



UiT The Arctic University of Norway

NORA summer school on multi-modal learning

Multi-Modal Learning: Beyond Vision and Language

Challenges in Remote Sensing

Rwiddhi Chakraborty

UiT Machine Learning Group and Visual Intelligence

Schedule Today

- 09 - 12: Beyond Vision-Language – I and II
- 12 – 13: Lunch
- 13 – 16: Group Project

In this talk

Domain
Challenges

Remote Sensing

Pre-training

Zero-Shot

In this talk

Domain
Challenges

Remote Sensing

Pre-training

Zero-Shot

Beyond Vision and Language

Pre-trained vision-language models lack fine-grained understanding

Unsuitable for domains where specificity is key – e.g. remote sensing, biology, etc.

Text as a data mode has its own issues

Beyond Vision and Language

Real world datasets exist in multiple views and modalities

Only vision-language as a pre-training objective is too restrictive

We want to leverage similar ideas for diverse tasks

In this talk

Domain
Challenges

Remote Sensing

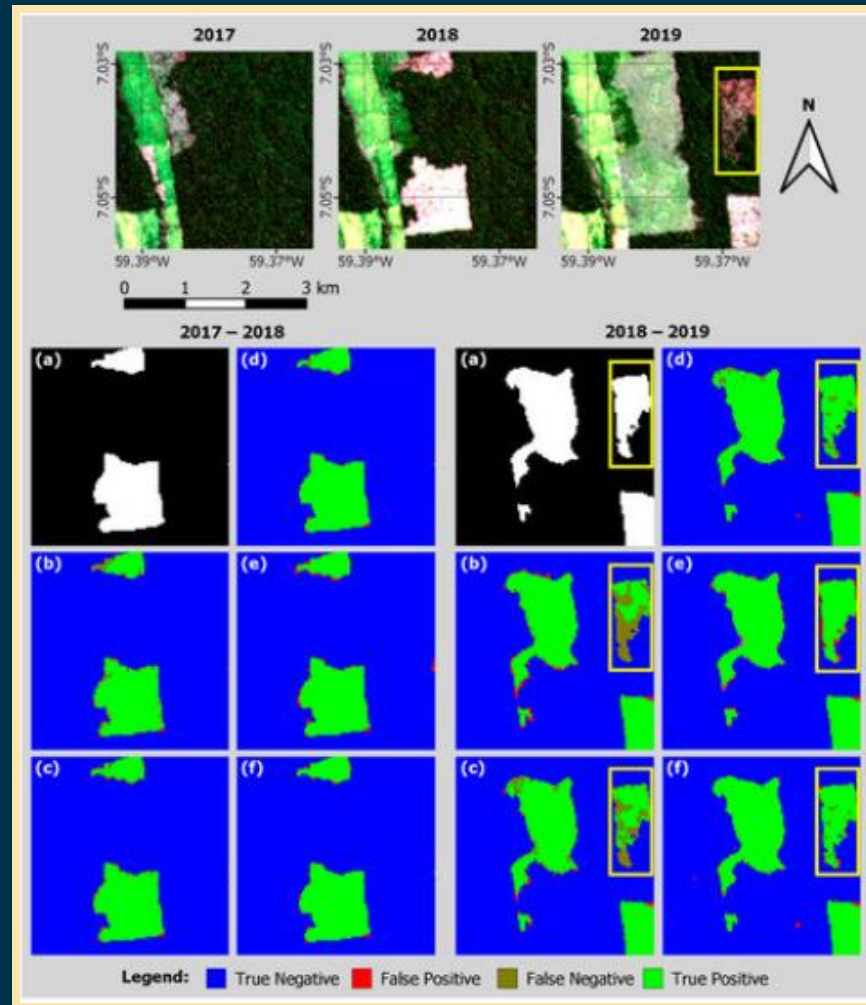
Pre-training

Zero-Shot

Monitoring Deforestation in the Amazon



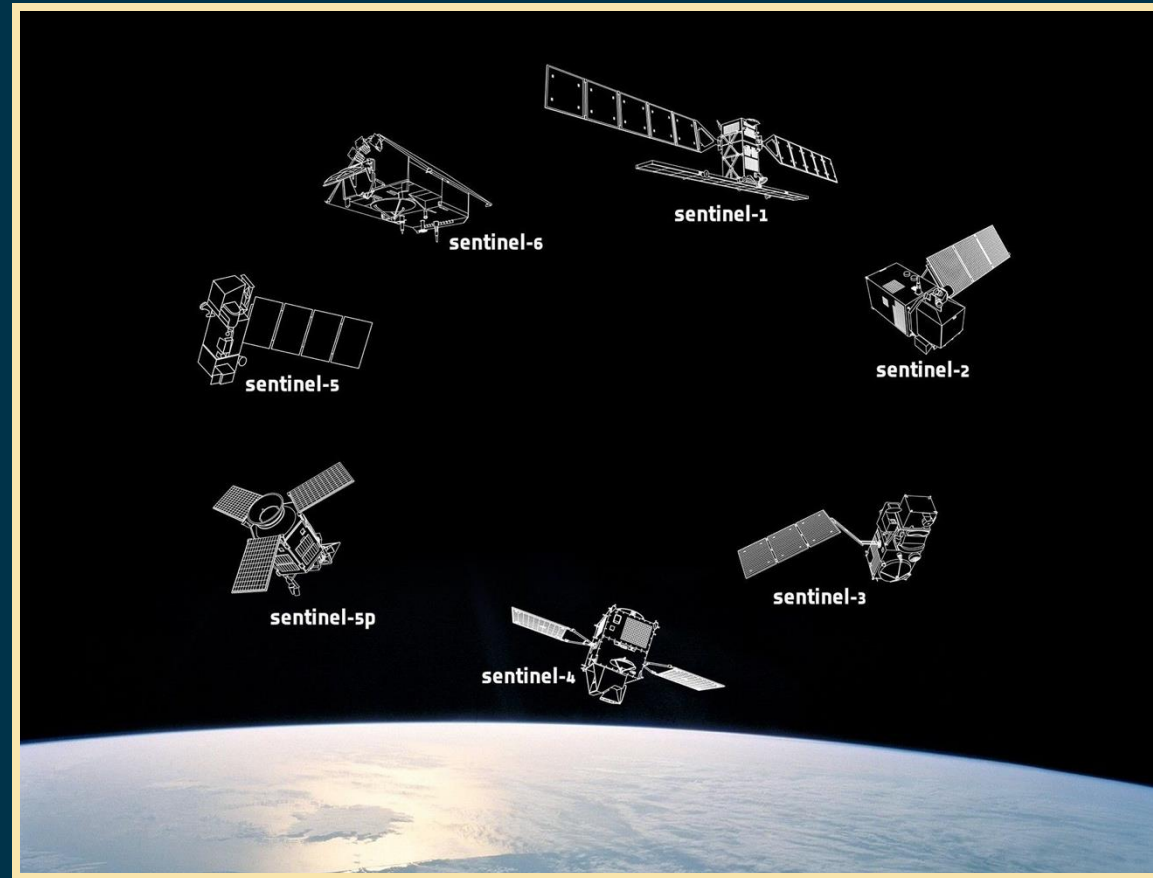
Monitoring Deforestation in the Amazon



Remote Sensing: The Sentinels

Part of Europe's
Copernicus program

Monitor changes on the
Earth's surface



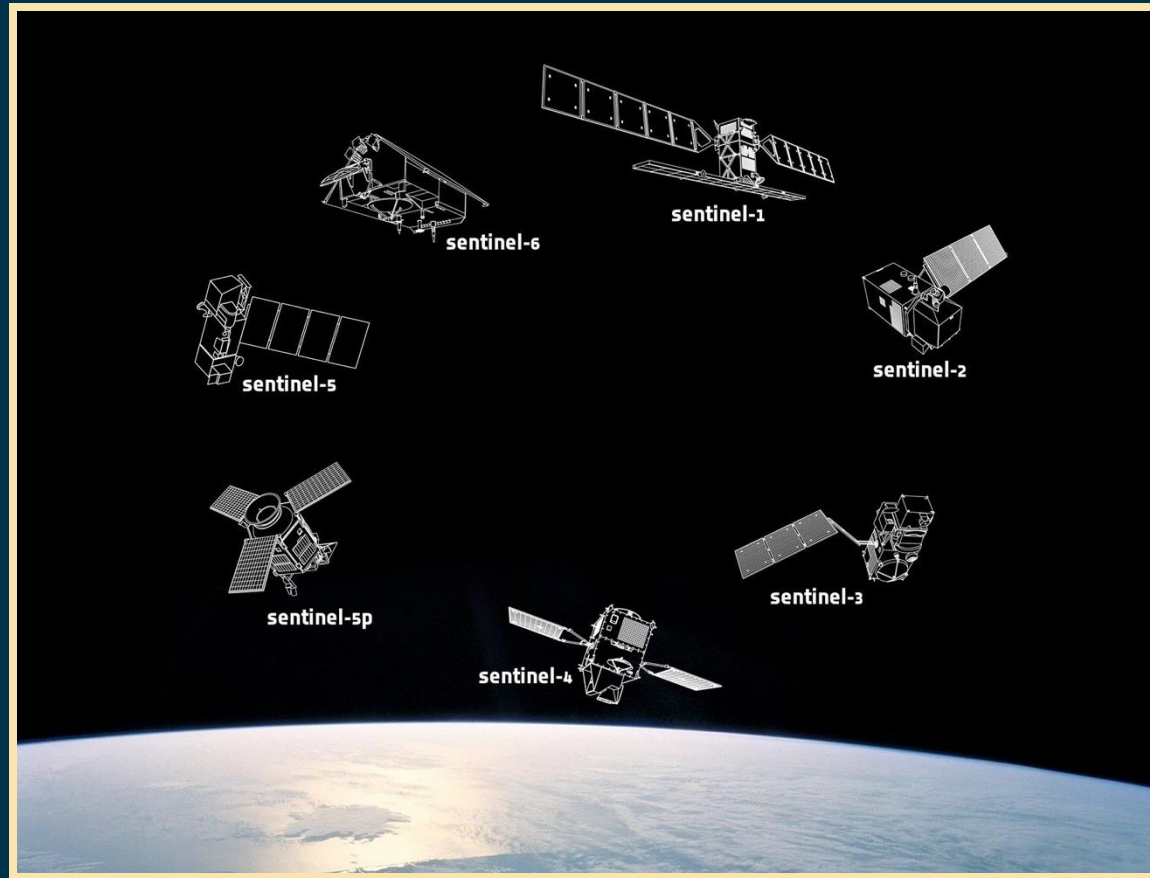
Multispectral optical
(MO) imagery

Synthetic Aperture
Radar (SAR) data

Remote Sensing: The Sentinels

MO imagery reveals
material composition of
objects (440-2200nm)

