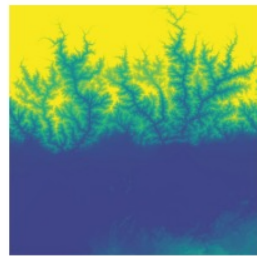


“Satellite Data is a Distinct Modality in ML”

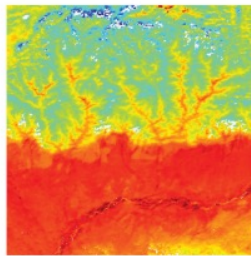
Tom Davies

Products



ALOS DEM

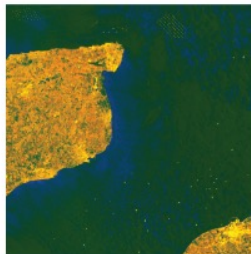
Eastern Himalayas,
English Channel



MODIS day temp.

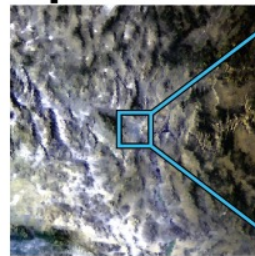


NDVI

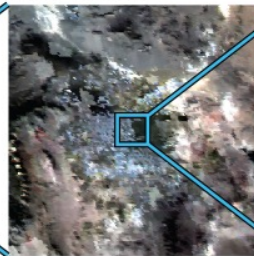


SAR (Sentinel 1)

Spatial resolutions



GOES-18 at 2000m/px



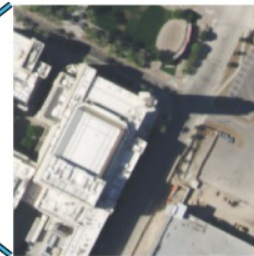
MODIS at 250m/px



Landsat 9 at 30m/px



Sentinel 2 at 10m/px



NAIP at 0.6m/px

Las Vegas, Nevada, USA

Time steps



Dec. 25, 1973



Dec. 3, 1982



Dec. 9, 1993



Dec. 24, 2001



Dec. 23, 2013



Dec. 28, 2023

Las Vegas, Nevada, USA

Context

- ICML ‘24 position paper

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

Esther Rolf^{* 1 2} Konstantin Klemmer³ Caleb Robinson⁴ Hannah Kerner^{* 5}

Abstract

Satellite data has the potential to inspire a seismic shift for machine learning—one in which we rethink existing practices designed for traditional data modalities. As machine learning for satellite data (SatML) gains traction for its real-world impact, our field is at a crossroads. We can either continue applying ill-suited approaches, or we can initiate a new research agenda that centers around the unique characteristics and challenges of satellite data. This position paper argues that satellite data constitutes a distinct modality for machine learning research and that we must recognize it as such to advance the quality and impact of SatML research across theory, methods, and deployment. We outline critical discussion questions and actionable suggestions to transform SatML from merely an intriguing application area to a dedicated research discipline that helps move the needle on big challenges for machine learning and society.

Satellite data presents challenges and opportunities distinct from other data modalities (Figure 1). Unlike natural images, the size of targets in satellite images span a logarithmic scale from < 1m (e.g., trees) to > 1km (e.g., forests). Temporal patterns in satellite time series also span logarithmic scales, from hours or days (e.g., floods) to years or decades (e.g., sea level rise). Data are acquired using a variety of sensors that capture diverse spectral channels (beyond 3-channel RGB) and precise measurements (beyond 8 bits). Satellites collect data over the entire surface of the Earth at fixed time intervals and spatial resolutions. Observations are acquired from an overhead perspective from fixed altitudes and lack a “natural” orientation, unlike natural images.

