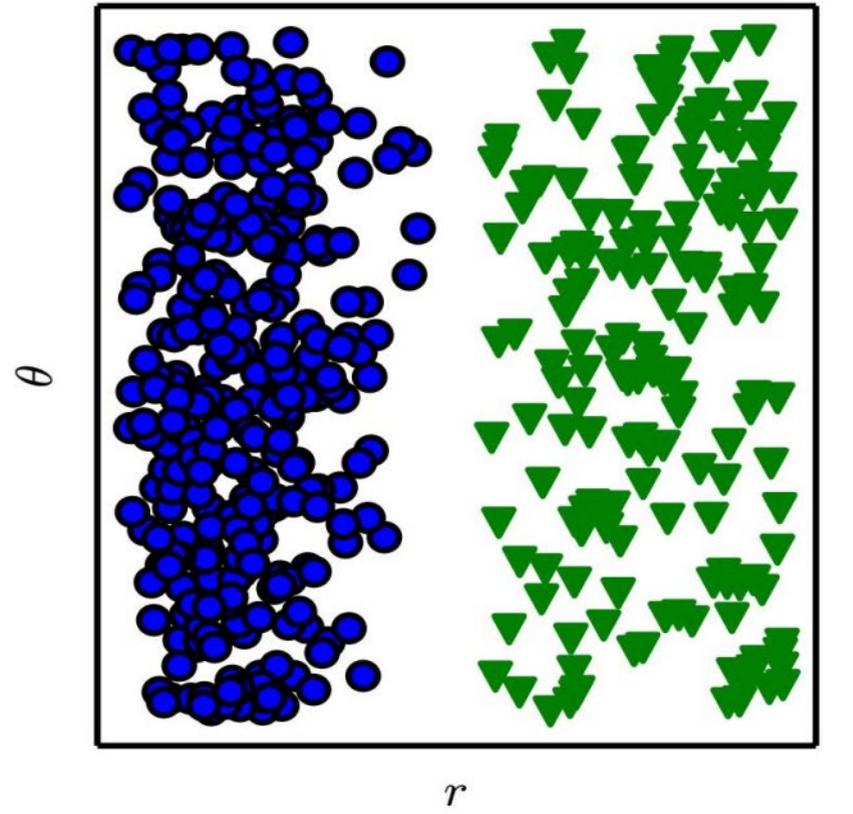
# Representations are important

- Rotation example

Cartesian coordinates

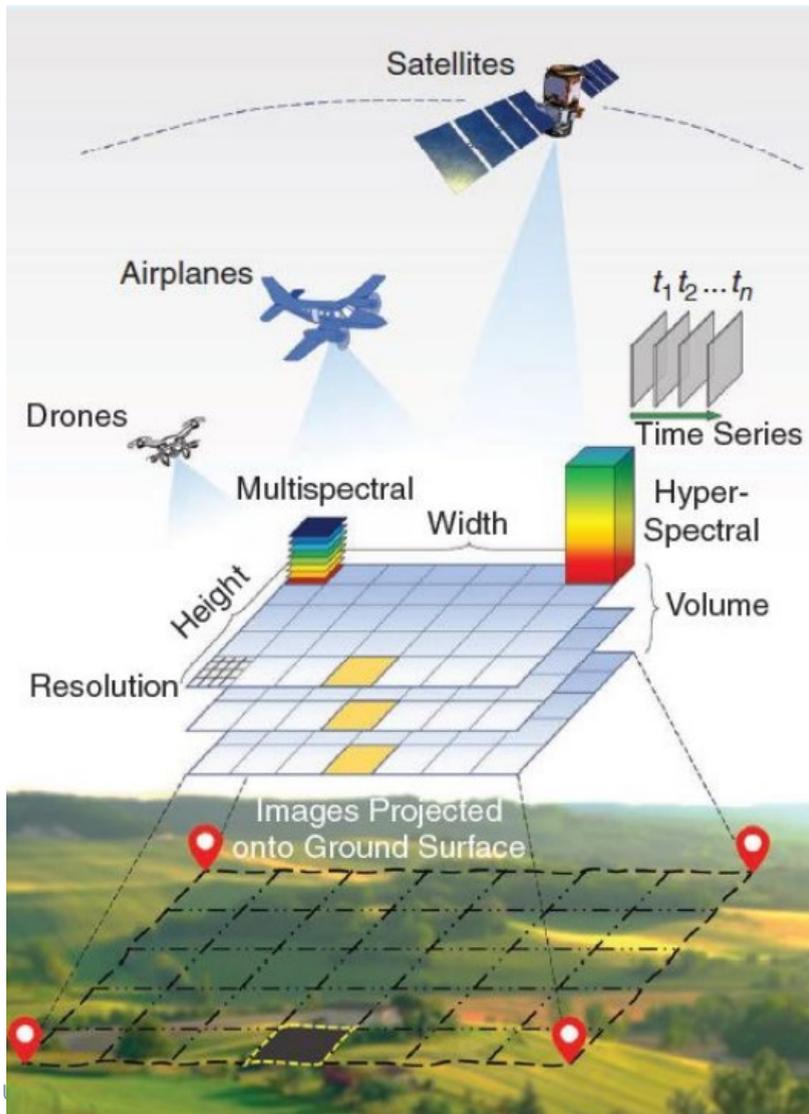


Polar coordinates



- We have:
  - Large dataset of unlabelled inputs  $x$
  - Small downstream dataset(s) of input, target pairs  $(x, y)$
- Pretext task:
  - Learn model  $f: x_{in} \rightarrow x_{out}$
  - For transformation  $g: x \rightarrow (x_{in}, x_{out})$
- Goal:
  - Obtain strong model (representations)  $f$  from large dataset
  - That improves supervised prediction by downstream model  $f_{down}$  on small dataset:  
 $y = f_{down}(f(x))$

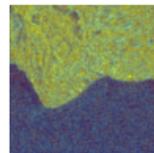
# Input modalities



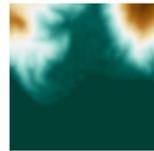
Sentinel-2 L2A  
Sentinel-2 L1C  
Sentinel-2 RGB



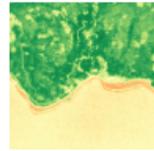
Sentinel-1 GRD  
Sentinel-1 RTC



Digital Elevation Model (DEM)



Vegetation index (NDVI)



Land-use/land-cover maps (LULC)

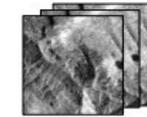


Image caption

A satellite image of a coastline with ...

Coordinates at 0.25 x 0.25 deg

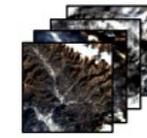
0°30'N  
122°45'E



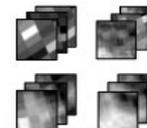
Sentinel-1



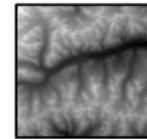
Sentinel-2



Sentinel-3



Sentinel-5P



GLO30 DEM

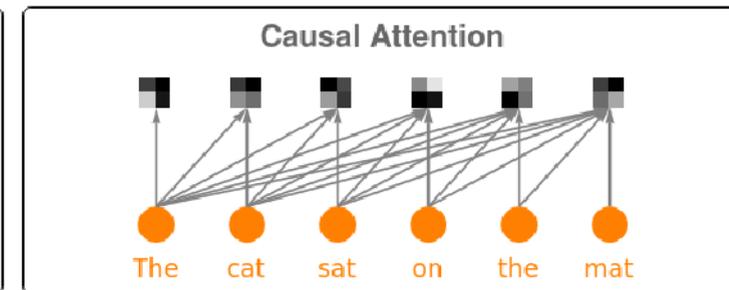
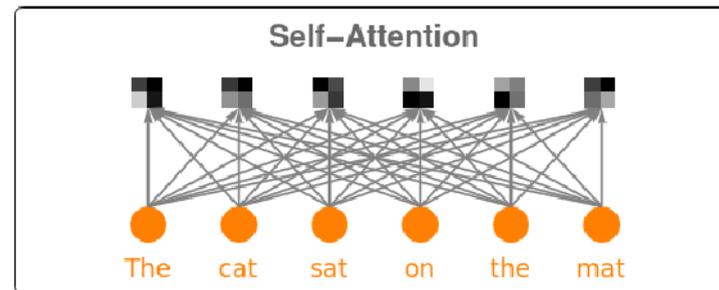
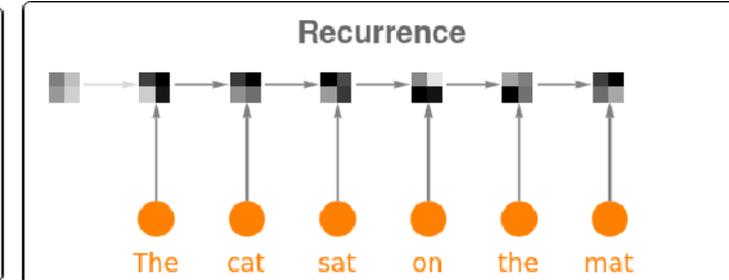
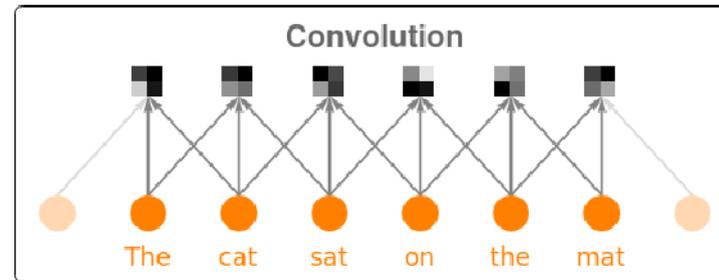
Wang, Yi, et al. "Towards a unified copernicus foundation model for earth vision." arXiv preprint arXiv:2503.11849 (2025).

Tseng, Gabriel, et al. "Lightweight, pre-trained transformers for remote sensing timeseries." arXiv preprint arXiv:2304.14065 (2023).



Jakubik, Johannes, et al. "Terramind: Large-scale generative multimodality for earth observation." arXiv preprint arXiv:2504.11171 (2025).





## Convolutional Neural Networks (CNNs):

Local, spatial feature extraction

## Recurrent Neural Networks (RNNs):

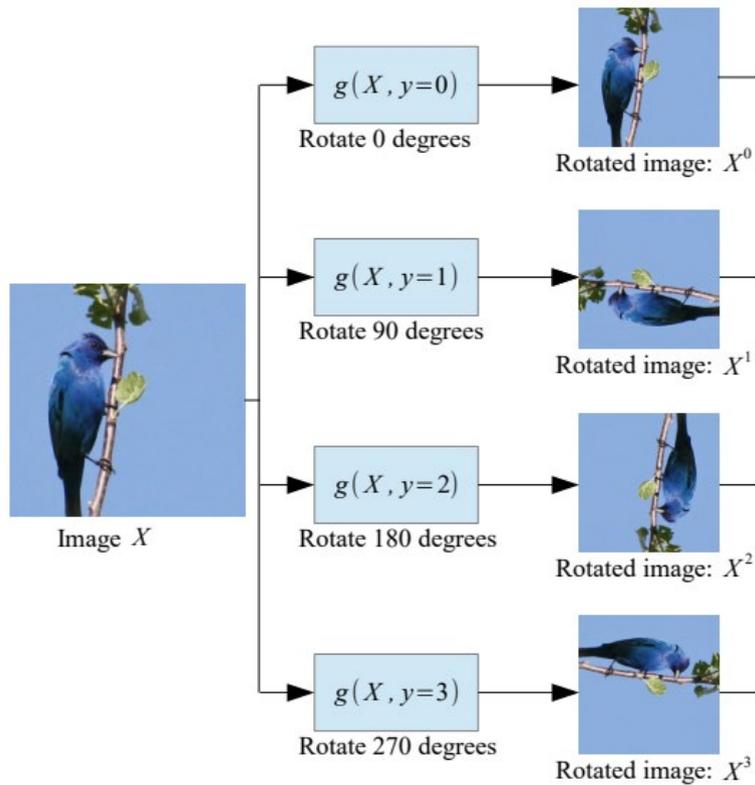
Process sequential data step by step, ideal for time-series analysis

## Transformers:

(Self-)attention to capture global spatial, temporal, inter-modality relationships

Train in parallel so especially fit for large datasets

<https://www.wolfram.com/language/12/neural-network-framework/use-transformer-neural-nets.html>



SSL:

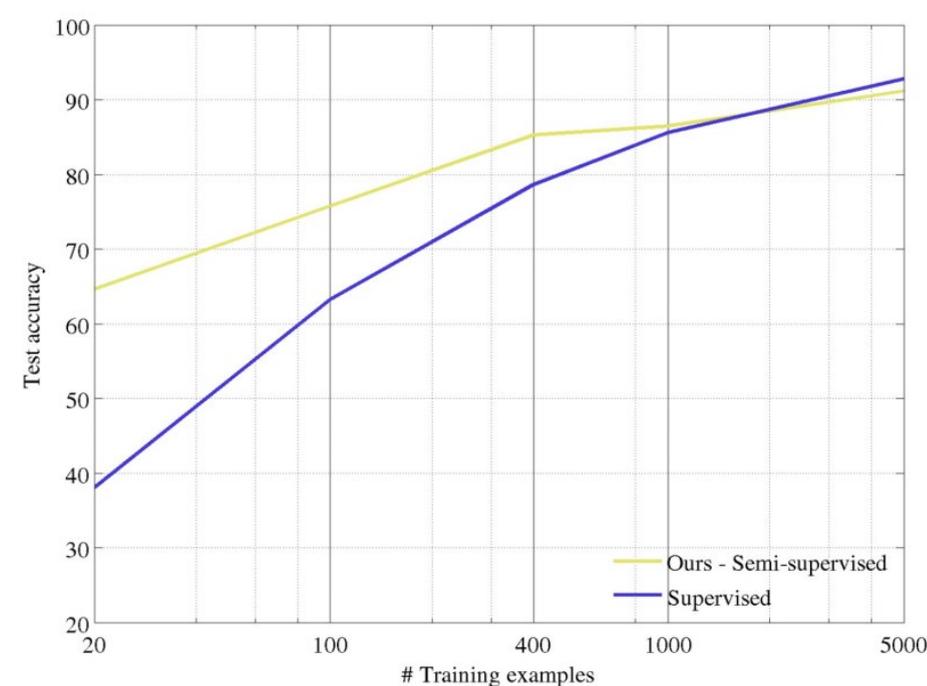
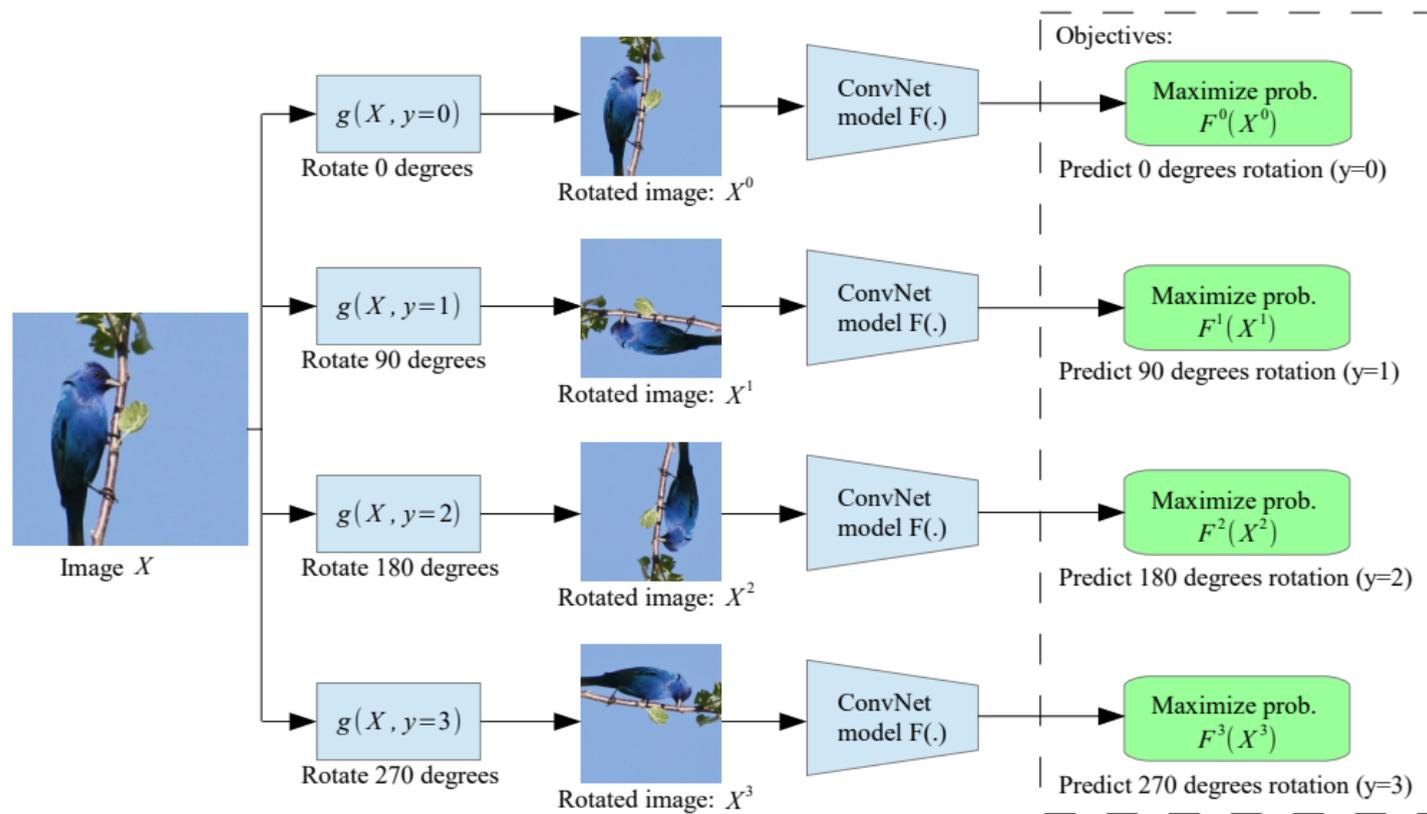
Learn model  $f: x_{in} \rightarrow x_{out}$

For transformation  $g: x \rightarrow (x_{in}, x_{out})$

Learning rotations:

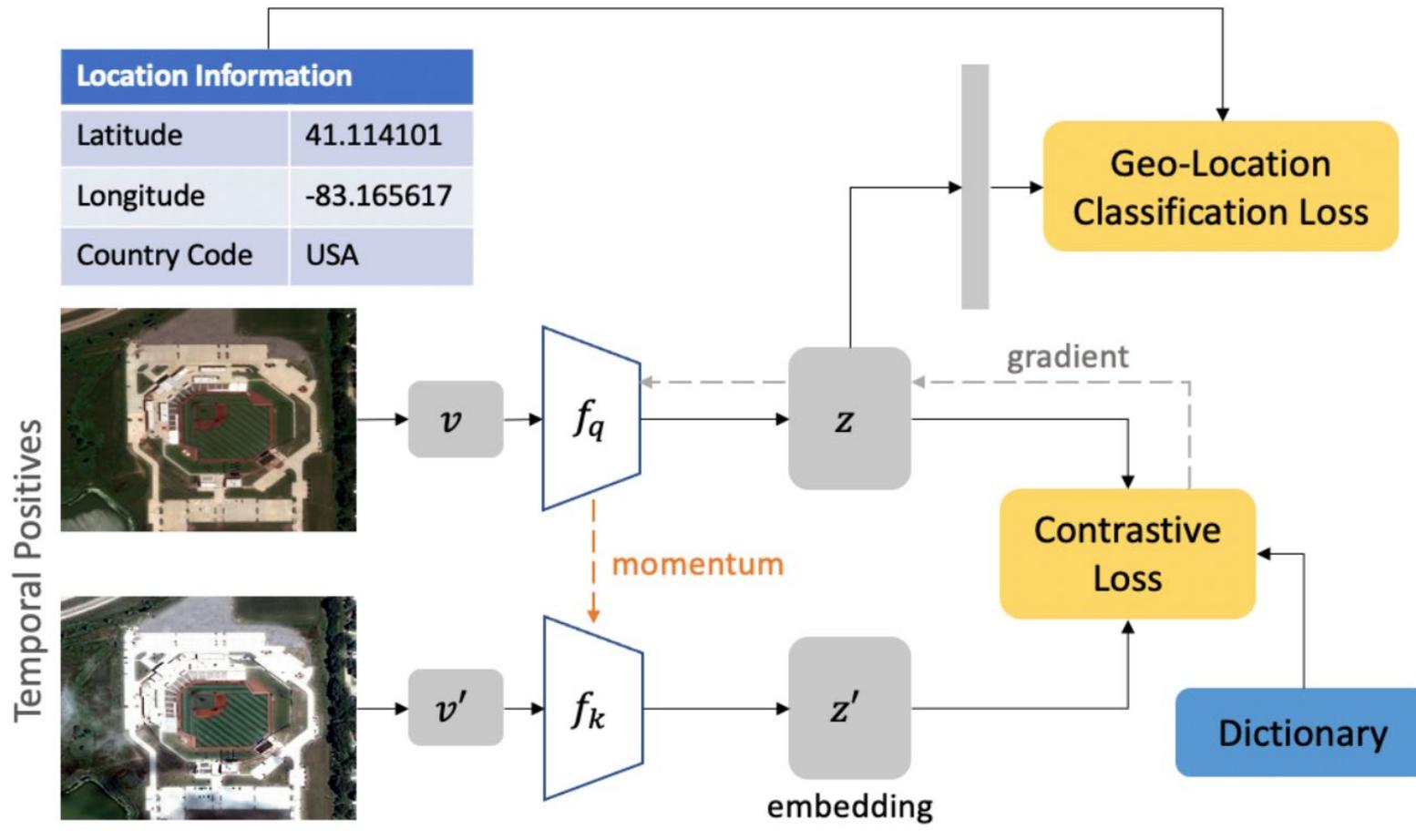
- $x_{out} = angle \in \{270, 180, 90, 0\}$
- $g: x \rightarrow (x_{rotated}, angle)$

# Learning metadata



Gidaris, Spyros, Praveer Singh, and Nikos Komodakis. "Unsupervised representation learning by predicting image rotations." arXiv preprint arXiv:1803.07728 (2018).

# Learning metadata: Geography-aware self-supervised learning



Ayush, Kumar, et al. "Geography-aware self-supervised learning." Proceedings of the IEEE/CVF International Conference on Computer Vision. 2021.

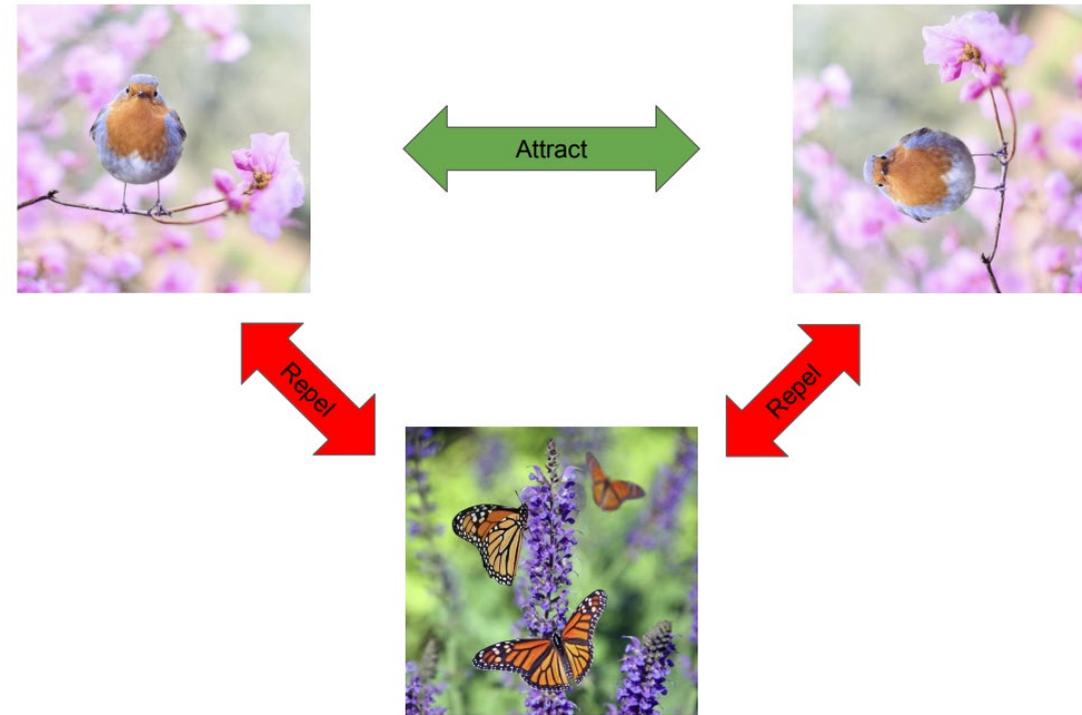
SSL:

Learn model  $f: x_{in} \rightarrow x_{out}$   
For transformation  $g: x \rightarrow (x_{in}, x_{out})$

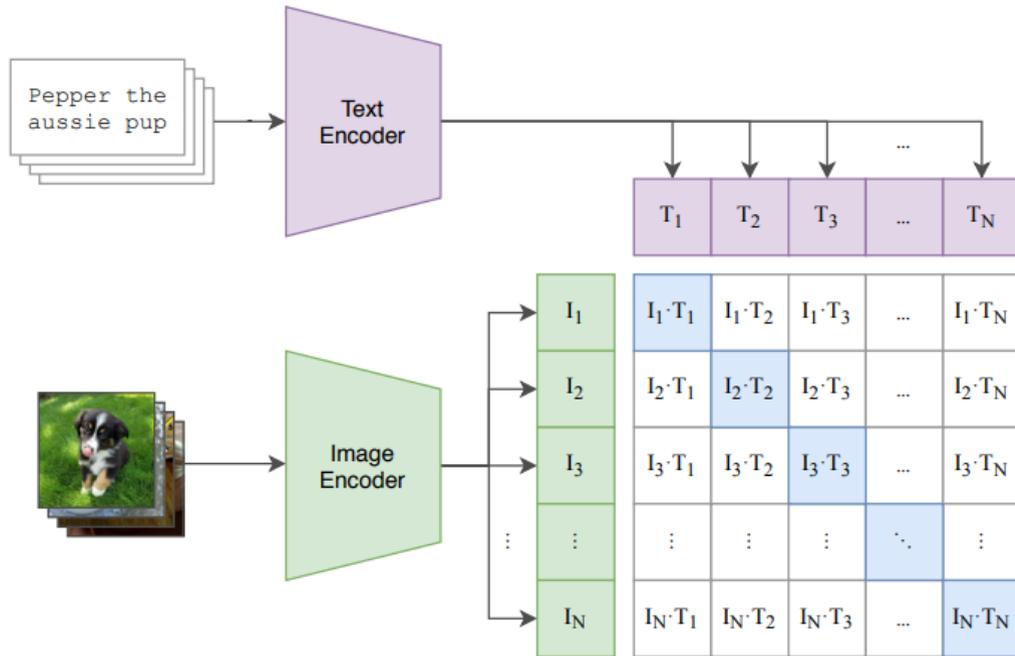
Contrastive:

Learn model  $f: x \rightarrow \mathbb{R}^D$   
For positive and negative view  $x^+, x^-$   
And similarity  $s: x \rightarrow \mathbb{R}^1$   
So that:

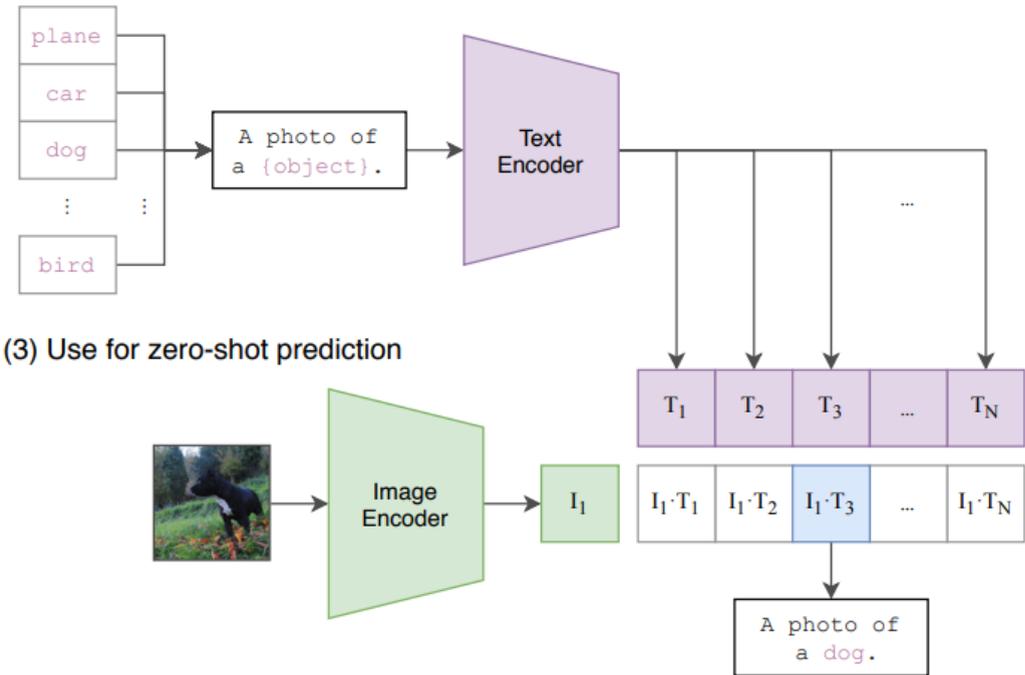
$$s(f(x), f(x^+)) > s(f(x), f(x^-))$$



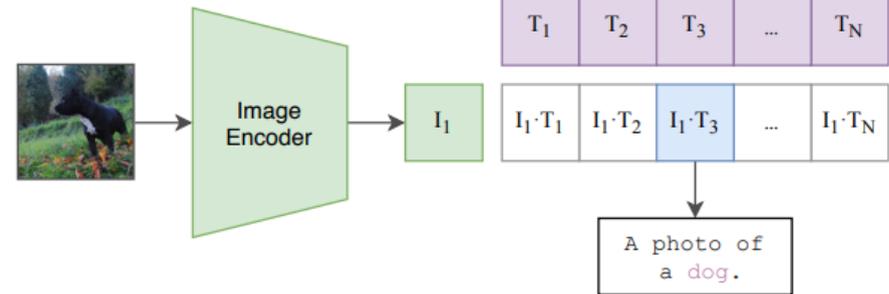
### (1) Contrastive pre-training



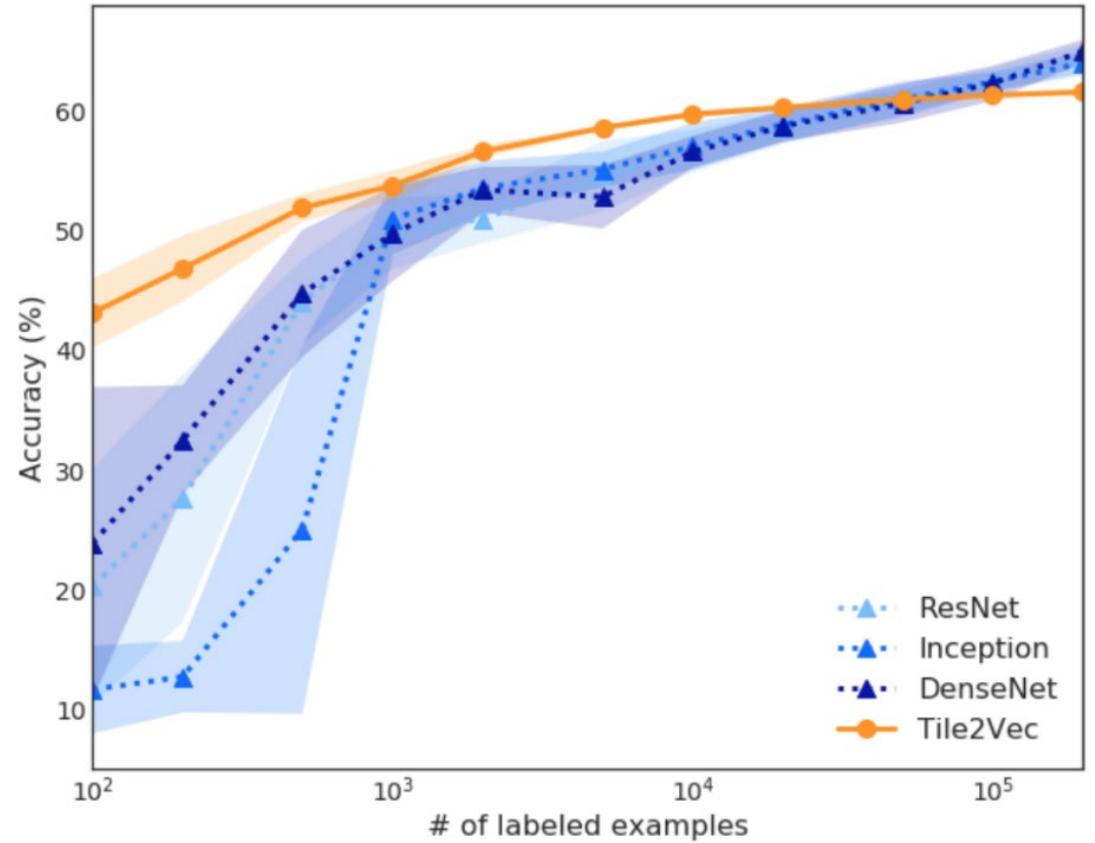
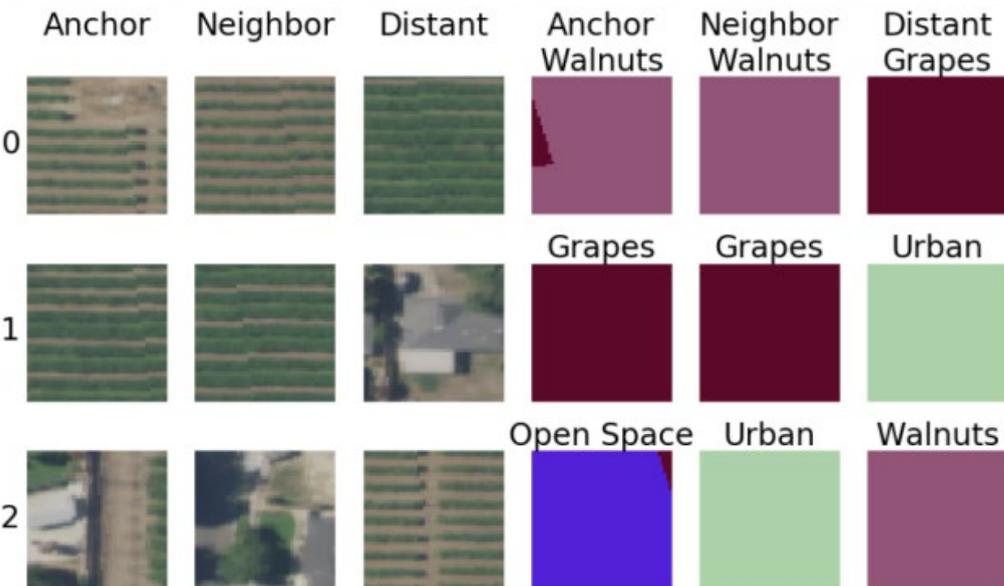
### (2) Create dataset classifier from label text