SAR reveals geometry,
roughness, and
electrical properties of
objects (5.5cm)



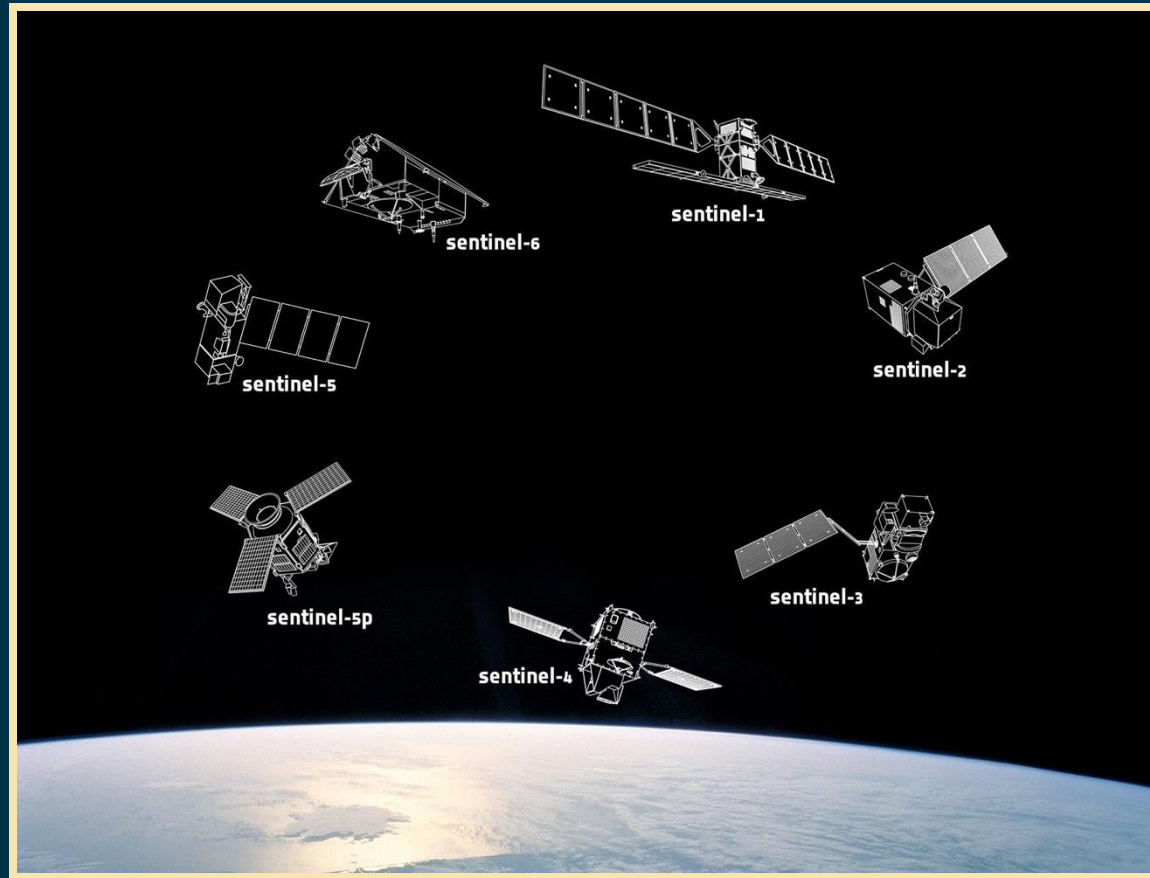
Agriculture, Climate
Change, Forestry....

Cornerstone of Earth
Observation

Remote Sensing: The Sentinels

Multiple views of the
same features on
ground

Spatially aligned

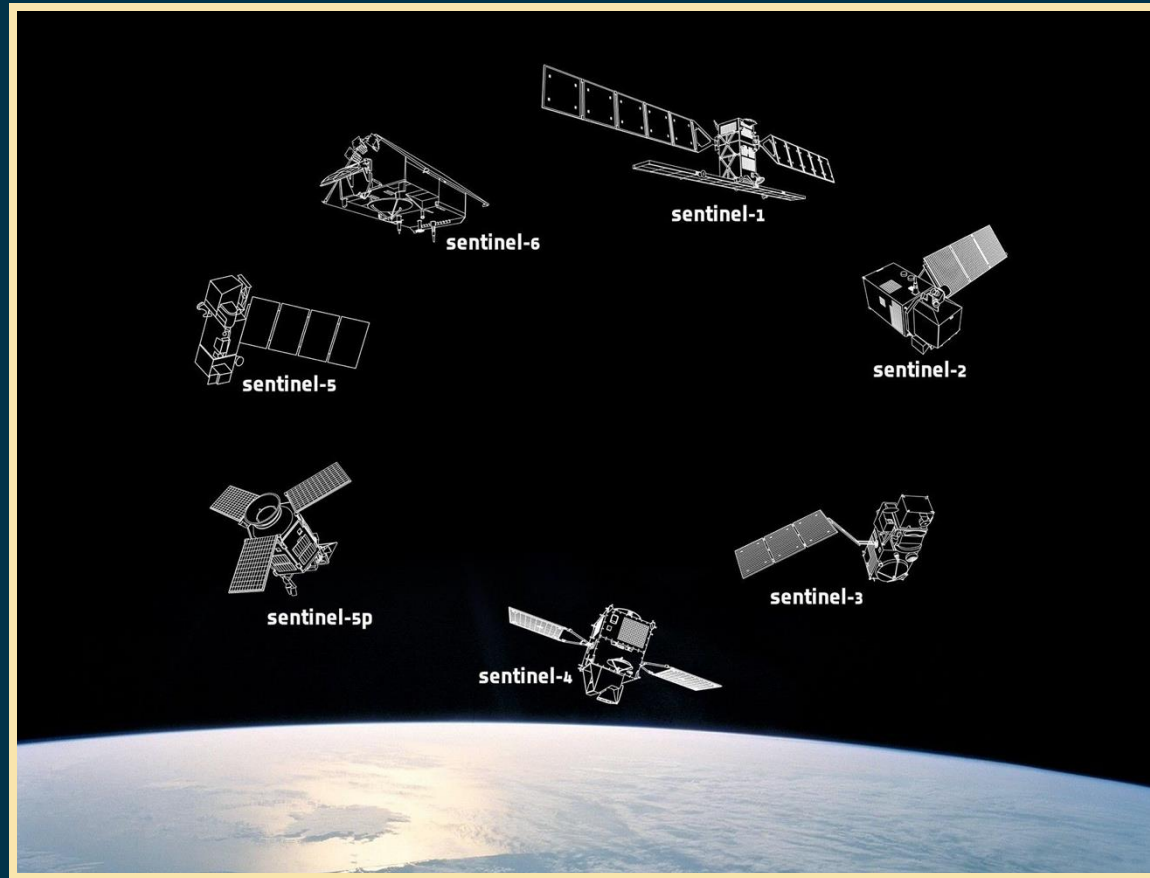


Can self-supervised
learning help?

Remote Sensing: The Sentinels

Multiple views of the
same features on
ground

Spatially aligned



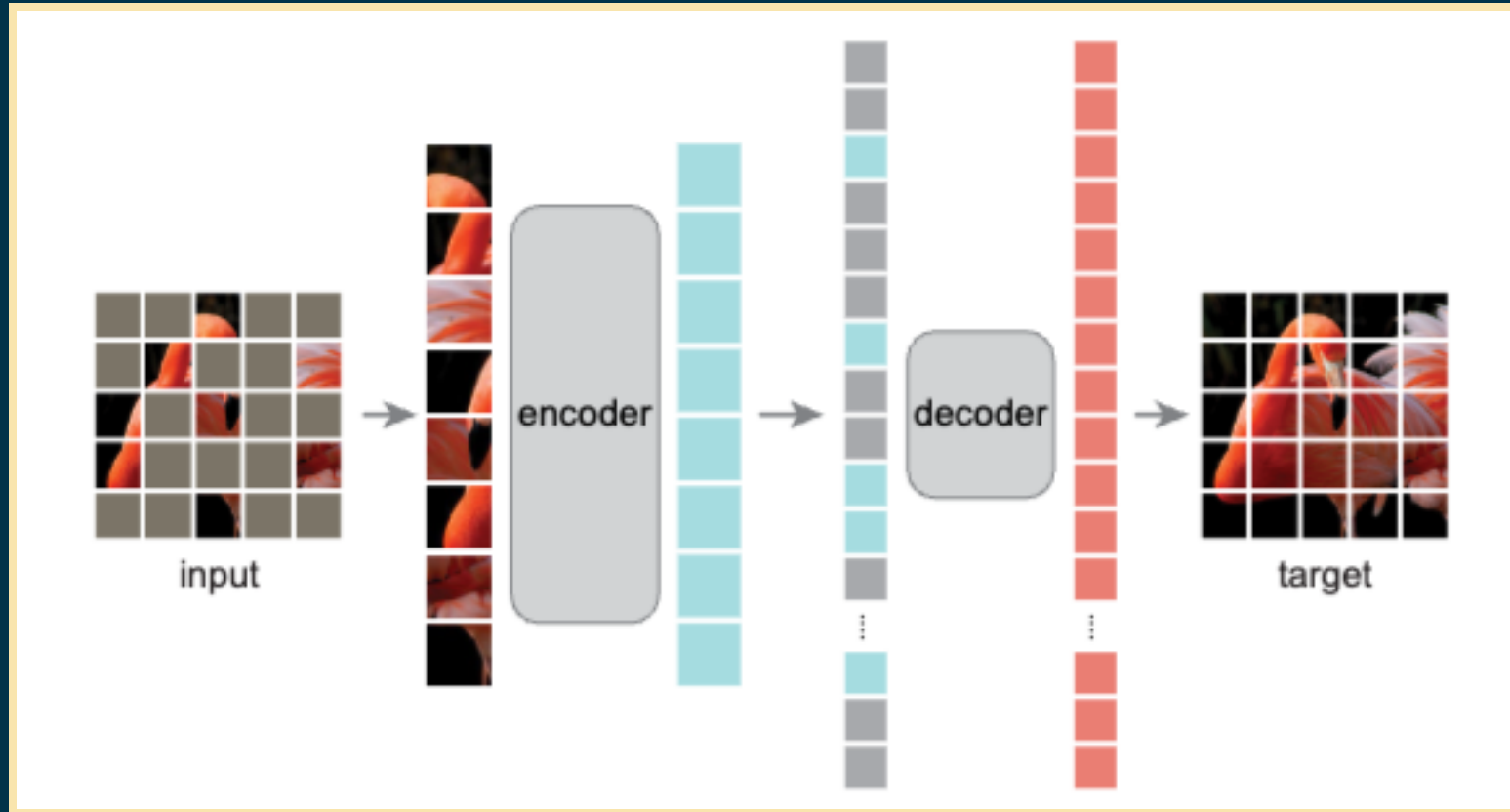
Can self-supervised
learning help?