While there has been increasing interest in ML for satellite data (SatML) (Zhu et al. (2017); Table A1), SatML research as a whole falls short on these challenges and opportunities. The mainstream approach to SatML has been to adapt or “lift and shift” solutions designed for other modalities, especially natural images, to satellite data with minimal tailoring. Rather than tackling novel or outstanding challenges, many studies propose ML solutions to well-resolved or low-priority problems in the field of remote sensing (Tuia

Context

- ICML ‘24 position paper
- Supported by Microsoft AI for social good lab – has a geospatial ML project
 - Includes TorchGeo

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

Esther Rolf^{* 1 2} Konstantin Klemmer³ Caleb Robinson⁴ Hannah Kerner^{* 5}

Abstract

Satellite data has the potential to inspire a seismic shift for machine learning—one in which we rethink existing practices designed for traditional data modalities. As machine learning for satellite data (SatML) gains traction for its real-world impact, our field is at a crossroads. We can either continue applying ill-suited approaches, or we can initiate a new research agenda that centers around the unique characteristics and challenges of satellite data. This position paper argues that satellite data constitutes a distinct modality for machine learning research and that we must recognize it as such to advance the quality and impact of SatML research across theory, methods, and deployment. We outline critical discussion questions and actionable suggestions to transform SatML from merely an intriguing application area to a dedicated research discipline that helps move the needle on big challenges for machine learning and society.

Satellite data presents challenges and opportunities distinct from other data modalities (Figure 1). Unlike natural images, the size of targets in satellite images span a logarithmic scale from < 1m (e.g., trees) to > 1km (e.g., forests). Temporal patterns in satellite time series also span logarithmic scales, from hours or days (e.g., floods) to years or decades (e.g., sea level rise). Data are acquired using a variety of sensors that capture diverse spectral channels (beyond 3-channel RGB) and precise measurements (beyond 8 bits). Satellites collect data over the entire surface of the Earth at fixed time intervals and spatial resolutions. Observations are acquired from an overhead perspective from fixed altitudes and lack a “natural” orientation, unlike natural images.

While there has been increasing interest in ML for satellite data (SatML) (Zhu et al. (2017); Table A1), SatML research as a whole falls short on these challenges and opportunities. The mainstream approach to SatML has been to adapt or “lift and shift” solutions designed for other modalities, especially natural images, to satellite data with minimal tailoring. Rather than tackling novel or outstanding challenges, many studies propose ML solutions to well-resolved or low-priority problems in the field of remote sensing (Tuia

Context

- ICML ‘24 position paper
- Supported by Microsoft AI for social good lab – has a geospatial ML project
 - Includes TorchGeo
- Paper in three parts:
 - Data
 - Tailored ML (adapt ML for sat)
 - Core ML (ML learns from sat)

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

Esther Rolf^{* 1 2} Konstantin Klemmer³ Caleb Robinson⁴ Hannah Kerner^{* 5}

Abstract

Satellite data has the potential to inspire a seismic shift for machine learning—one in which we rethink existing practices designed for traditional data modalities. As machine learning for satellite data (SatML) gains traction for its real-world impact, our field is at a crossroads. We can either continue applying ill-suited approaches, or we can initiate a new research agenda that centers around the unique characteristics and challenges of satellite data. This position paper argues that satellite data constitutes a distinct modality for machine learning research and that we must recognize it as such to advance the quality and impact of SatML research across theory, methods, and deployment. We outline critical discussion questions and actionable suggestions to transform SatML from merely an intriguing application area to a dedicated research discipline that helps move the needle on big challenges for machine learning and society.

Satellite data presents challenges and opportunities distinct from other data modalities (Figure 1). Unlike natural images, the size of targets in satellite images span a logarithmic scale from < 1m (e.g., trees) to > 1km (e.g., forests). Temporal patterns in satellite time series also span logarithmic scales, from hours or days (e.g., floods) to years or decades (e.g., sea level rise). Data are acquired using a variety of sensors that capture diverse spectral channels (beyond 3-channel RGB) and precise measurements (beyond 8 bits). Satellites collect data over the entire surface of the Earth at fixed time intervals and spatial resolutions. Observations are acquired from an overhead perspective from fixed altitudes and lack a “natural” orientation, unlike natural images.