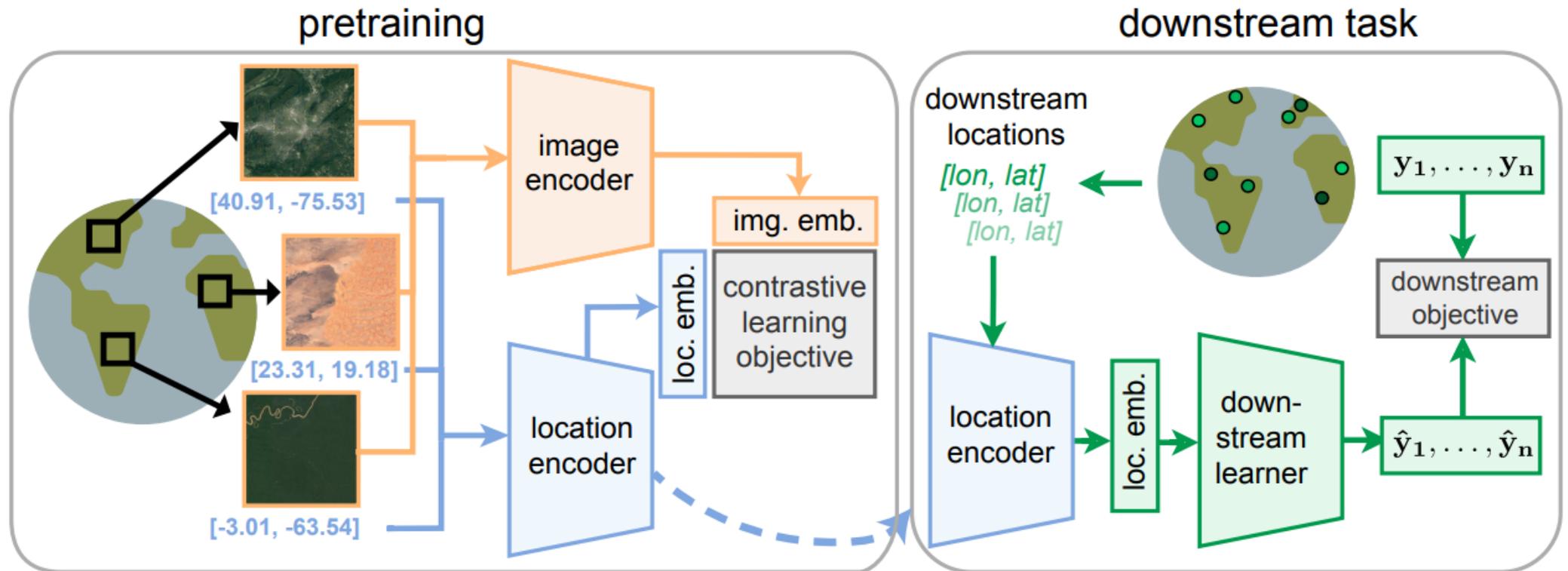
### (3) Use for zero-shot prediction



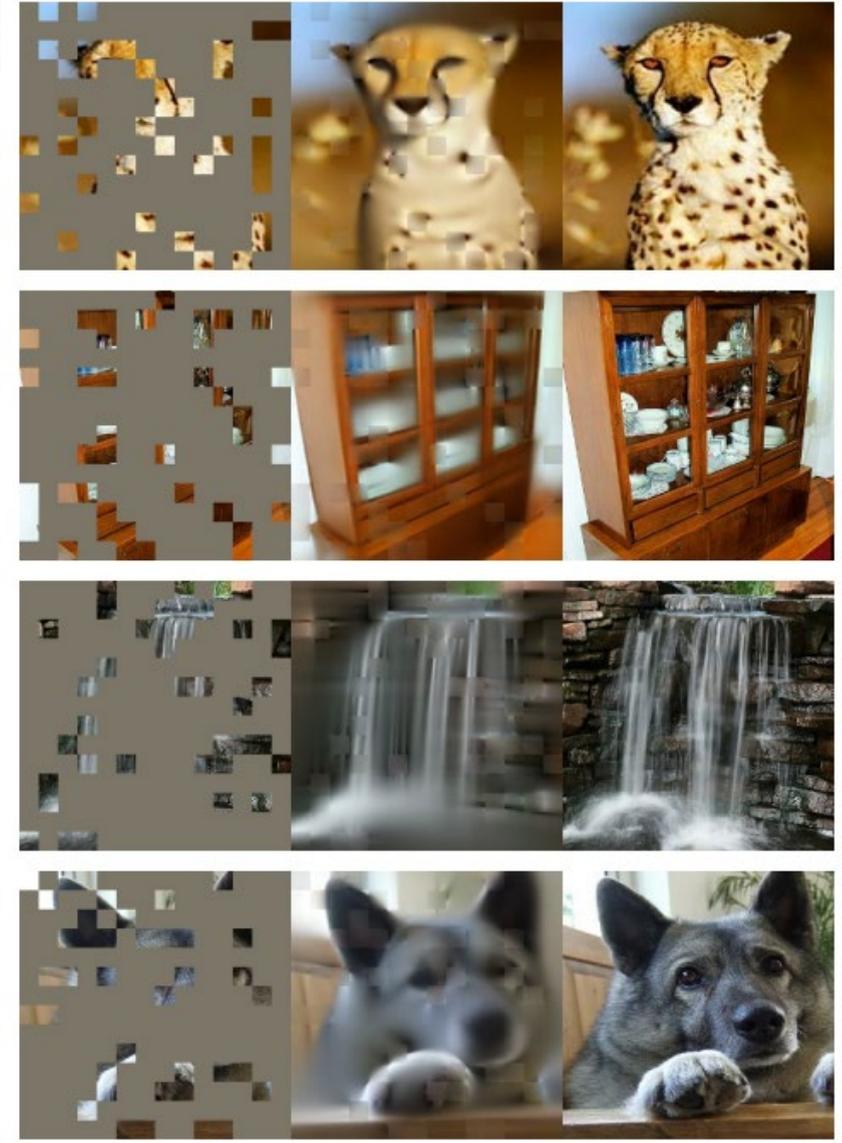
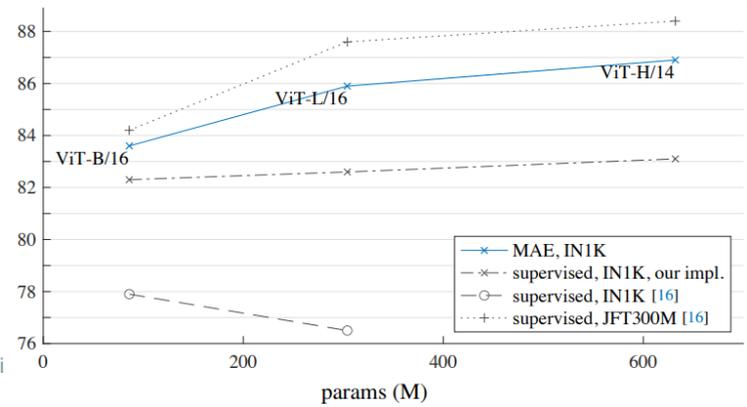
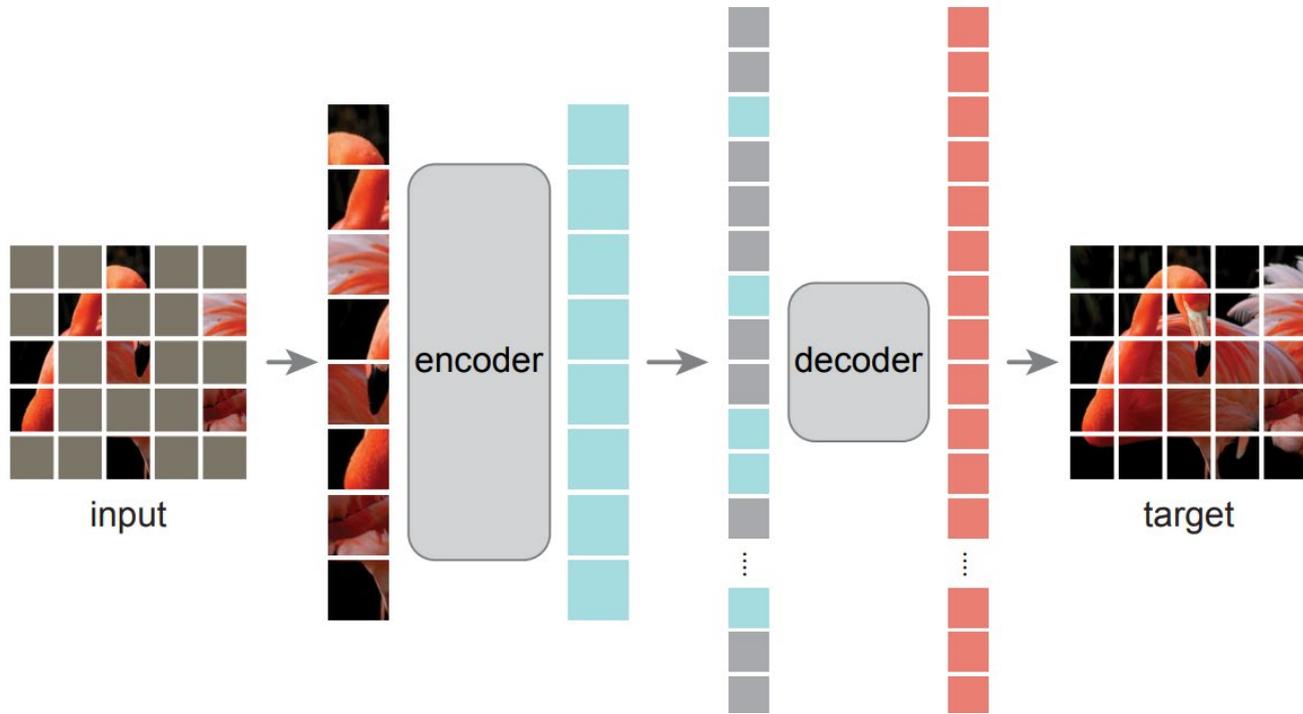
Radford, Alec, et al. "Learning transferable visual models from natural language supervision." International conference on machine learning. PmLR, 2021.



Jean, Neal, et al. "Tile2vec: Unsupervised representation learning for spatially distributed data." Proceedings of the AAAI Conference on Artificial Intelligence. Vol. 33. No. 01. 2019.

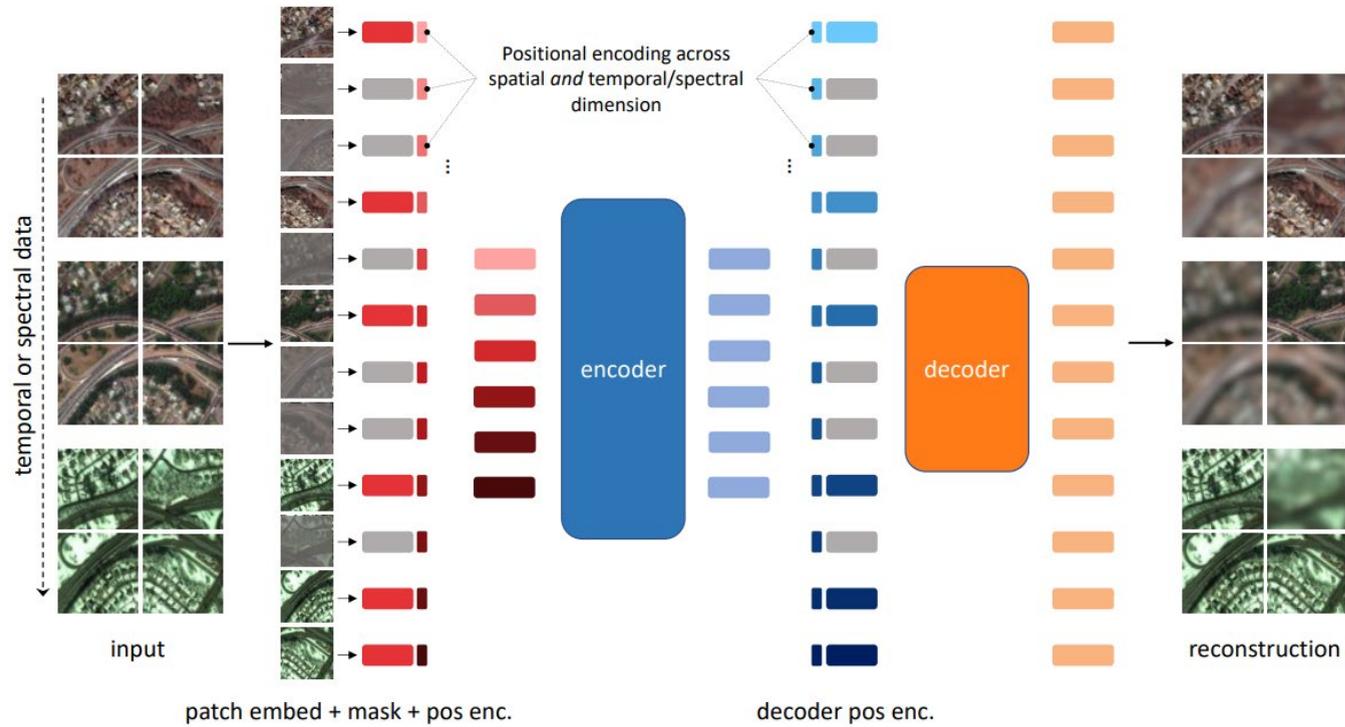


# Reconstruction

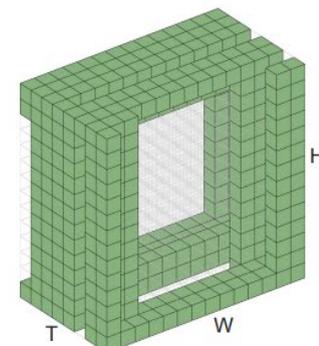
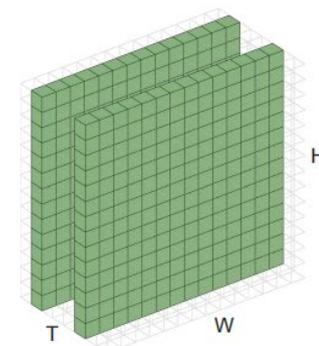
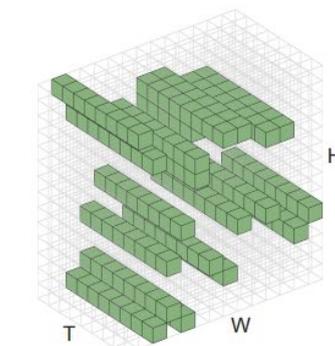
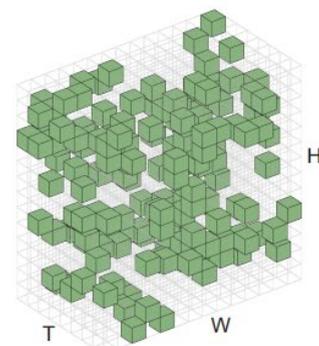


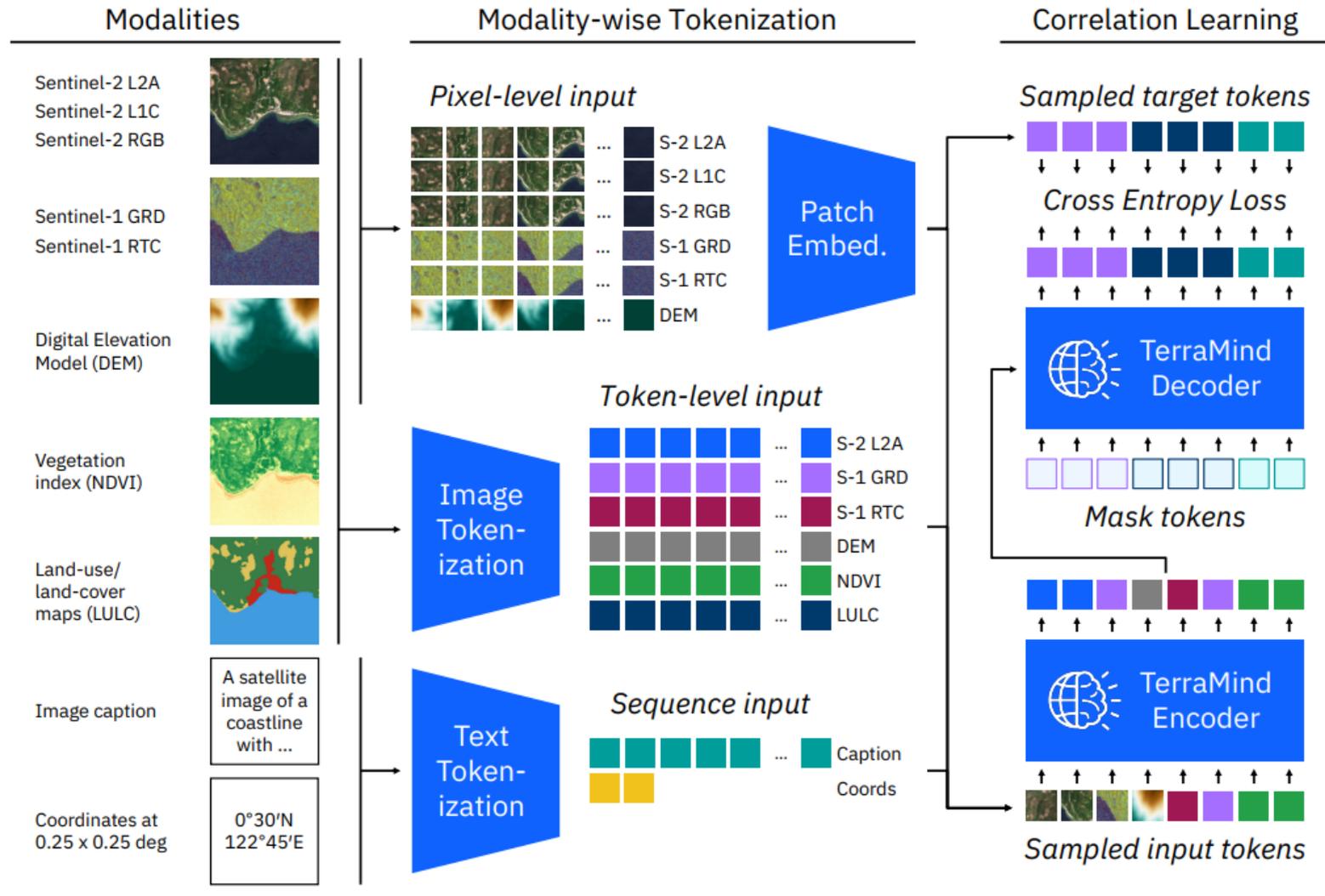
He, Kaiming, et al. "Masked autoencoders are scalable vision learners." Proceedings of the IEEE/CVF conference on computer vision and pattern recognition. 2022.