Yes

In this talk



The Masked Autoencoder



MAE in Remote Sensing?

Good starting point

MAE in Remote Sensing?

Good starting point

BUT

MAE in Remote Sensing?

Good starting point

BUT

Does *not* leverage scale

Absolute/Relative positional encoding

MAE in Remote Sensing?

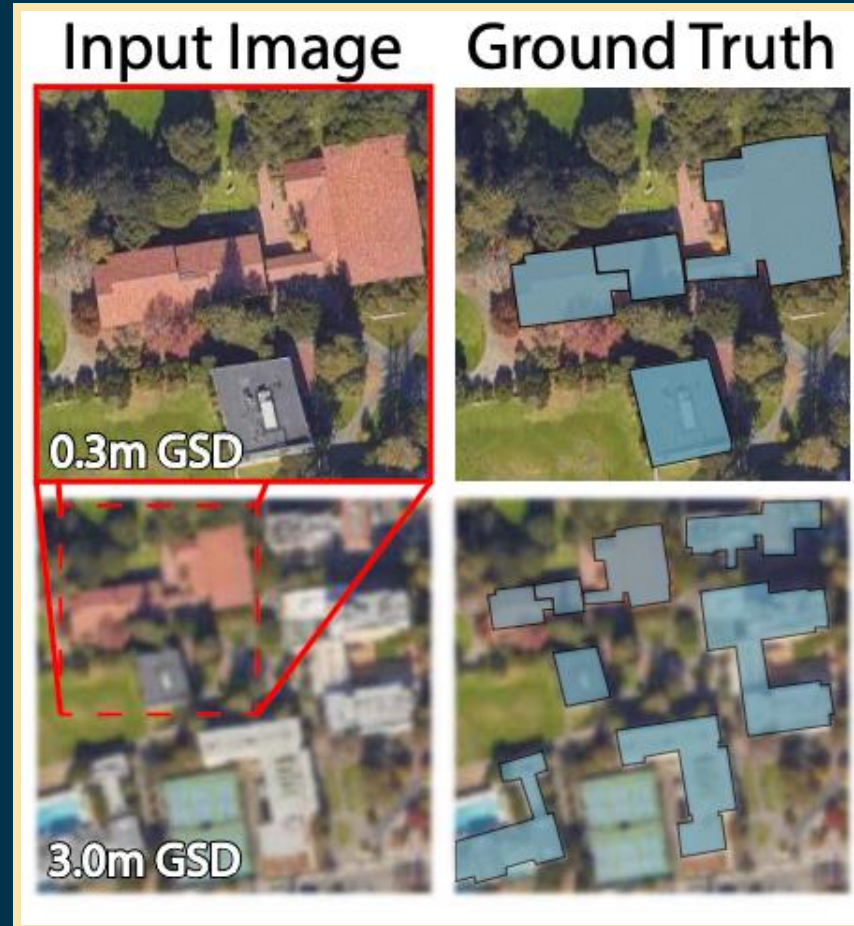
Good starting point

BUT

Does *not* leverage scale

Absolute/Relative positional encoding

MAE in Remote Sensing?



MAE in Remote Sensing?

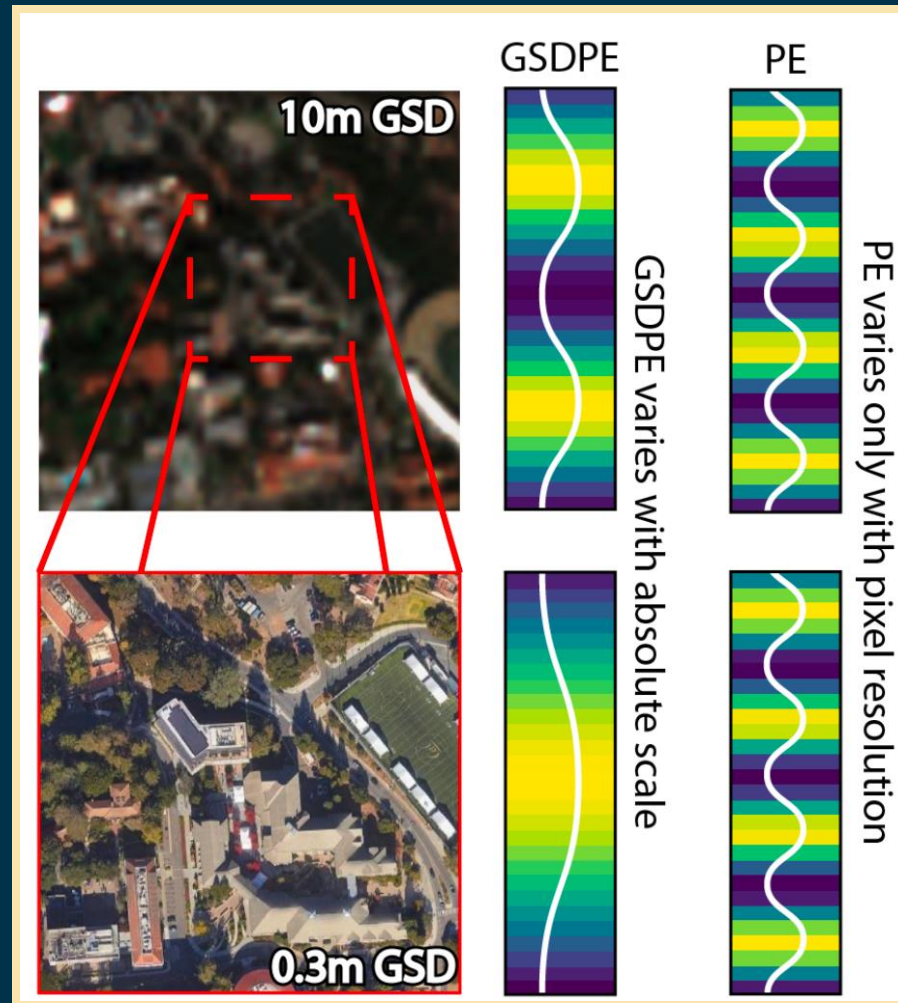
Good starting point

BUT

Does *not* leverage scale

Absolute/Relative positional encoding

MAE in Remote Sensing?

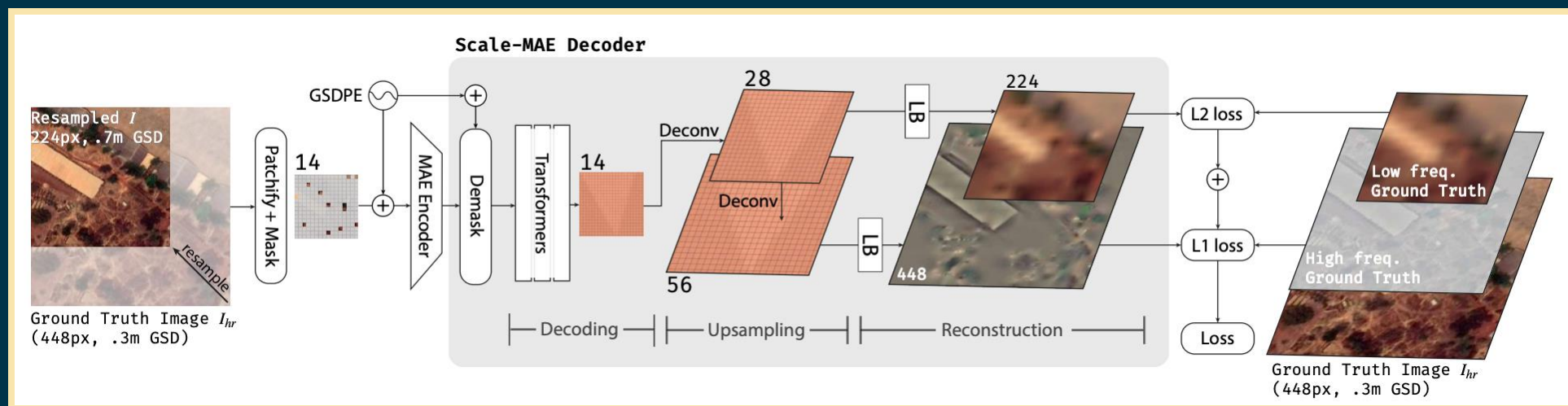


MAE in Remote Sensing?