While there has been increasing interest in ML for satellite data (SatML) (Zhu et al. (2017); Table A1), SatML research as a whole falls short on these challenges and opportunities. The mainstream approach to SatML has been to adapt or “lift and shift” solutions designed for other modalities, especially natural images, to satellite data with minimal tailoring. Rather than tackling novel or outstanding challenges, many studies propose ML solutions to well-resolved or low-priority problems in the field of remote sensing (Tuia

Context

- ICML ‘24 position paper
- Supported by Microsoft AI for social good lab – has a geospatial ML project
 - Includes TorchGeo
- Paper in three parts:
 - Data
 - Tailored ML (adapt ML for sat)
 - Core ML (ML learns from sat)
- Satellite data: excludes aerial images, climate sensors

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

Esther Rolf^{* 1 2} Konstantin Klemmer³ Caleb Robinson⁴ Hannah Kerner^{* 5}

Abstract

Satellite data has the potential to inspire a seismic shift for machine learning—one in which we rethink existing practices designed for traditional data modalities. As machine learning for satellite data (SatML) gains traction for its real-world impact, our field is at a crossroads. We can either continue applying ill-suited approaches, or we can initiate a new research agenda that centers around the unique characteristics and challenges of satellite data. This position paper argues that satellite data constitutes a distinct modality for machine learning research and that we must recognize it as such to advance the quality and impact of SatML research across theory, methods, and deployment. We outline critical discussion questions and actionable suggestions to transform SatML from merely an intriguing application area to a dedicated research discipline that helps move the needle on big challenges for machine learning and society.

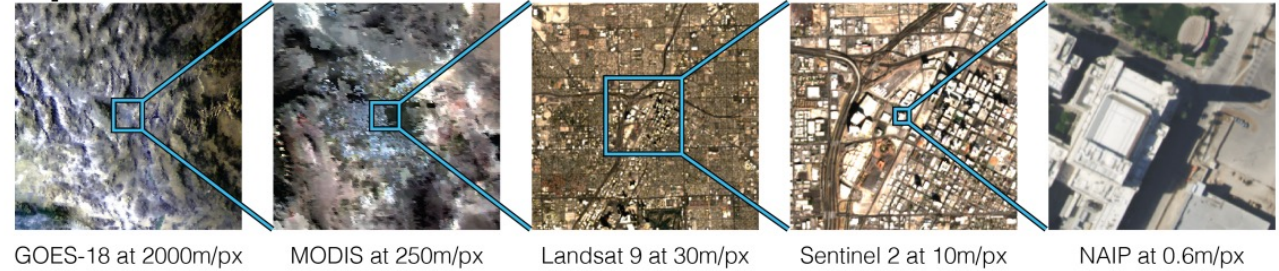
Satellite data presents challenges and opportunities distinct from other data modalities (Figure 1). Unlike natural images, the size of targets in satellite images span a logarithmic scale from < 1m (e.g., trees) to > 1km (e.g., forests). Temporal patterns in satellite time series also span logarithmic scales, from hours or days (e.g., floods) to years or decades (e.g., sea level rise). Data are acquired using a variety of sensors that capture diverse spectral channels (beyond 3-channel RGB) and precise measurements (beyond 8 bits). Satellites collect data over the entire surface of the Earth at fixed time intervals and spatial resolutions. Observations are acquired from an overhead perspective from fixed altitudes and lack a “natural” orientation, unlike natural images.