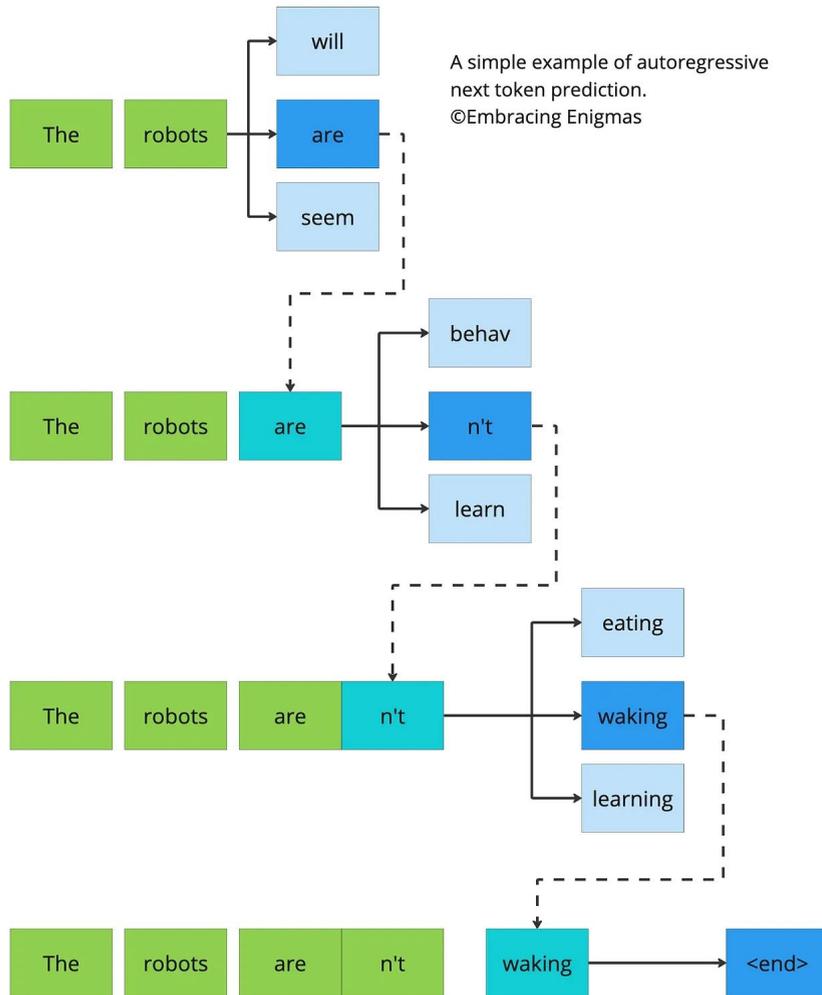




Feichtenhofer, Christoph, Yanghao Li, and Kaiming He. "Masked autoencoders as spatiotemporal learners." Advances in neural information processing systems 35 (2022): 35946-35958.







Learn a model  $p_{\theta}(x)$  that approximates the true underlying probability distribution  $p_{data}(x)$

Generate new data by sampling from the learned model:  $x_{new} \sim p_{\theta}(x)$

Foundation model: use learned internal representations as inputs for various tasks

E.g. autoregressive language models:

$$p_{\theta}(x) = \prod_i p_{\theta}(x_i | x_0, \dots, x_{i-1})$$

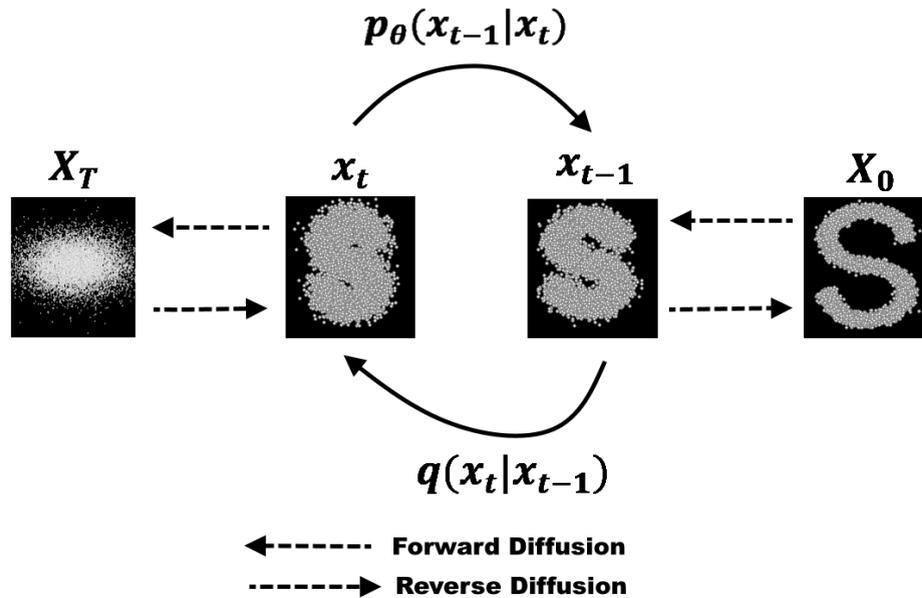
$$p_{\theta}(the, robots, are) = p_{\theta}(the) \cdot p_{\theta}(robots | the) \cdot p_{\theta}(are | the, robots)$$

<https://embracingenigmas.substack.com/p/next-token-prediction-is-a-fundamental>

Learn to reverse a process of adding noise

### Forward (Fixed):

Start with a real image  $x_0$   
Gradually add Gaussian noise until image equals pure noise  $x_T$



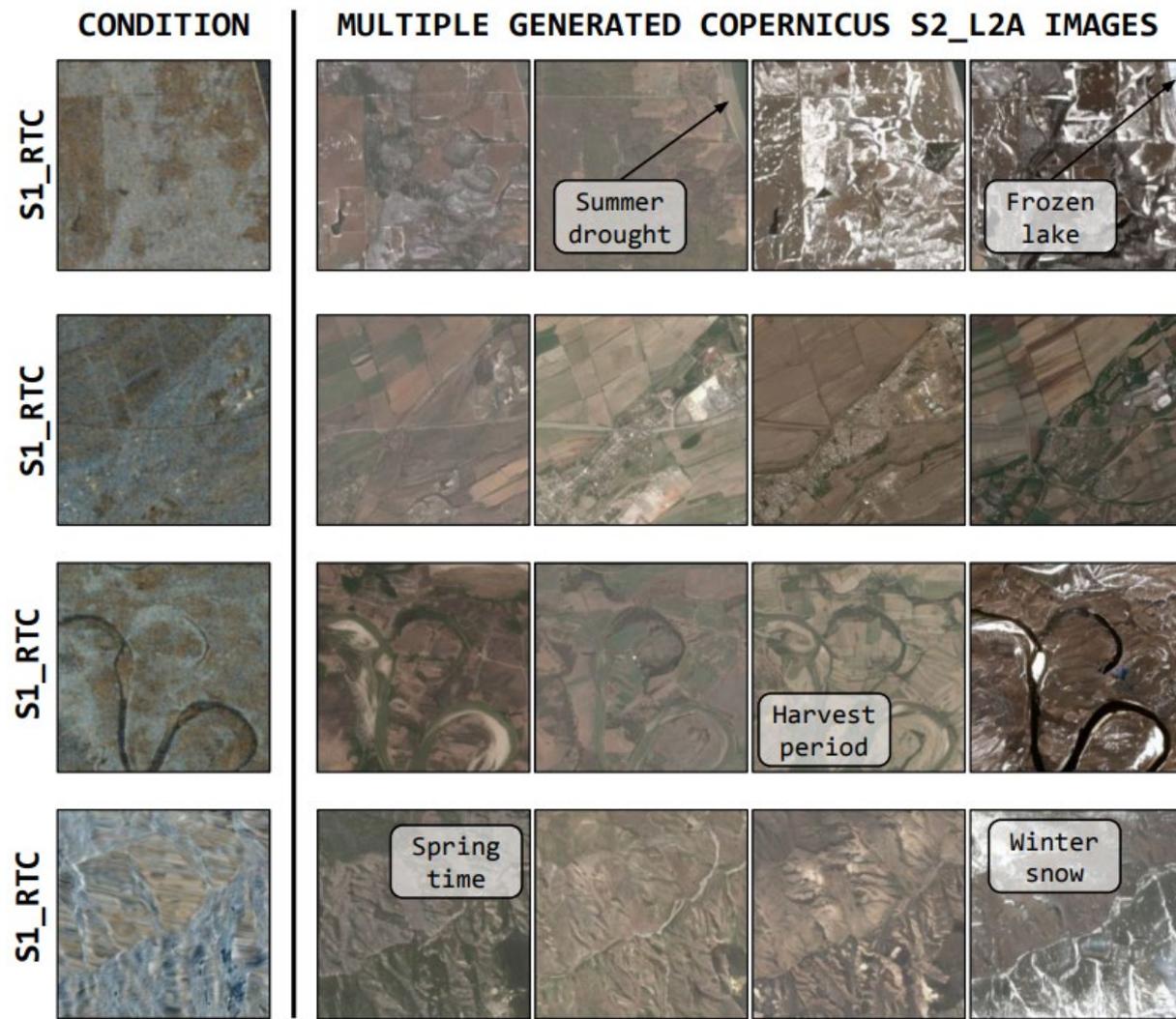
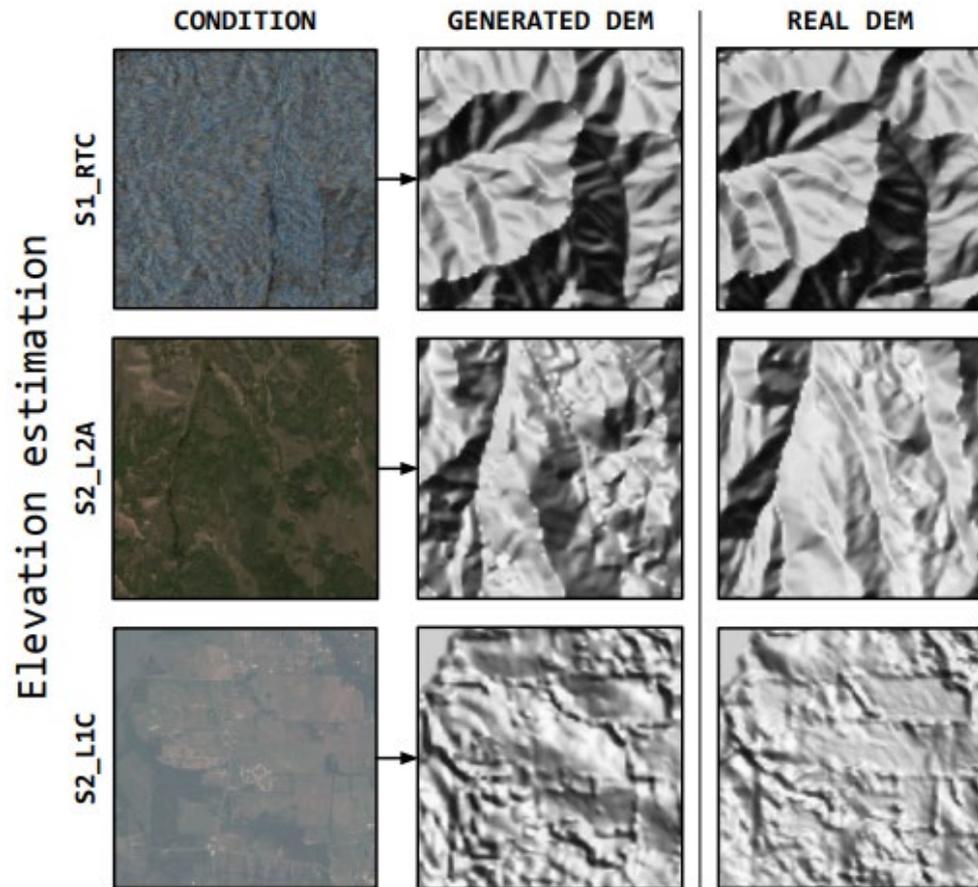
### Reverse (Learned):

Train a neural network  $p_\theta(x)(x_{t-1}|x_t)$  to denoise the image  
Model learns to predict a slightly less noisy image  $x_{t-1}$  from  $x_t$ .  
After T steps, it generates a clean image.