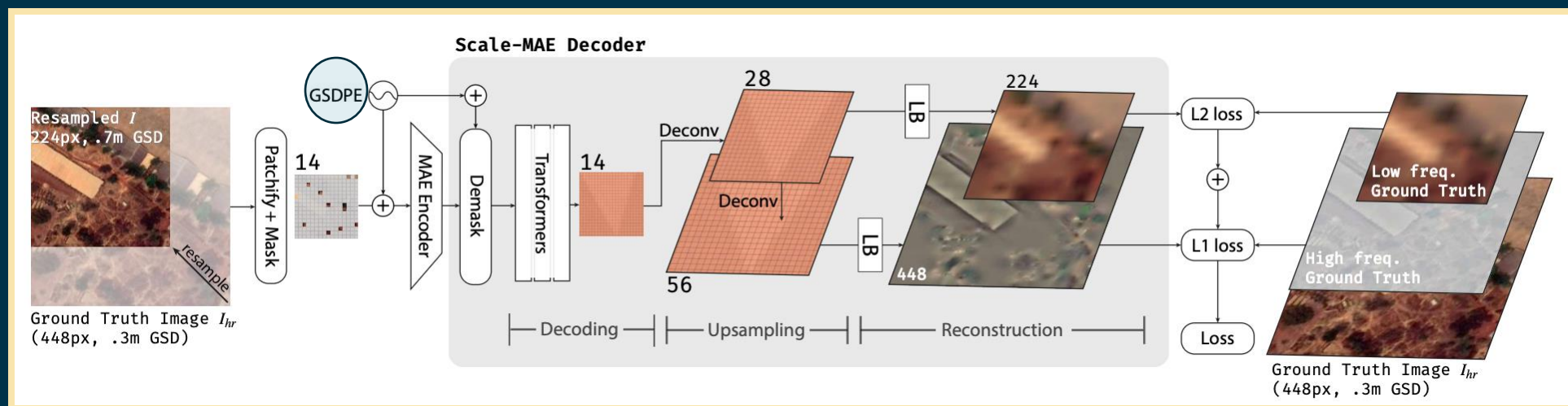
In essence, we need a decoder that has:

- (i) Scale-aware positional encoding
- (ii) Multi-resolution sampling

ScaleMAE for Remote Sensing



ScaleMAE for Remote Sensing



ScaleMAE for Remote Sensing

$$\begin{aligned}v_x(pos, 2i) &= \sin \frac{pos}{10000^{\frac{2i}{D}}} \\v_y(pos, 2i + 1) &= \cos \frac{pos}{10000^{\frac{2i}{D}}}\end{aligned}$$

Before

ScaleMAE for Remote Sensing

$$\begin{aligned}v_x(pos, 2i) &= \sin \frac{pos}{10000^{\frac{2i}{D}}} \\v_y(pos, 2i + 1) &= \cos \frac{pos}{10000^{\frac{2i}{D}}}\end{aligned}$$

Before

$$\begin{aligned}v_{gsd,x}(pos, 2i) &= \sin \frac{g}{G} \frac{pos}{10000^{\frac{2i}{D}}} \\v_{gsd,y}(pos, 2i + 1) &= \cos \frac{g}{G} \frac{pos}{10000^{\frac{2i}{D}}}\end{aligned}$$

Now

ScaleMAE for Remote Sensing

$$\begin{aligned}v_x(pos, 2i) &= \sin \frac{pos}{10000^{\frac{2i}{D}}} \\v_y(pos, 2i + 1) &= \cos \frac{pos}{10000^{\frac{2i}{D}}}\end{aligned}$$

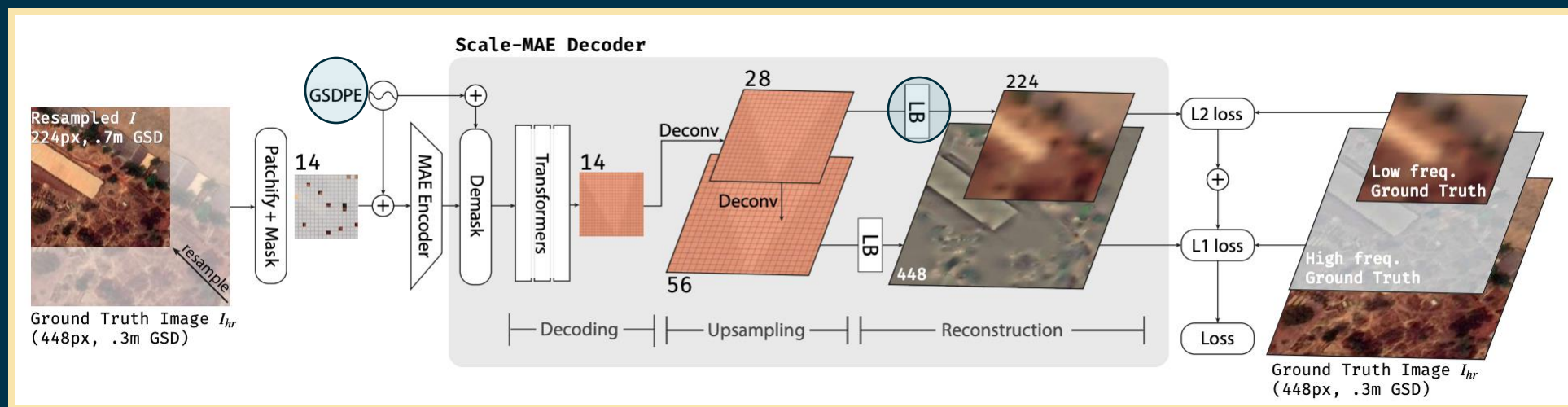
Before

$$\begin{aligned}v_{gsd,x}(pos, 2i) &= \sin \left(\frac{g}{G} \right) \frac{pos}{10000^{\frac{2i}{D}}} \\v_{gsd,y}(pos, 2i + 1) &= \cos \frac{g}{G} \frac{pos}{10000^{\frac{2i}{D}}}\end{aligned}$$

Now

An object at a finer resolution has more pixels representing it
The same object at a coarser resolution must then map to fewer pixels

ScaleMAE for Remote Sensing



ScaleMAE for Remote Sensing

The Laplacian Block (LB) helps sample at a specified frequency (resolution)

Reconstruct low res images for low frequencies, high res images for high frequencies

Final image reconstruction combines these two resolutions

ScaleMAE for Remote Sensing

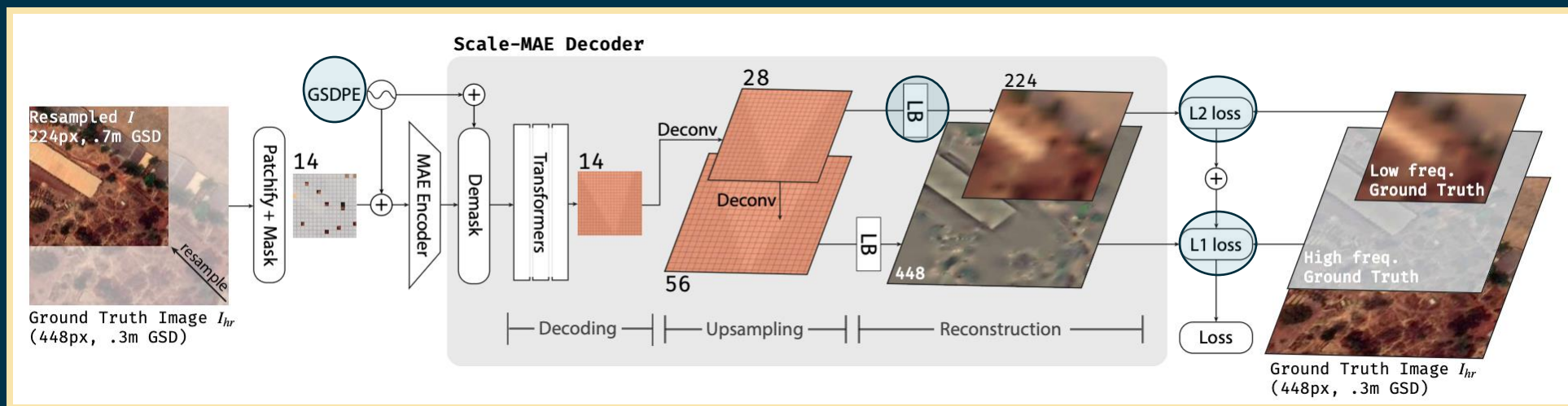
The Laplacian Block (LB) helps sample at a specified frequency (resolution)

Reconstruct low res images for low frequencies, high res images for high frequencies