While there has been increasing interest in ML for satellite data (SatML) (Zhu et al. (2017); Table A1), SatML research as a whole falls short on these challenges and opportunities. The mainstream approach to SatML has been to adapt or “lift and shift” solutions designed for other modalities, especially natural images, to satellite data with minimal tailoring. Rather than tackling novel or outstanding challenges, many studies propose ML solutions to well-resolved or low-priority problems in the field of remote sensing (Tuia

Satellite data is a distinct modality

- Logarithmic feature scales
 - Spatially: from meters (trees) to KMs (forests)

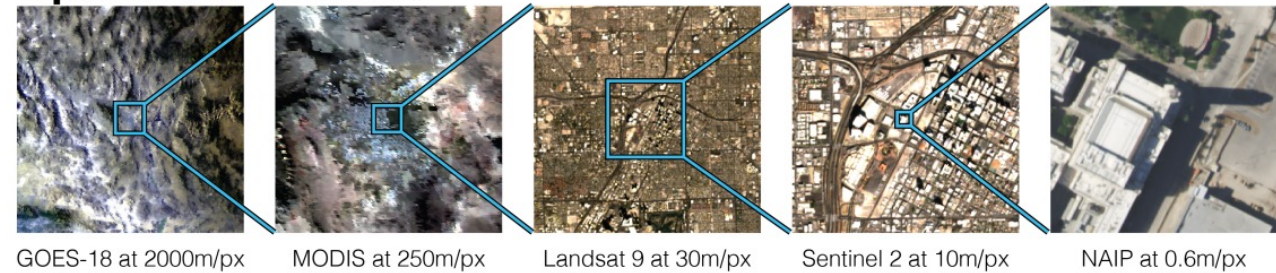
Spatial resolutions



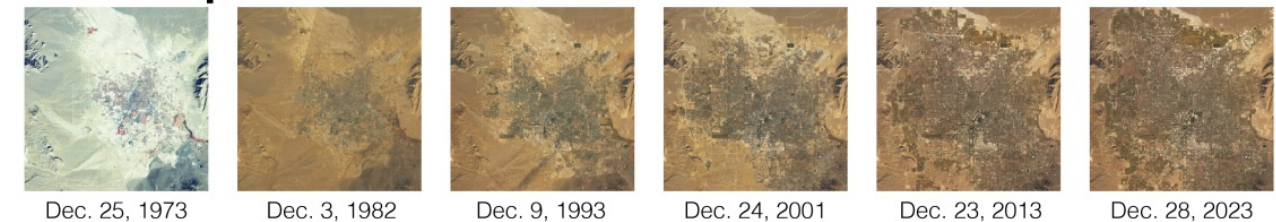
Satellite data is a distinct modality

- Logarithmic feature scales
 - Spatially: from meters (trees) to KMs (forests)
 - Temporally: hours (earthquakes), weeks (construction), seasons (crops), years (glacial retreat), decades (sea levels)