<https://towardsdatascience.com/diffusion-models-made-easy-8414298ce4da/>



# Generative models



Espinosa, Miguel, et al. "COP-GEN-Beta: Unified Generative Modelling of COPernicus Imagery Thumbnails." Proceedings of the Computer Vision and Pattern Recognition Conference. 2025.

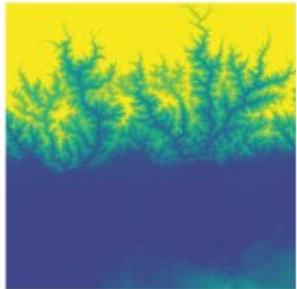
ESA UNCLASSIFIED - For ESA Official Use Only



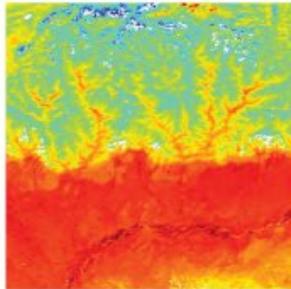
## Mission Critical – Satellite Data is a Distinct Modality in Machine Learning

### Products

Eastern Himalayas,  
English Channel



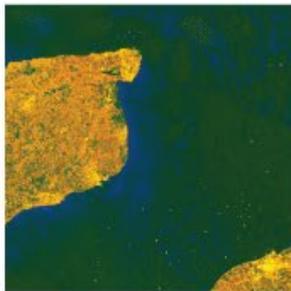
ALOS DEM



MODIS day temp.



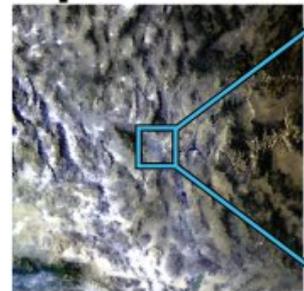
NDVI



SAR (Sentinel 1)

### Spatial resolutions

Las Vegas, Nevada, USA



GOES-18 at 2000m/px



MODIS at 250m/px



Landsat 9 at 30m/px



Sentinel 2 at 10m/px



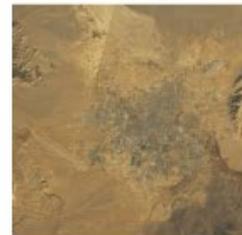
NAIP at 0.6m/px

### Time steps

Las Vegas, Nevada, USA



Dec. 25, 1973



Dec. 3, 1982



Dec. 9, 1993



Dec. 24, 2001



Dec. 23, 2013

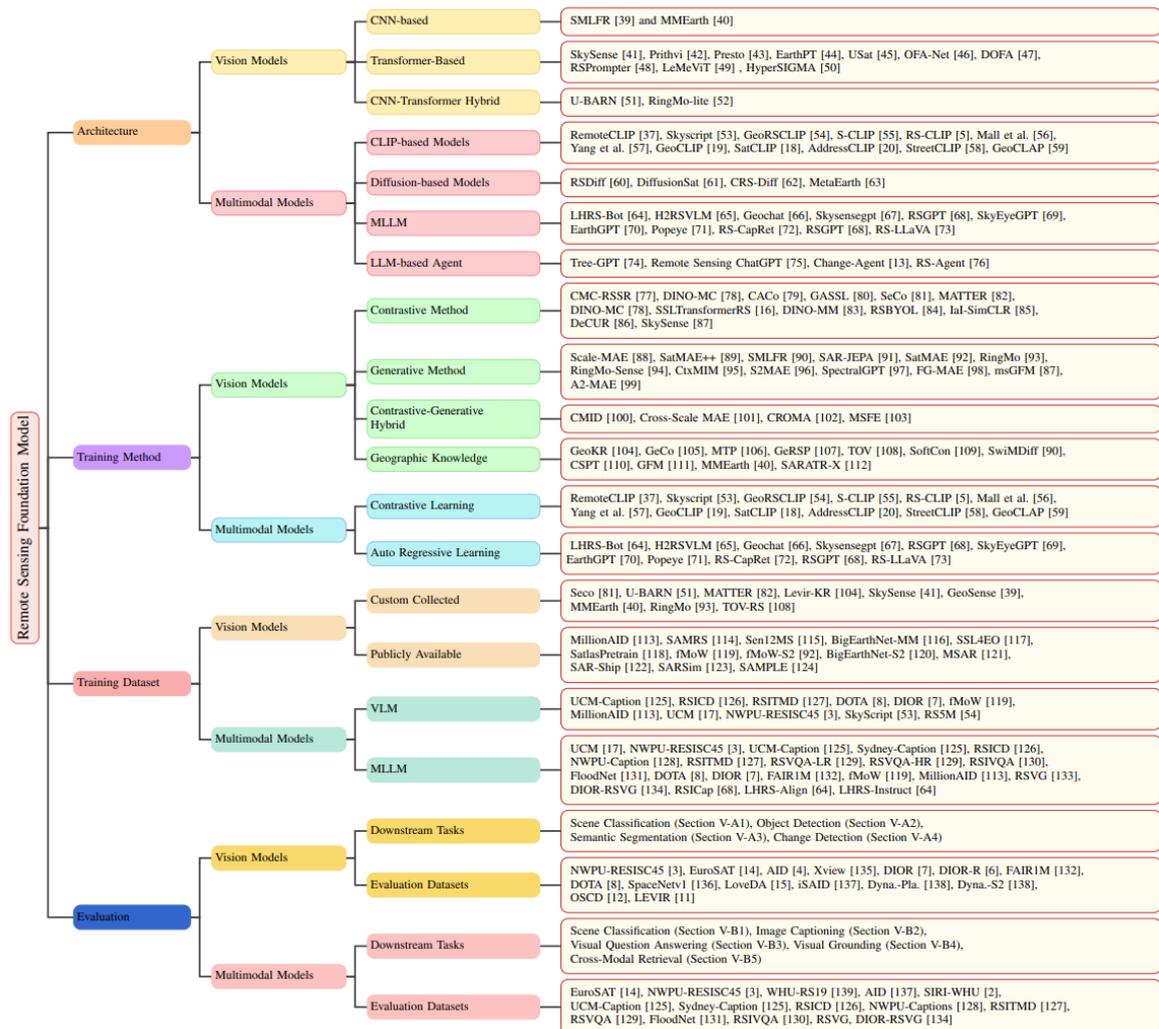


Dec. 28, 2023

# BENCHMARKING GEOSPATIAL FOUNDATION MODELS

---

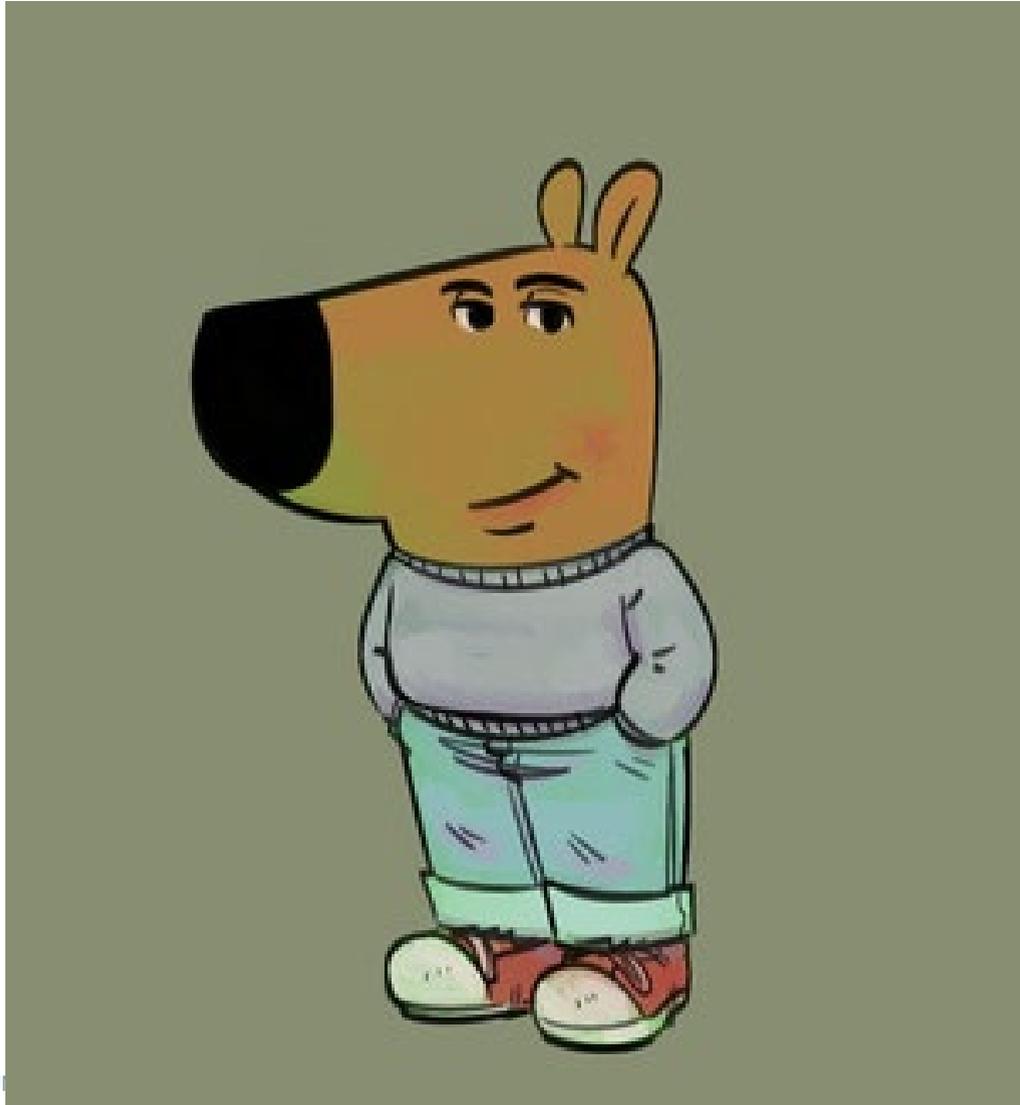
# Which is the best among all available GFM's?



# Which is the best among all available GFM's?



Φ-lab



Let's benchmark them and find out



ESA U



THE EUROPEAN SPACE AGENCY

## 1. Performance Evaluation

## 2. Fairness & Robustness

## 3. Guiding Improvements

### **Research themes:**

Effective benchmarking can enhance the capabilities of GFM's themselves, improving the easiness of adoption and finding the limitations to be considered