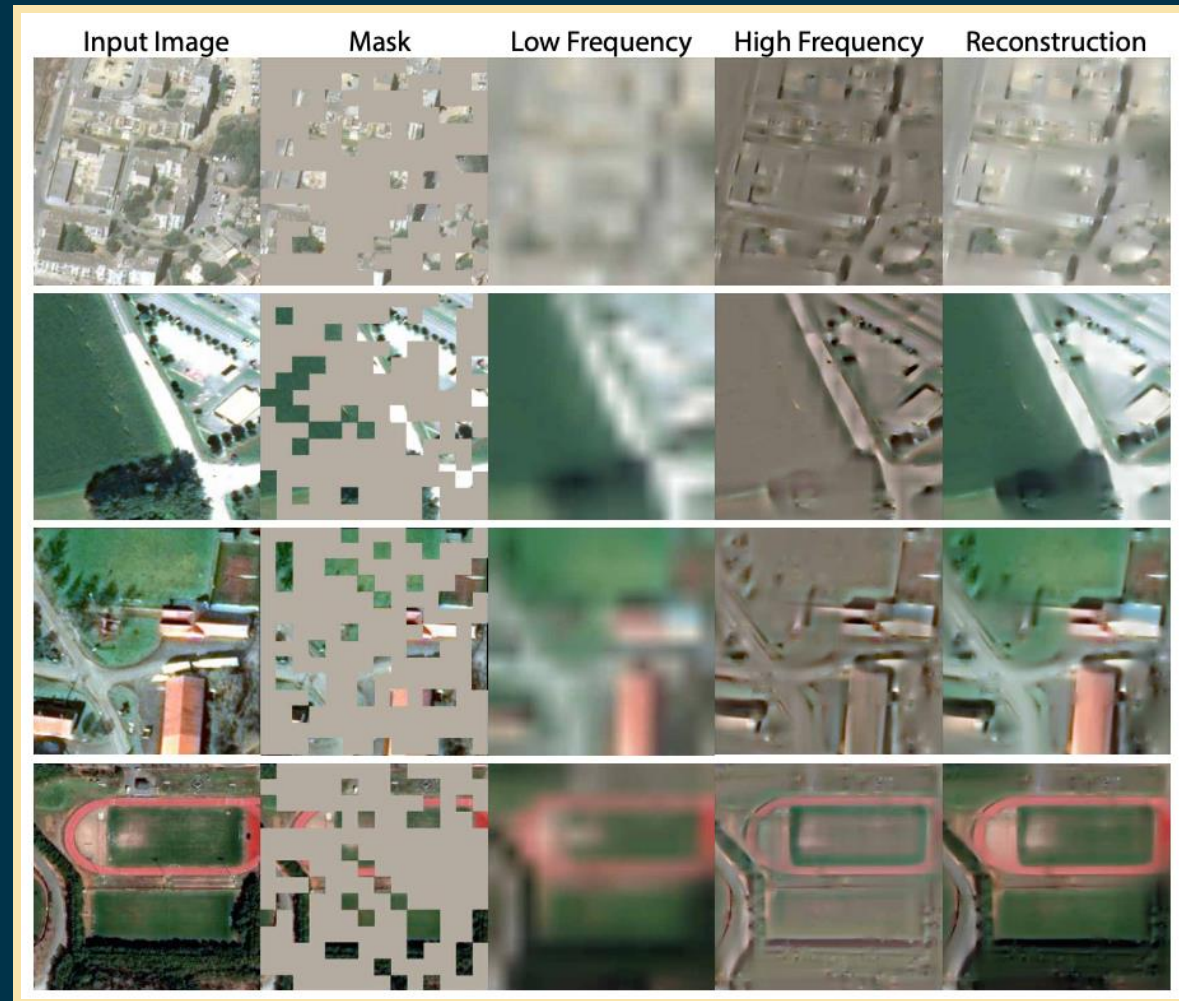
Final image reconstruction combines these two resolutions

Multi-scale reconstruction like this leads to the term Laplacian *pyramid*

ScaleMAE for Remote Sensing



ScaleMAE for Remote Sensing



In this talk



Contrastive ideas as well?

Masked reconstruction algorithms are great:

- (i) Scalable
- (ii) No need for paired data

BUT

Contrastive ideas as well?

Masked reconstruction algorithms are great:

- (i) Scalable
- (ii) No need for paired data

BUT

Significant fine-tuning required

Contrastive ideas as well?

Masked reconstruction algorithms are great:

- (i) Scalable
- (ii) No need for paired data

BUT

Significant fine-tuning required

Contrastive algorithms are great:

- (i) Rich view information
- (ii) Great for downstream tasks

BUT

Contrastive ideas as well?

Masked reconstruction algorithms are great:

- (i) Scalable
- (ii) No need for paired data

BUT

Significant fine-tuning required

Contrastive algorithms are great:

- (i) Rich view information
- (ii) Great for downstream tasks

BUT

Data and compute hungry
Sensitive to paired info

Contrastive ideas as well?

Masked reconstruction algorithms are great:

- (i) Scalable
- (ii) No need for paired data

BUT

Significant fine-tuning required



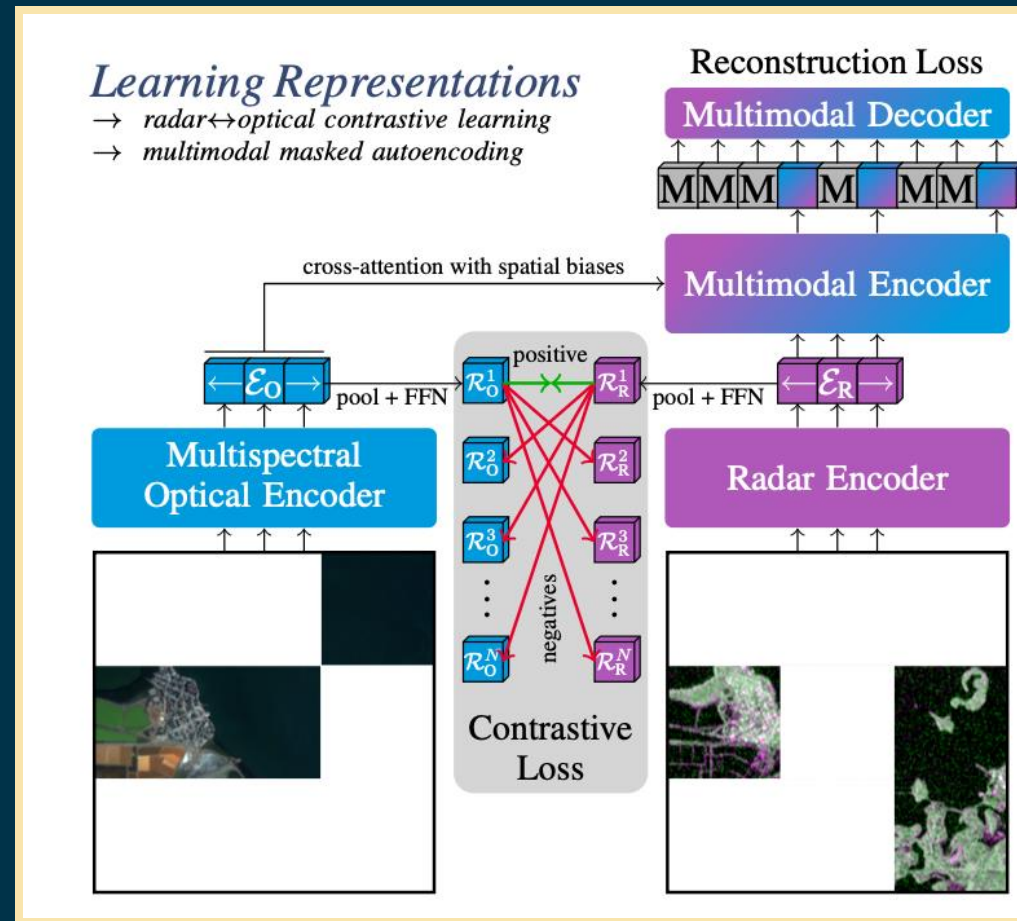
Contrastive algorithms are great:

- (i) Rich view information
- (ii) Great for downstream tasks

BUT

Data and compute hungry
Sensitive to paired info

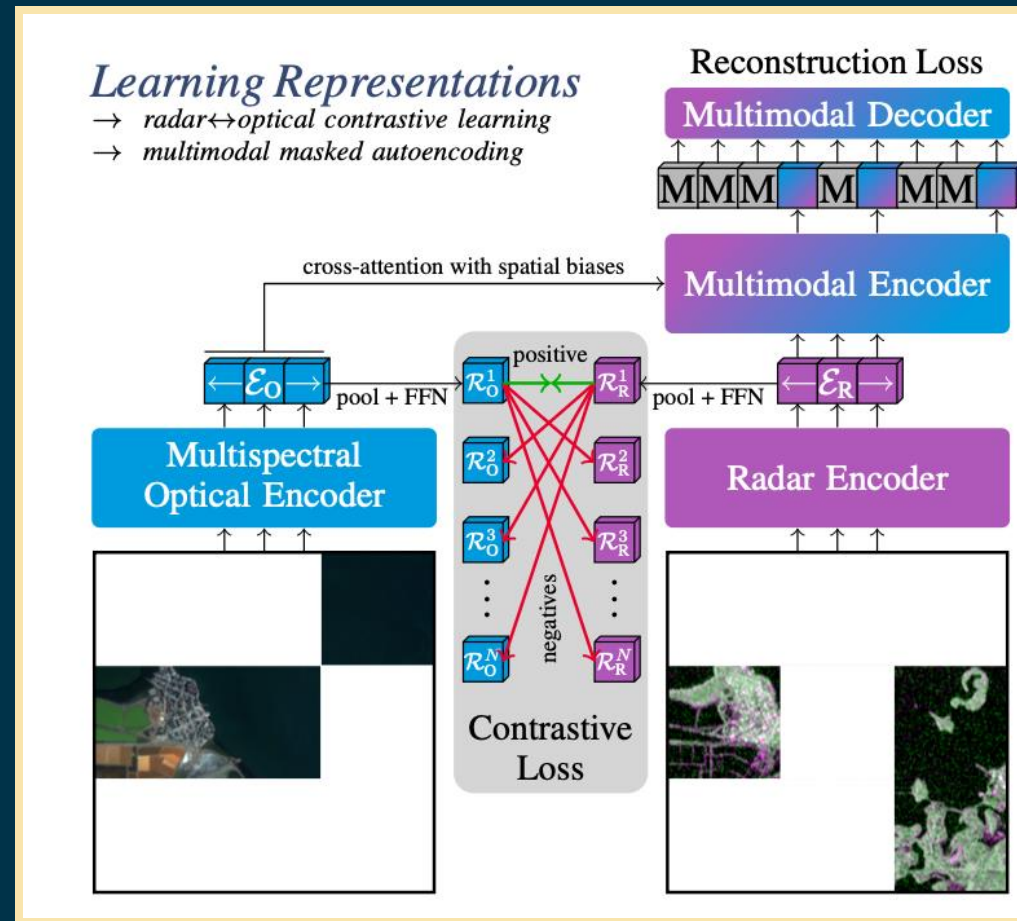
Contrastive + Masked Reconstruction!



Contrastive + Masked Reconstruction!