Spatial resolutions

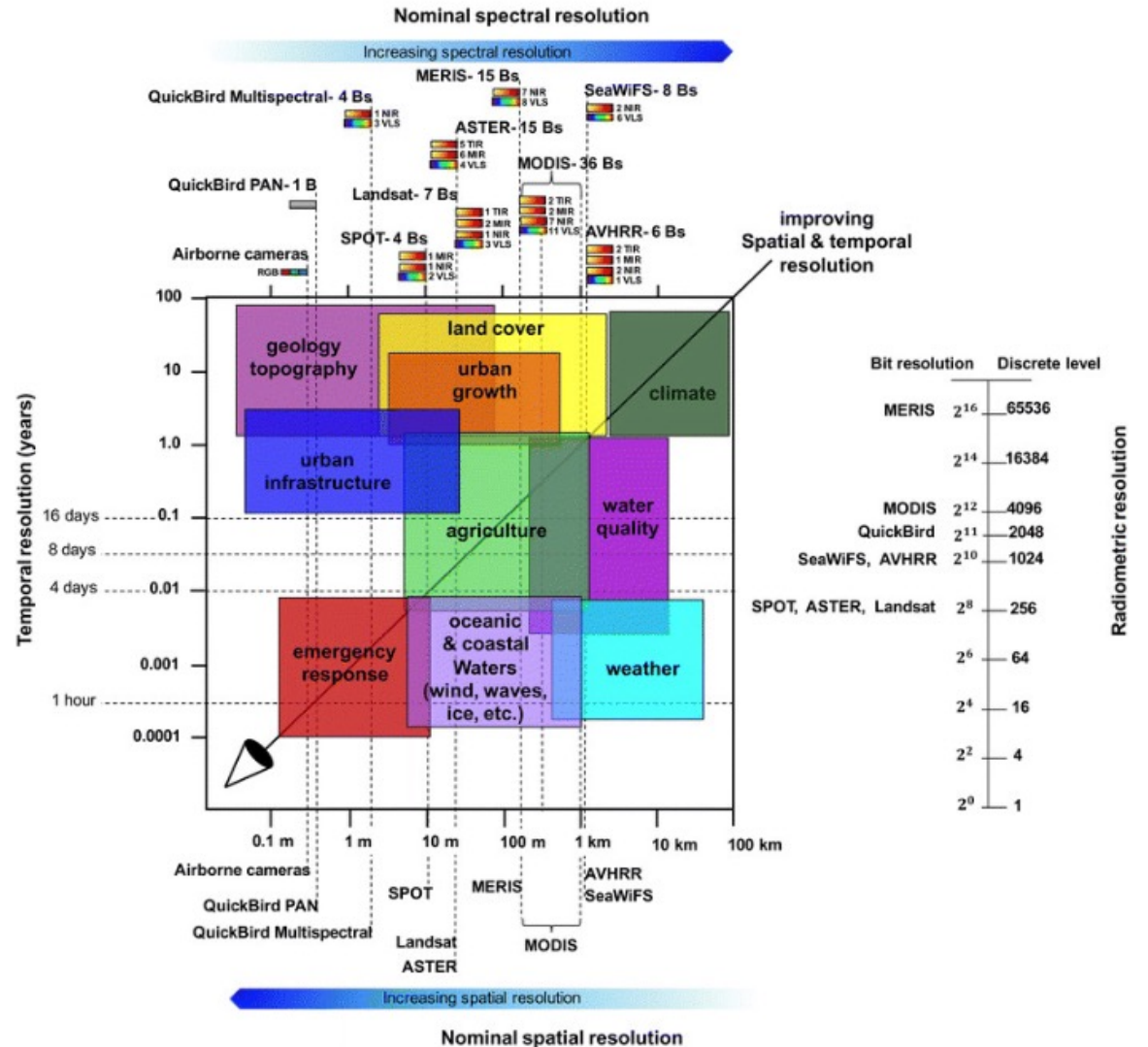


Time steps



Satellite data is a distinct modality

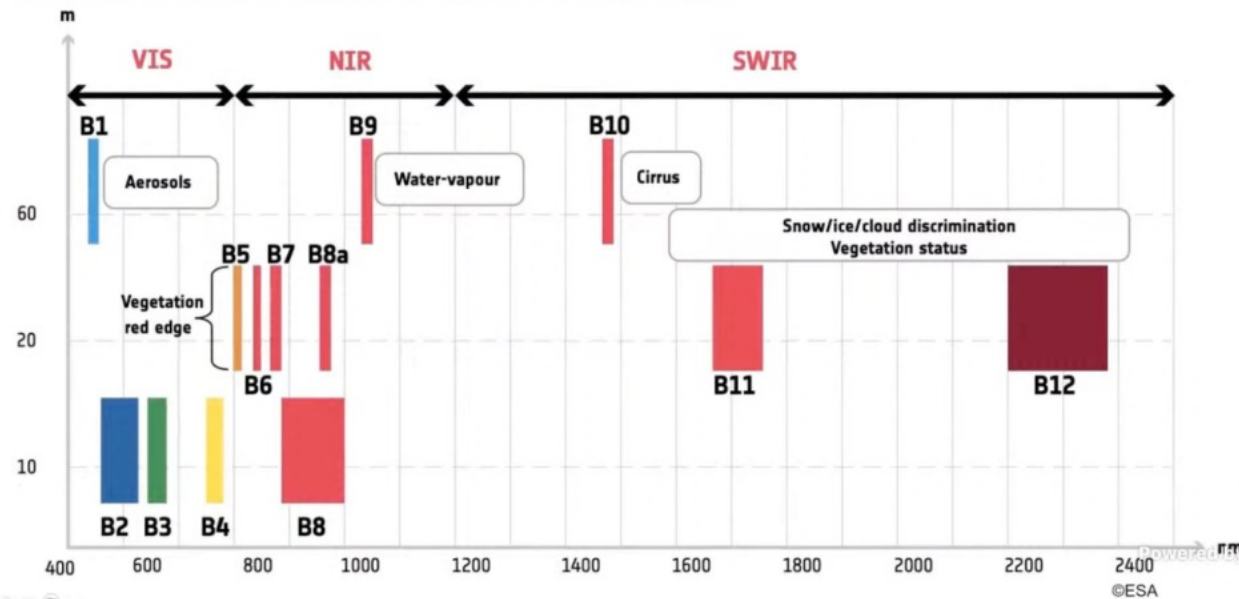
- Logarithmic feature scales
- Varying T+S resolutions
 - Instruments and orbits vary massively between satellites