# PANGAEA



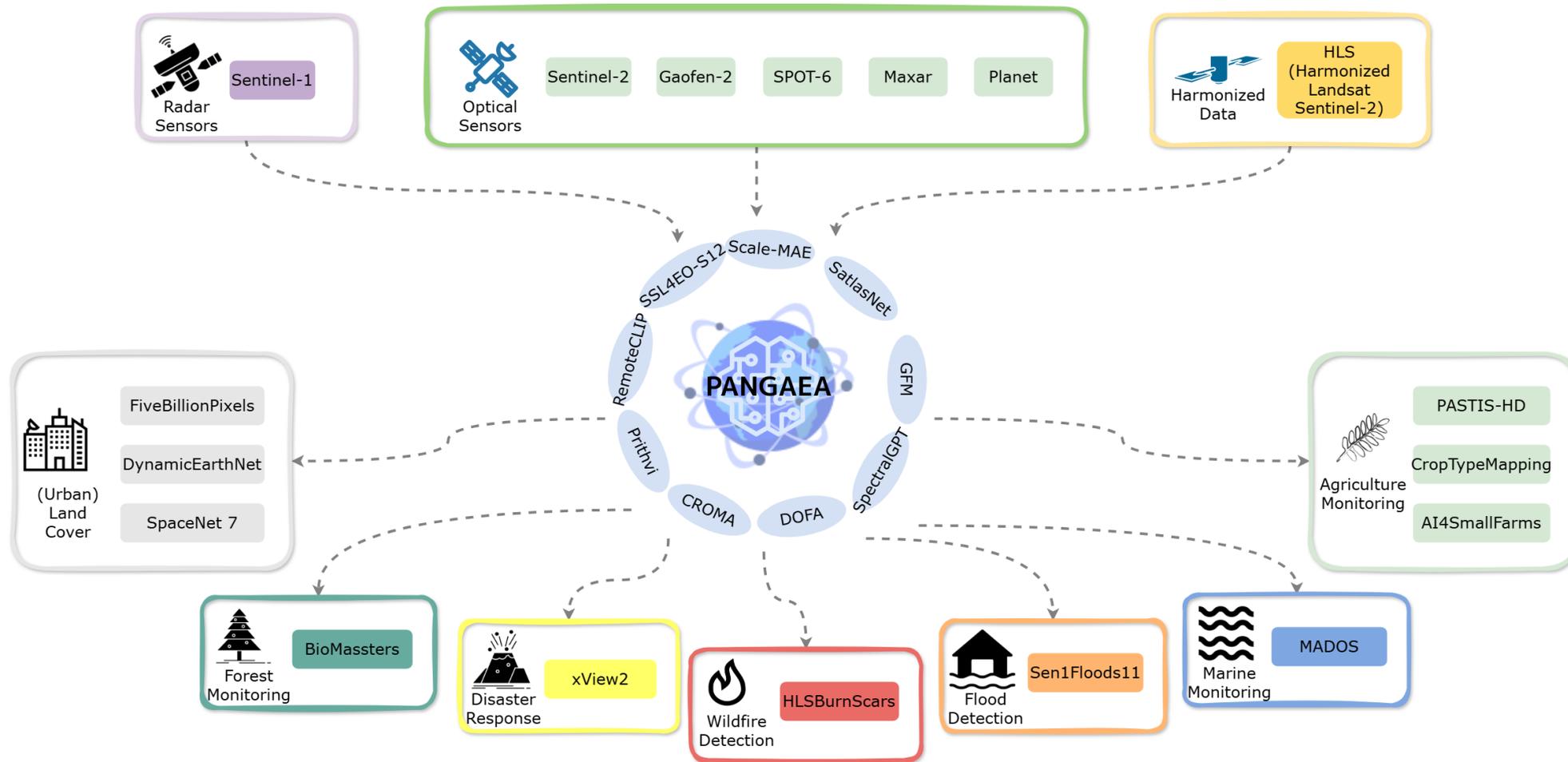
Φ-lab



ESA U



→ THE EUROPEAN SPACE AGENCY



# PANGAEA: the datasets

BurnScars



HLS

MADOS



S-2 L2R

PASTIS



S-2 L2A

Sen1Fl11



S-2 L1C

FBP



Gaofen-2

DynamicEN



Planet

CTM-SS



S-2 L2A

SpaceNet7

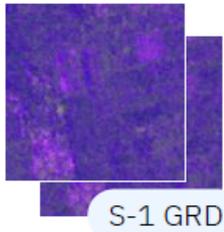


Planet

AI4Farms



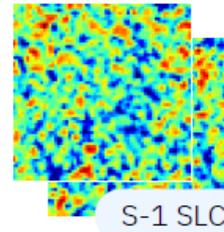
S-2 L2A



S-1 GRD



S-1 GRD



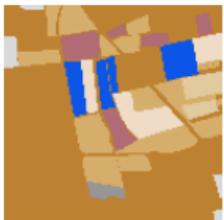
S-1 SLC



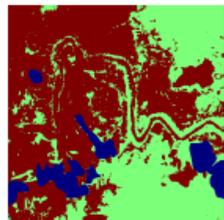
Wildfire



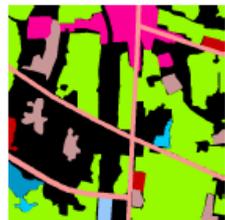
Marine



Agriculture



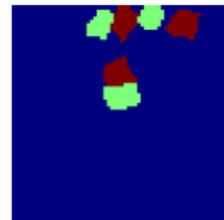
Flood



Land cover



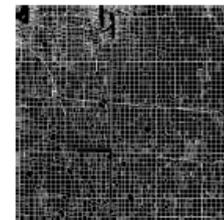
Land cover



Agriculture

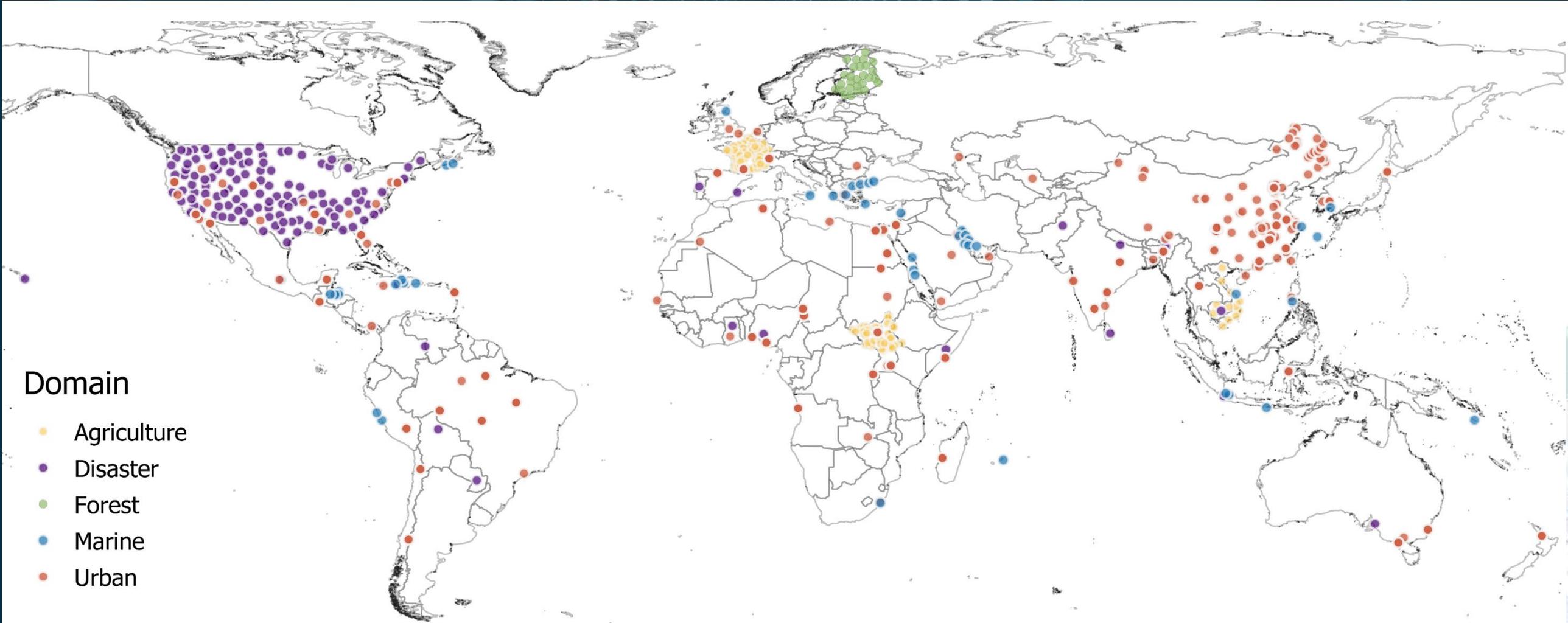


Change det.



Agriculture

# PANGAEA: the datasets



## Domain

- Agriculture
- Disaster
- Forest
- Marine
- Urban

ESA UNCLASSIFIED - For ESA Official Use Only



Table 3: Overview of the pretraining datasets and the number of patches used by the selected GFMs. For Prithvi, the data volume is reported.

Model	Pretraining Images	Patches/Volume
CROMA	Sentinel-1, Sentinel-2	3M
DOFA	Sentinel-1, Sentinel2, Gaofen-2, NAIP, EnMAP	8.08M
GFM-Swin	NAIP, RSD46-WHU, MLRSNet, RESISC45, PatternNet	600K
Prithvi	Harmonized Landsat Sentinel-2 (HLS)	1TB
RemoteCLIP	SEG-4, DET-10, RET-3	165K
SatlasNet	Sentinel-2, NAIP	856K
Scale-MAE	FMoW-RGB	363.6K
SpectralGPT	fMoW-S2, BigEarthNet	1.47M
SSL4EO-S12	Sentinel-1, Sentinel-2	3M

# Some results

Table 5: Performance evaluation of Geospatial Foundation Models across 11 benchmark datasets using 100% of the data. For semantic segmentation and change detection tasks, the mIoU  $\uparrow$  is reported. For regression task, RMSE  $\downarrow$  is reported. #Top2 indicates the number of datasets where the models achieve top-2 performance across all evaluated datasets.

Model	HLS Burns	MADOS	PASTIS	Sen1Floods11	xView2	FBP	DynEarthNet	CropMap	SN7	AI4Farms	BioMassters	#Top2
CROMA	82.42	<b>67.55</b>	32.32	<u>90.89</u>	53.27	51.83	38.29	49.38	59.28	25.65	36.81	2
DOFA	80.63	59.58	30.02	89.37	<u>59.64</u>	43.18	<u>39.29</u>	51.33	61.84	27.07	42.81	2
GFM-Swin	76.90	<u>64.71</u>	21.24	72.60	59.15	67.18	34.09	46.98	60.89	27.19	46.83	1
Prithvi	<u>83.62</u>	49.98	33.93	90.37	49.35	46.81	27.86	43.07	56.54	26.86	39.99	1
RemoteCLIP	76.59	60.00	18.23	74.26	57.41	<b>69.19</b>	31.78	<u>52.05</u>	57.76	25.12	49.79	2
SatlasNet	79.96	55.86	17.51	90.30	52.23	50.97	36.31	46.97	61.88	25.13	41.67	0
Scale-MAE	76.68	57.32	24.55	74.13	<b>60.72</b>	<u>67.19</u>	35.11	25.42	<b>62.96</b>	21.47	47.15	3
SpectralGPT	80.47	57.99	35.44	89.07	48.40	33.42	37.85	46.95	58.86	26.75	<u>36.11</u>	1
S12-MoCo	81.58	51.76	34.49	89.26	51.59	53.02	35.44	48.58	57.64	25.38	40.21	0
S12-DINO	81.72	49.37	<u>36.18</u>	88.61	50.56	51.15	34.81	48.66	56.47	25.62	41.23	1
S12-MAE	81.91	49.90	32.03	87.79	50.44	51.92	34.08	45.8	57.13	24.69	41.07	0
S12-Data2Vec	81.91	44.36	34.32	88.15	51.36	48.82	35.90	<b>54.03</b>	58.23	24.23	41.91	1
UNet Baseline	<b>84.51</b>	54.79	31.60	<b>91.42</b>	58.68	60.47	<b>39.46</b>	47.57	<u>62.09</u>	<b>46.34</b>	<b>35.67</b>	6
ViT Baseline	81.58	48.19	<b>38.53</b>	87.66	57.43	59.32	36.83	44.08	52.57	<u>38.37</u>	38.55	2



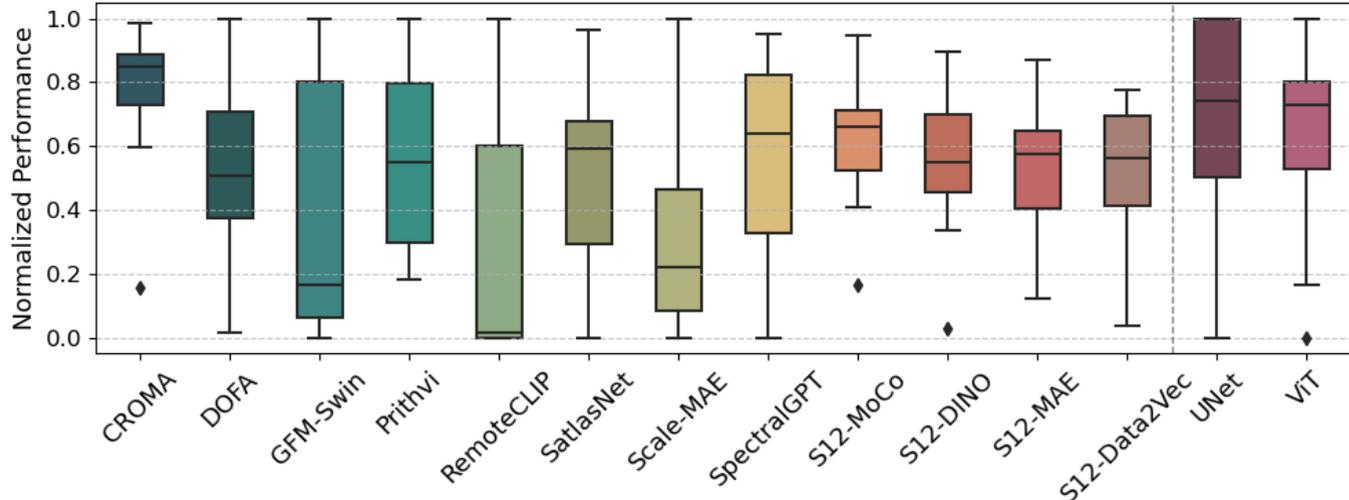
# Some results

HLS	82.4	80.6	76.9	83.6	76.6	80.0	76.7	80.5	81.6	81.7	81.9	81.9	84.5	81.6
S1/S2	60.0	57.6	51.4	54.3	51.1	52.7	45.4	57.4	56.0	55.7	53.9	55.2	56.3	54.6
Maxar	53.3	59.6	59.2	49.4	57.4	52.2	60.7	48.4	51.6	50.6	50.4	51.4	58.7	57.4
Gaofen-2	51.8	43.2	67.2	46.8	69.2	51.0	67.2	33.4	53.0	51.2	51.9	48.8	60.5	59.3
PlanetFusion	38.3	39.3	34.1	27.9	31.8	36.3	35.1	37.8	35.4	34.8	34.1	35.9	39.5	36.8
Planet	59.3	61.8	60.9	56.5	57.8	61.9	63.0	58.9	57.6	56.5	57.1	58.2	62.1	52.6
	CROMA	DOFA	GFM-Swin	Prithvi	RemoteCLIP	SatlasNet	Scale-MAE	SpectralGPT	S12-MoCo	S12-DINO	S12-MAE	S12-Data2Vec	UNet Baseline	ViT Baseline