Optical Data:
12-channel MO from Sentinel-2

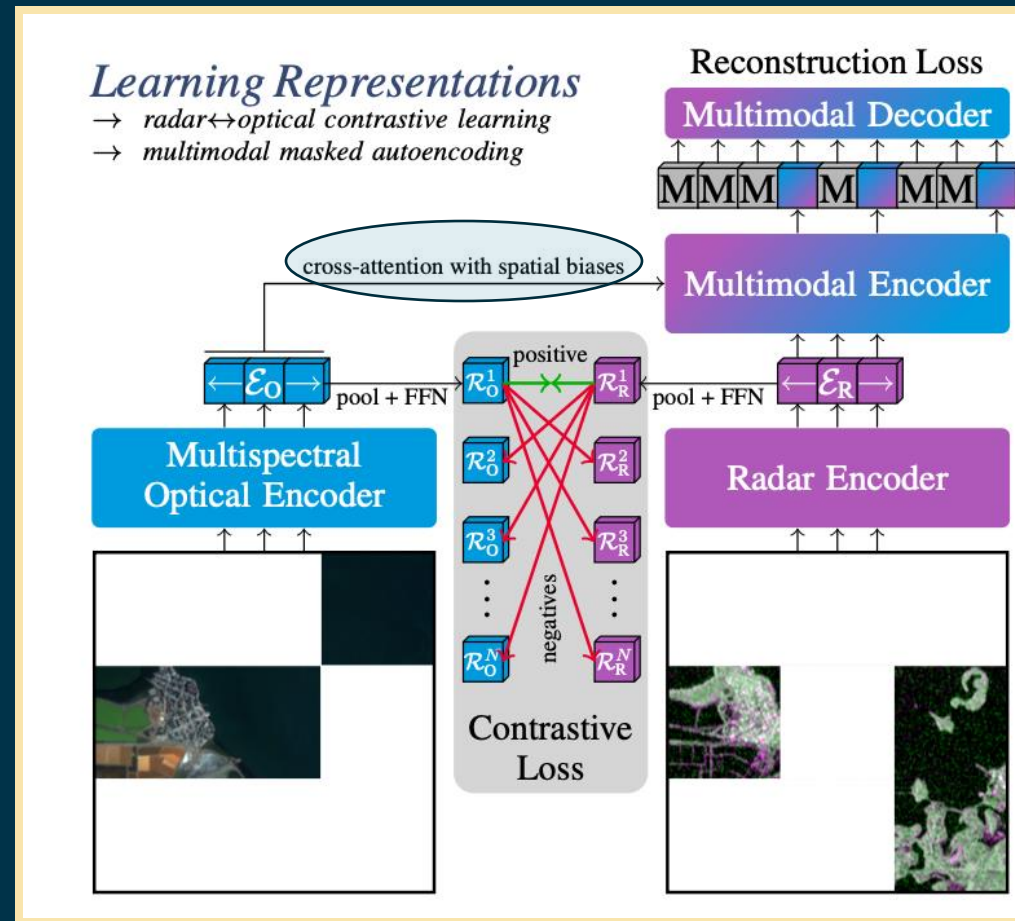
Radar data:
2-channel backscatter from Sentinel-1



Contrastive + Masked Reconstruction!

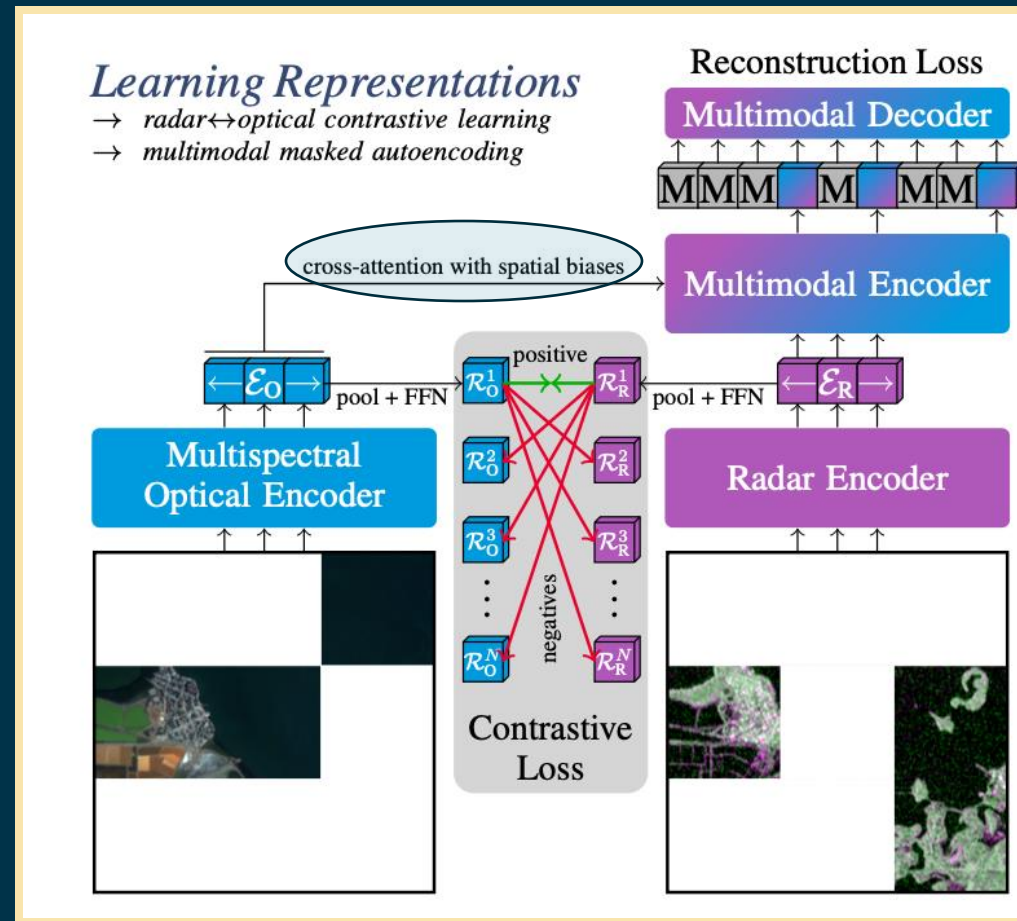
Optical Data:
12-channel MO from Sentinel-2

Radar data:
2-channel backscatter from Sentinel-1



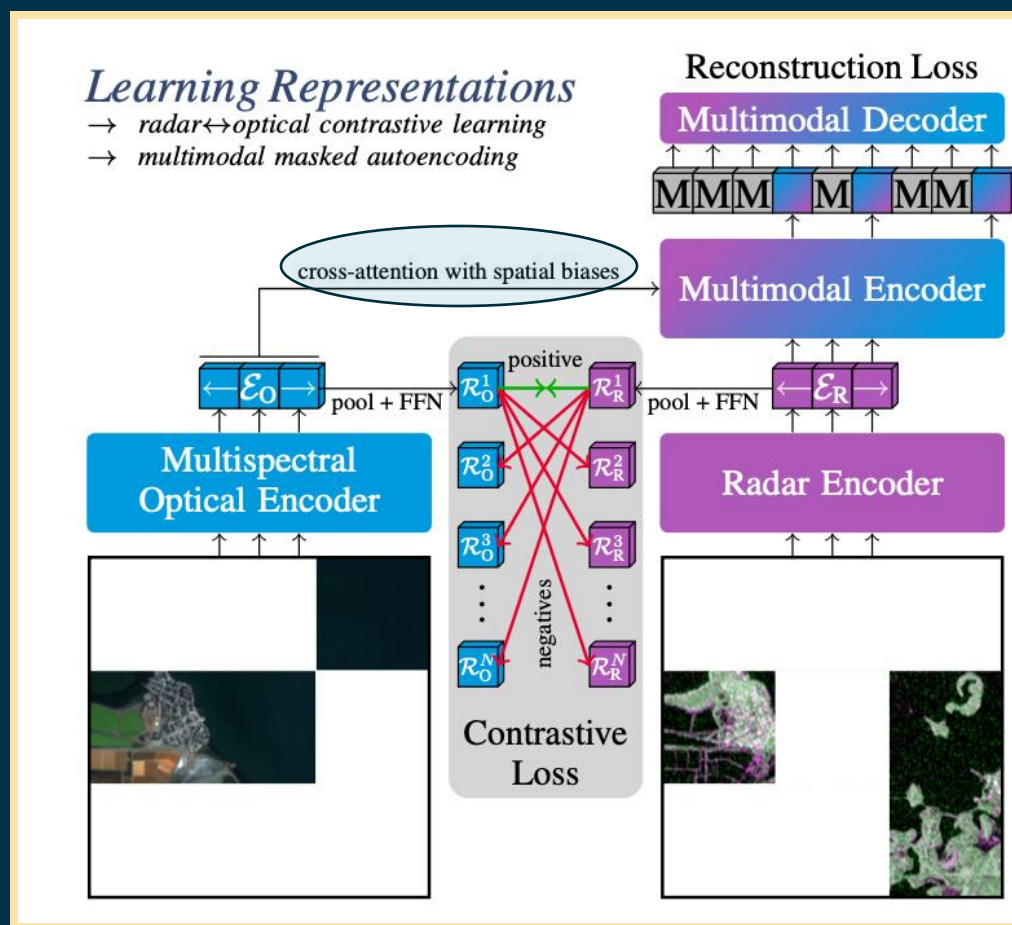
Contrastive + Masked Reconstruction!

Each ViT backbone has its self-attention heads biased based on the Euclidean distance between query-key pairs



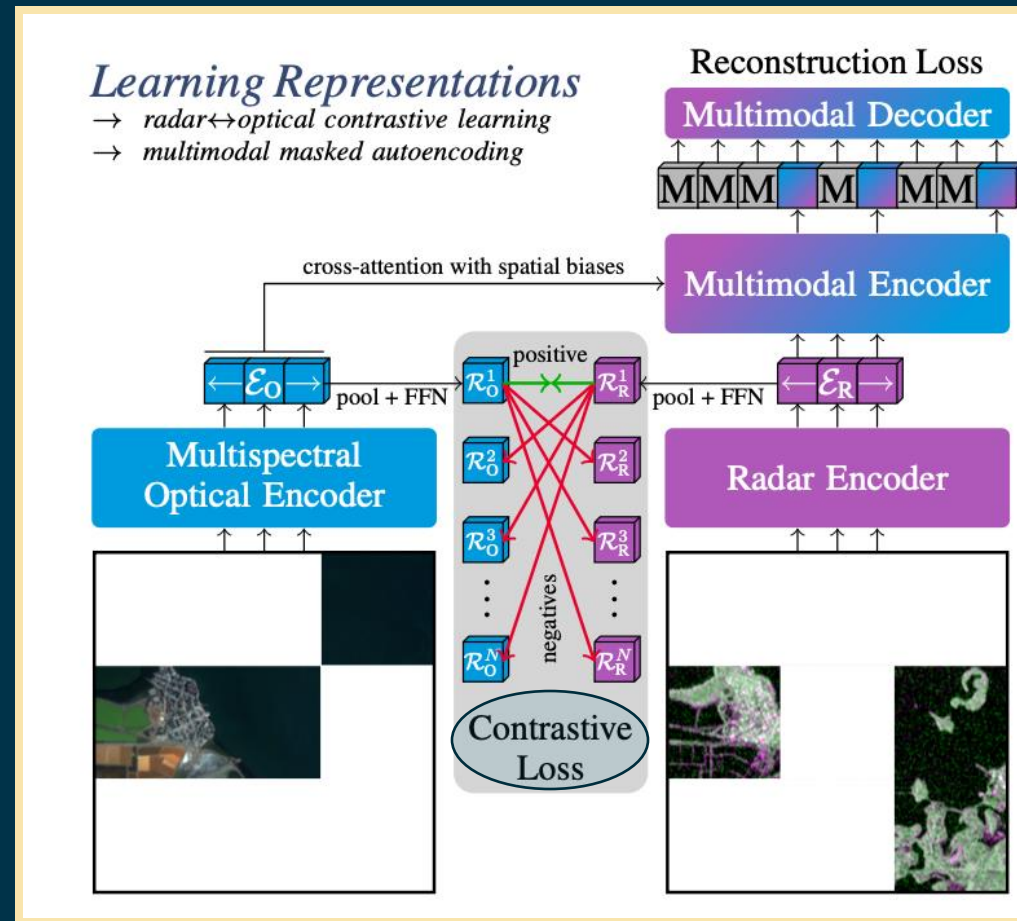
Contrastive + Masked Reconstruction!

Each ViT backbone has its self-attention heads biased based on the Euclidean distance between query-key pairs



The cross-attention matrix is modified by the Euclidean distance between *cross-modal* key-query pairs

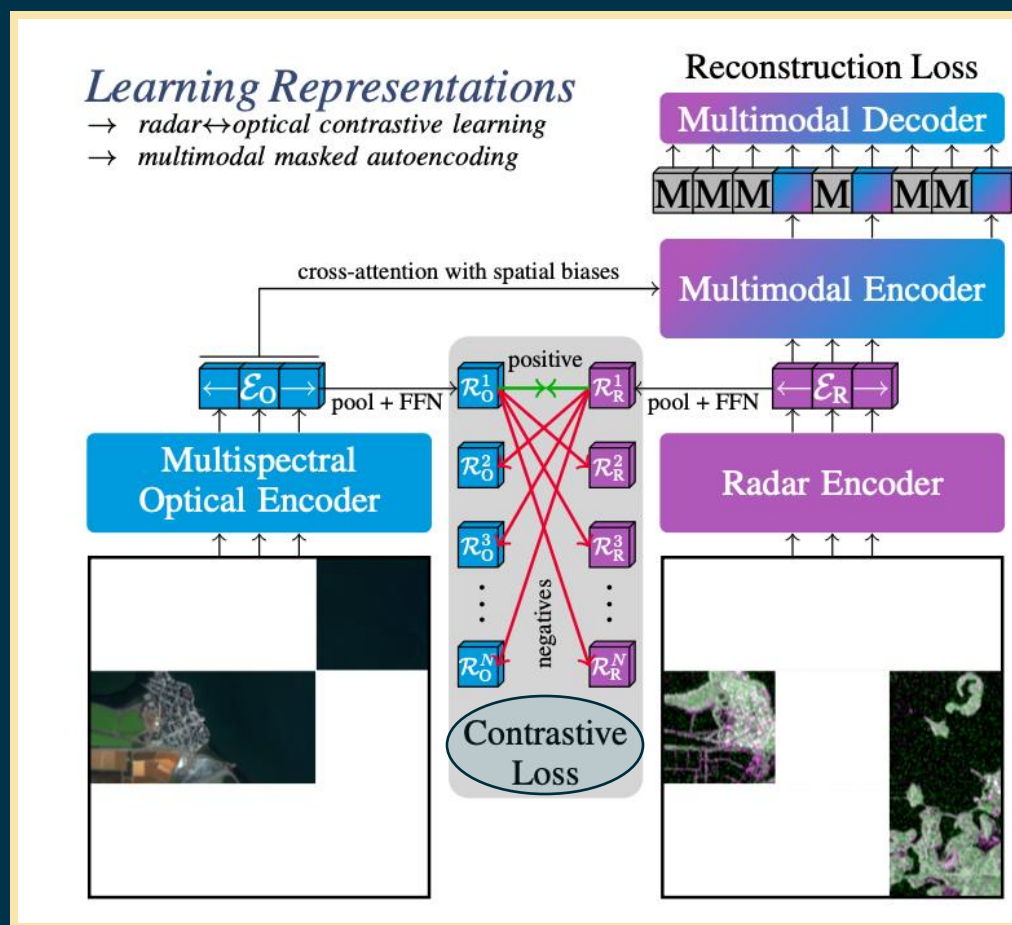
Contrastive + Masked Reconstruction!



Contrastive + Masked Reconstruction!

For optical anchor image, positive sample is the geographic, temporal match in the radar sample

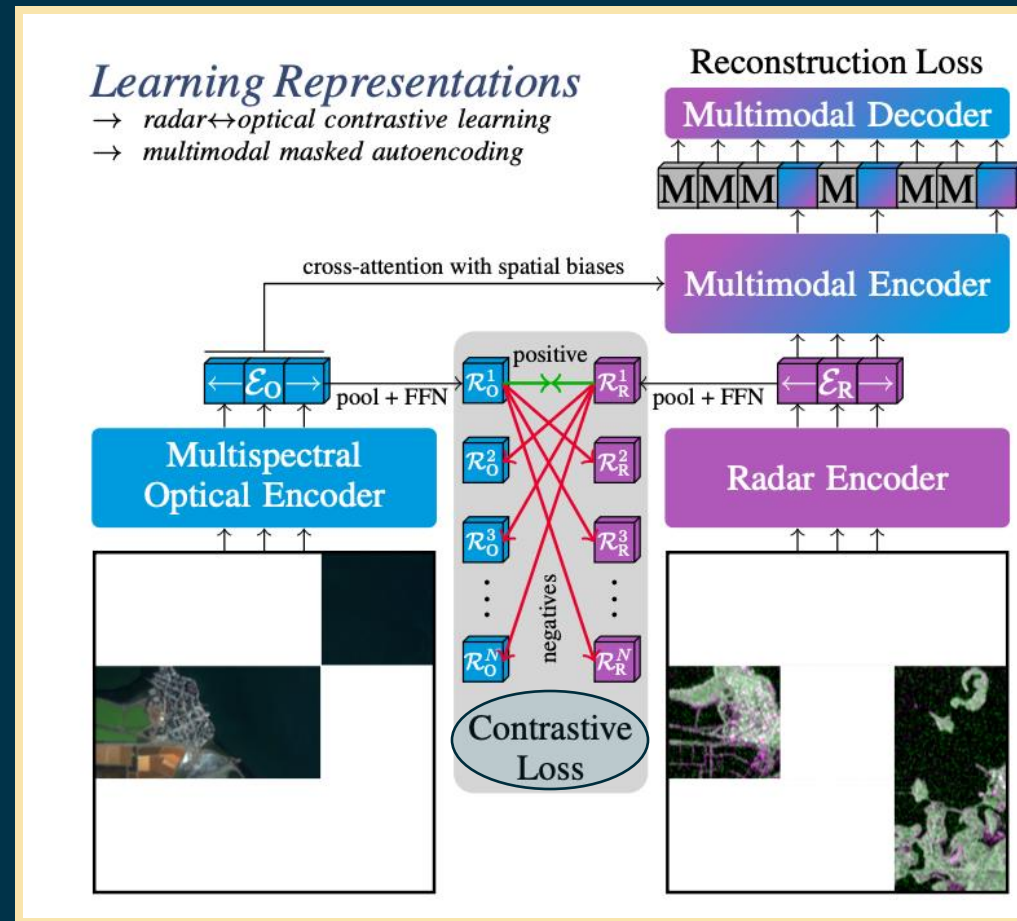
Negative sample? *All other radar samples in the batch!*



Contrastive + Masked Reconstruction!