ESA UNCLASSIFIED - For I

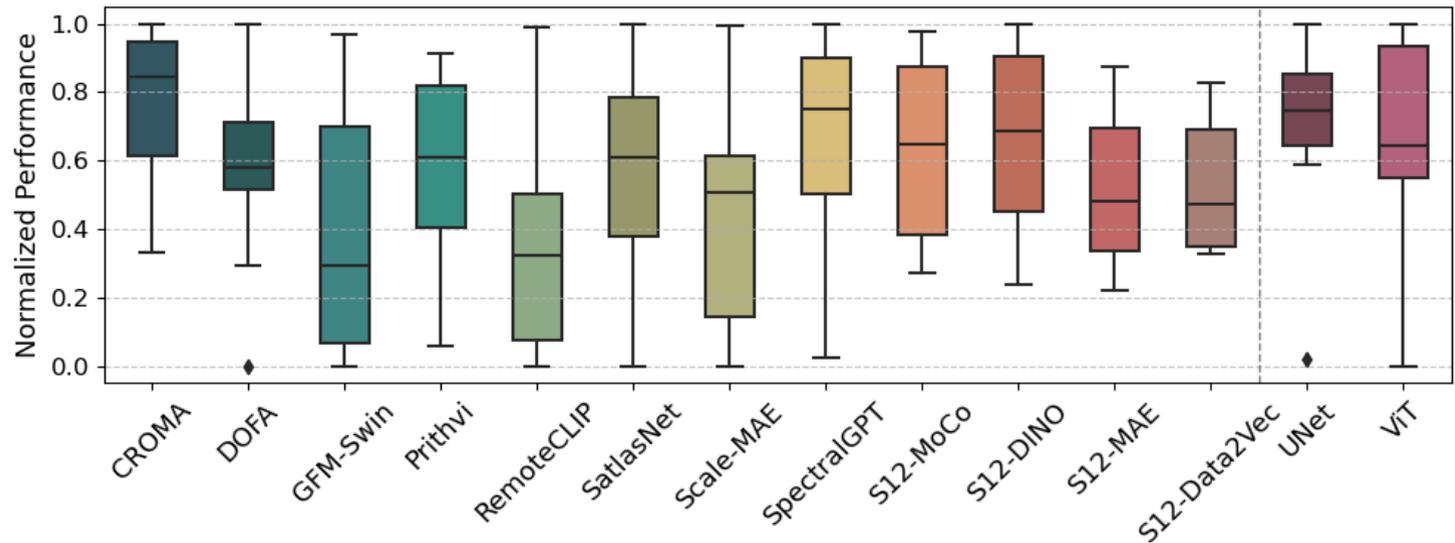


# Some results: about data scarcity



50% labels

10% labels



setup.py Minor update of the team name 7 months ago

README GPL-3.0 license

Tests passing

## PANGAEA: A Global and Inclusive Benchmark for Geospatial Foundation Models

### News

- [23/04/2025] we pushed a new version of the code, fixing different bugs (e.g. commands are working for all the datasets now, metric computation with ignore\_index is fixed, etc...). In the next month, we will provide: all downloadable datasets and models, downloadable stratified subsamples for all the datasets, classification. Stay tuned!
- [22/04/2025] on EarthDay, PANGAEA was officialy adopted to benchmark TerraMind. Read the [news](#) and the [pre-print](#). We will release the benchmarking code in PANGAEA very soon!
- [05/12/2024] the [pre-print](#) is out!

Contributors 12



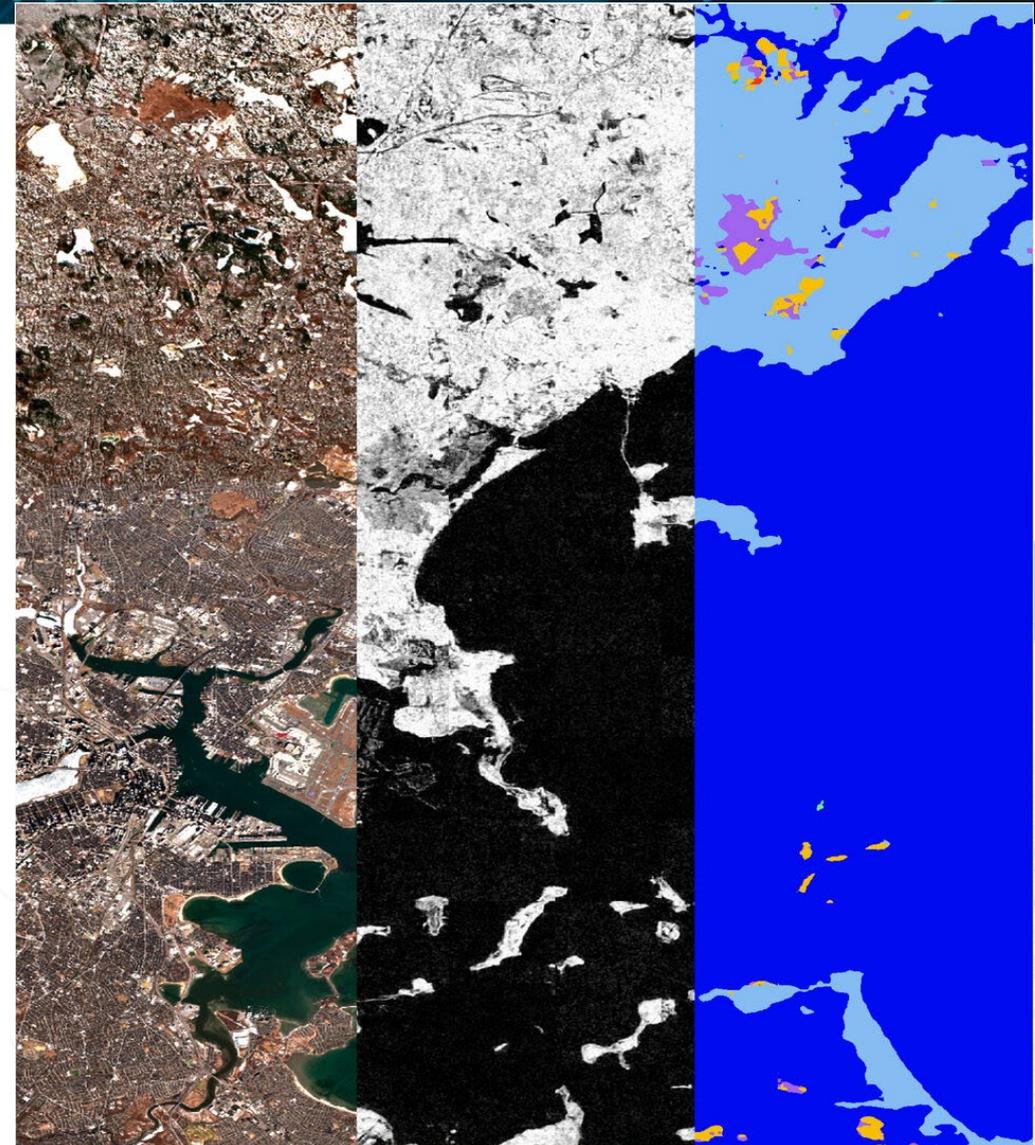
Languages

Python 100.0%

>200 stars

# Some nice achievements

- **AnySat** [CVPR 2025 Highlight]:  
One Earth Observation Model for Many Resolutions, Scales, and Modalities
- **TerraMind** [ICCV 2025]:  
Large-Scale Generative Multimodality for Earth Observation
- Several newly added features  
(e.g. kNN, classification, thematic benchmarks, challenges)



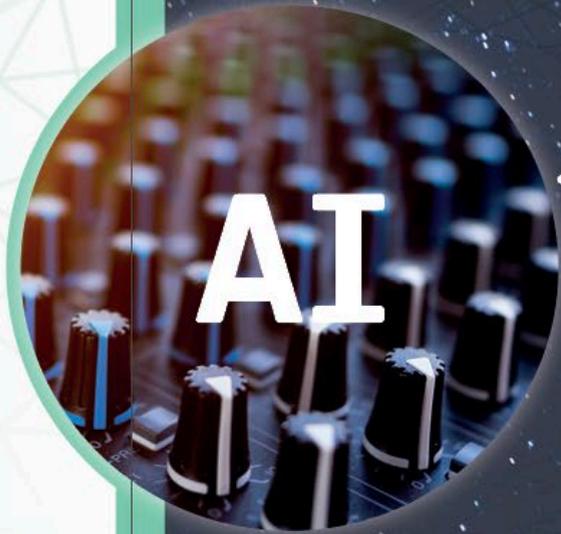
- Geospatial Foundation Models: hype or the new big thing?
- Importance of benchmarking
- Importance of an open-source project



# ESA Φ-lab

---

# The ESA $\Phi$ -lab – Why ?



## from Earth Observation to Earth Action

### From data to actionable information

ESA UNCLASSIFIED - For ESA Official Use Only



$\Phi$ -lab aims to become “the reference” for the transformational innovation and a key influencer (by reputation and authority) in the Earth Observation ecosystem



## Catalyst

- **Attract EO academic and industrial researchers to generate transformative ideas**
- **Exploit fail fast ethos, rapidly prototyping concepts**
- **An informal but rigorous, multi-disciplinary, collaborative environment**
- **Act as facilitator to foster competitiveness growth and entrepreneurial initiatives**
- **Implement investment actions from ESA MSs or in the investors industry**

## Bridge

- **Be the bridge between the European start-ups, academic and industrial researchers, New Space operators, Investors, ICT players, EO world leaders, and ESA**
- **Act as hub stimulating, connecting, and developing a growing ecosystem of talents and capabilities across Europe**



ESA UNCLASSIFIED - For ESA Official Use Only

# Innovation Technologies axis and Applications



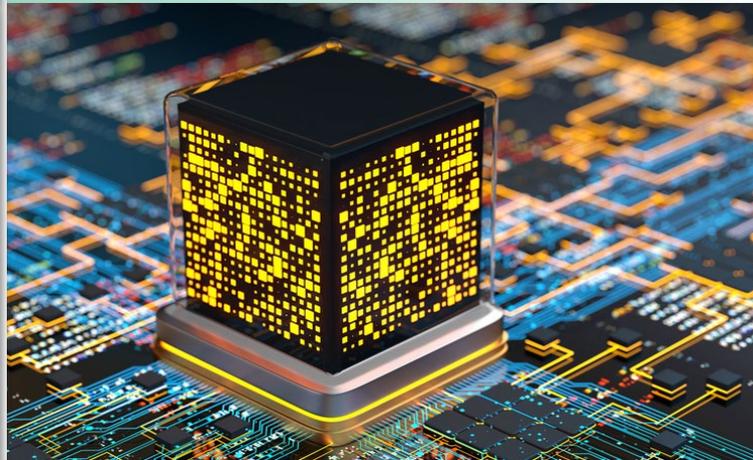
Φ-lab

AXIS I  
Augmented Intelligence



Foundation Models  
Digital Assistant and Twins  
Generative AI  
Decision Intelligence, Agentic AI  
Explainability (xAI)  
Physics-Informed ML  
AI4EO for Climate, Health, and Human

AXIS II  
Innovative Computing Paradigms



On-Board AI  
Quantum Computing  
Hybrid HPC Computing  
Neuromorphic  
Biocomputing and others

AXIS III  
Innovative Computing Paradigms



Cognitive Space  
VR/AR Immersive Visualisation  
Web 3.0  
IOT  
Distributed Ledgers/Blockchain

ESA UNCL



# Collaboration opportunities at $\Phi$ -lab



$\Phi$ -lab



Shared mutual interest

Join the  $\Phi$ -lab to explore disruptive ideas

as a Visiting Researcher (industry, academia),  
Visiting Professor, Research Fellow, PhD, YGT, etc.

Funded

1.  $\Phi$ -lab's [Invitation To Tender](#) on ESA-STARS
  - Foundation Models, Generative AI, QC4EO, Edge computing, Web 3.0, etc..
2. [InCubed](#) : partnership development of commercial products or services
3. [Open Space Innovation Platform](#) : co-funded research or researchers
4. [EO Science4Society](#) : no SOW, 100/200K, 6/18 months
5. ESA Technology Programmes like [GSTP](#) and [TDE](#)

Visiting Professor

Visiting Researchers (Industrial and Scientific)

ESA Research Fellowships

ESA Co-funded PhD

ESA Early Graduate Traineeships (EGT), Internships, National trainee



# Join ESA $\Phi$ -lab through CIN



$\Phi$ -lab

The **Collaborative Innovation Network (CIN)** by **ESA  $\Phi$ -lab**, provides to leading researchers and University Professors the **opportunity** to join ESA  $\Phi$ -lab and be actively involved in **accelerating the future of Earth Observation with ESA**.

The ESA  $\Phi$ -lab CIN aims to:

- Establish a global network through which researchers and innovators can **JOIN ESA  $\Phi$ -LAB**
- **Promote knowledge** sharing and develop groundbreaking EO solutions

Check the open calls



Follow CIN on LinkedIn



# What to go deeper into $\Phi$ -lab disruptive innovation?

Don't miss our spotlight session at Big Data from Space  
(Day 2 – October 2nd at 17:50)