For optical anchor image, positive sample is the geographic, temporal match in the radar sample

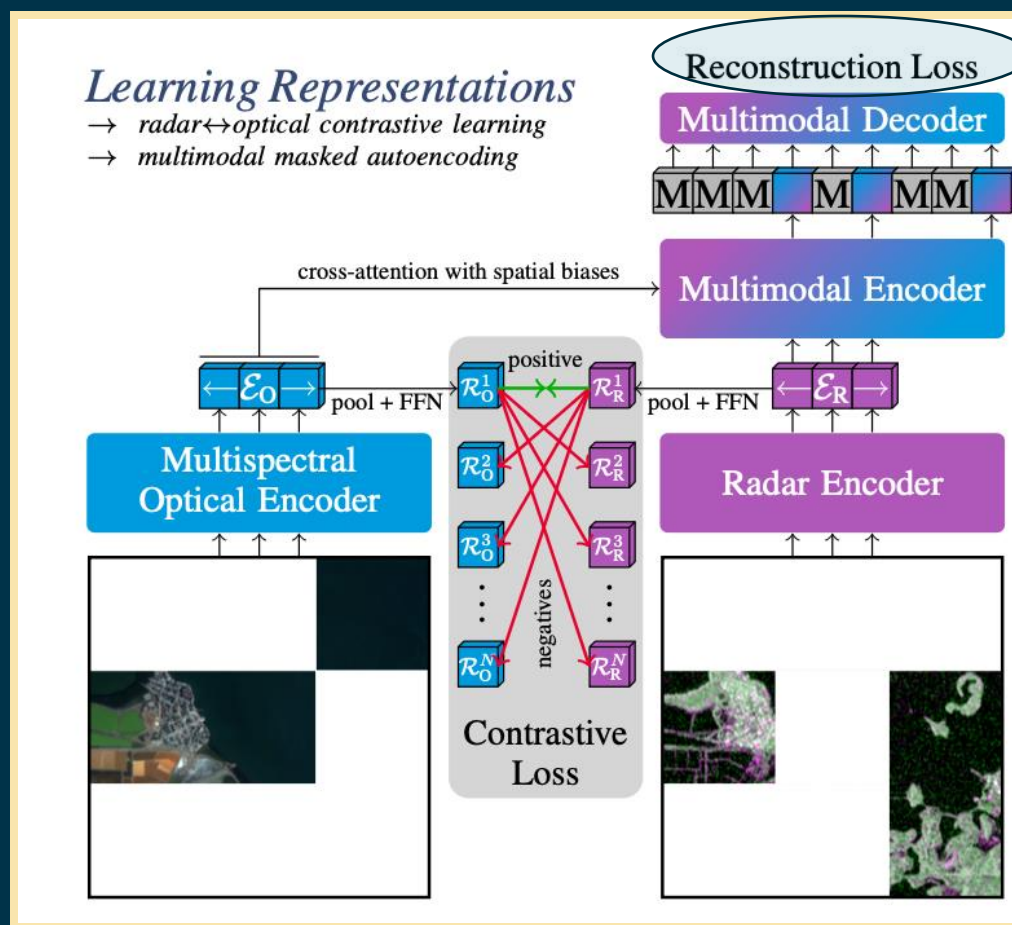
Negative sample? *All other radar samples in the batch!*



And vice versa!

Contrastive + Masked Reconstruction!

Independently
mask 75% of
radar and optical
patches



Loss term
incorporates
both radar and
optical
reconstruction
distances

In this talk

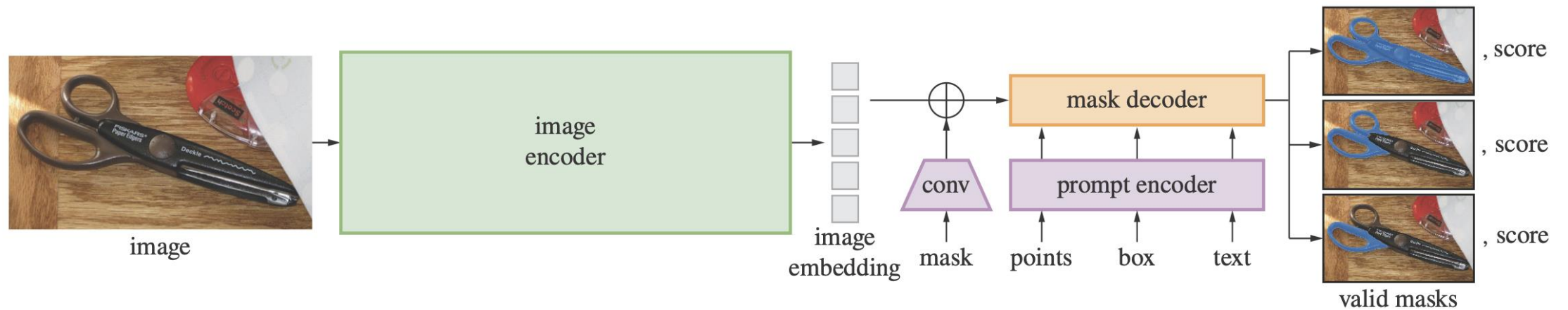
Domain
Challenges

Remote Sensing

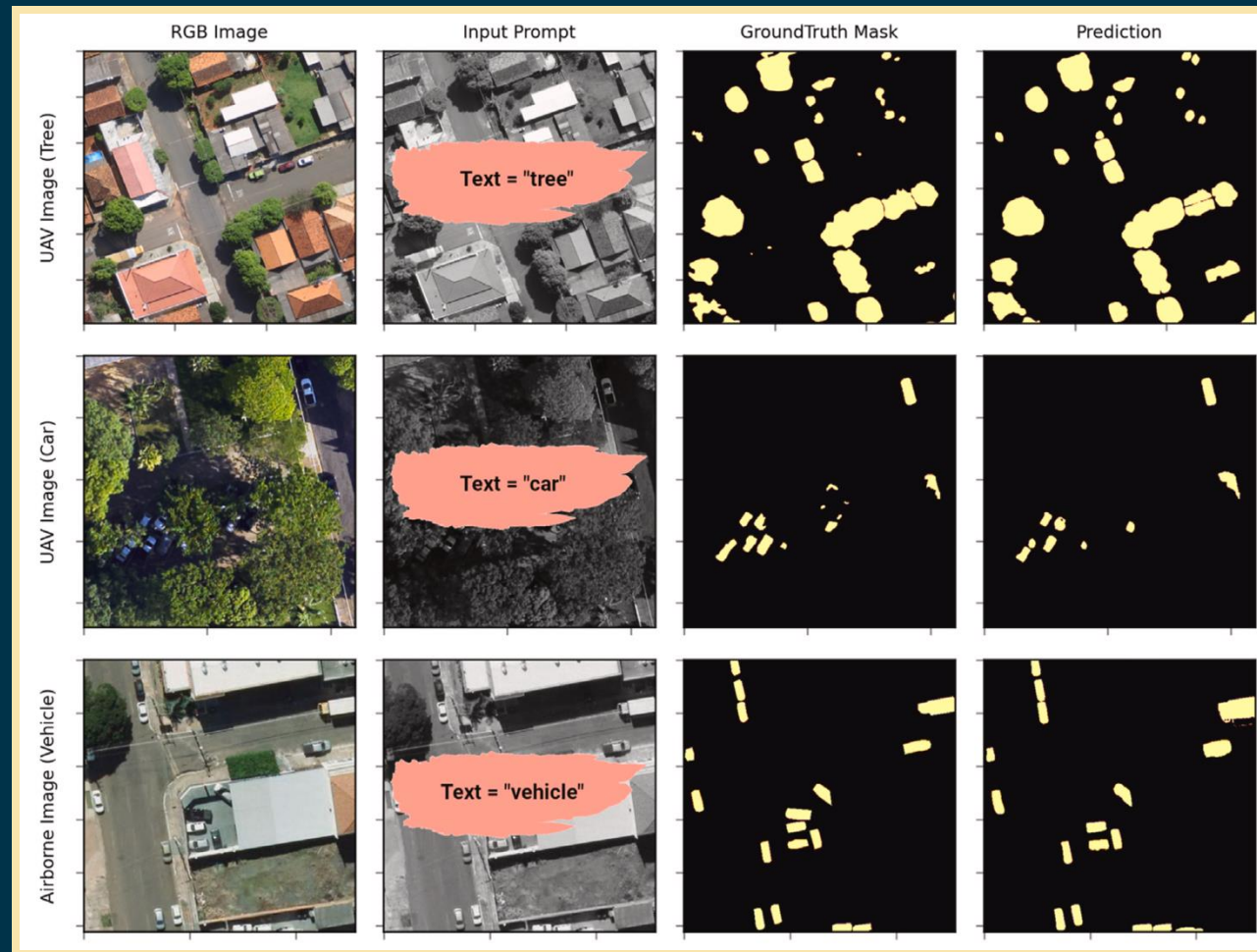
Pre-training

Zero-Shot

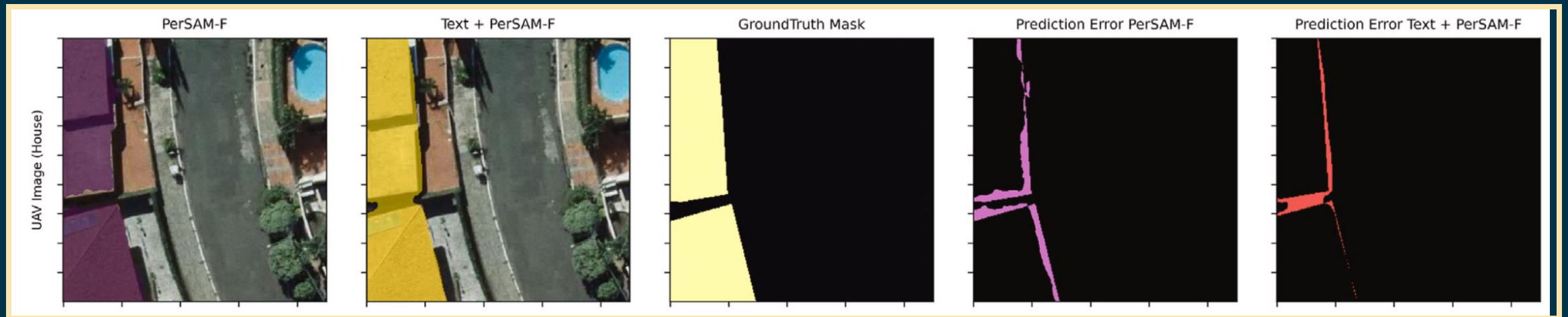
Segment Anything



Zero-Shot with SAM



One-Shot with SAM



Caution!

Concepts are still broad – special cases may be challenging

Tends to overestimate object boundaries – may be a prompt issue

Summary

Domain
Challenges

Remote Sensing

Pre-training

Zero-Shot