Satellite data is a distinct modality

- Logarithmic feature scales
- Varying T+S resolutions
- Non-standard image channels
 - Sensors correspond to different wavelengths or wildly different formats (SAR, LIDAR, altimetry).
 - Higher radiometric resolution than 8 bit (Sentintel-2 has 12 bit)
 - Leads to a variety of data formats

ESA Sentinel-2 satellite spectral channels



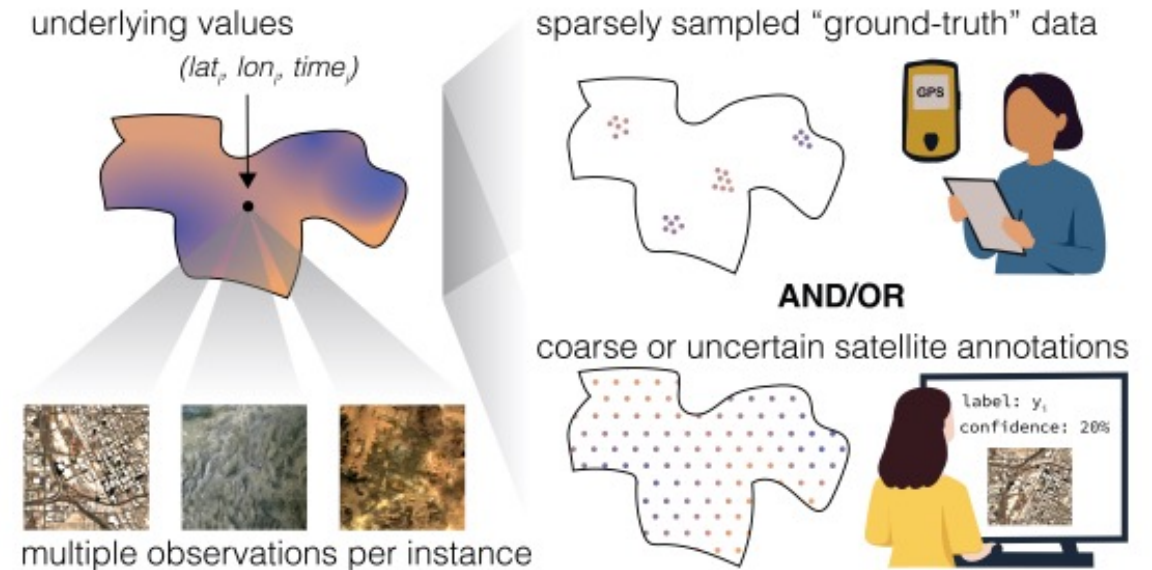
Satellite data is a distinct modality

- Logarithmic feature scales
- Varying T+S resolutions
- Non-standard image channels
- Vast amounts of data
 - Petabyte scales of data
 - GPT-3 trained on 570GB of text
 - Largest text-image dataset is 220TB
 - Upcoming NISAR is 85 TB *per day*



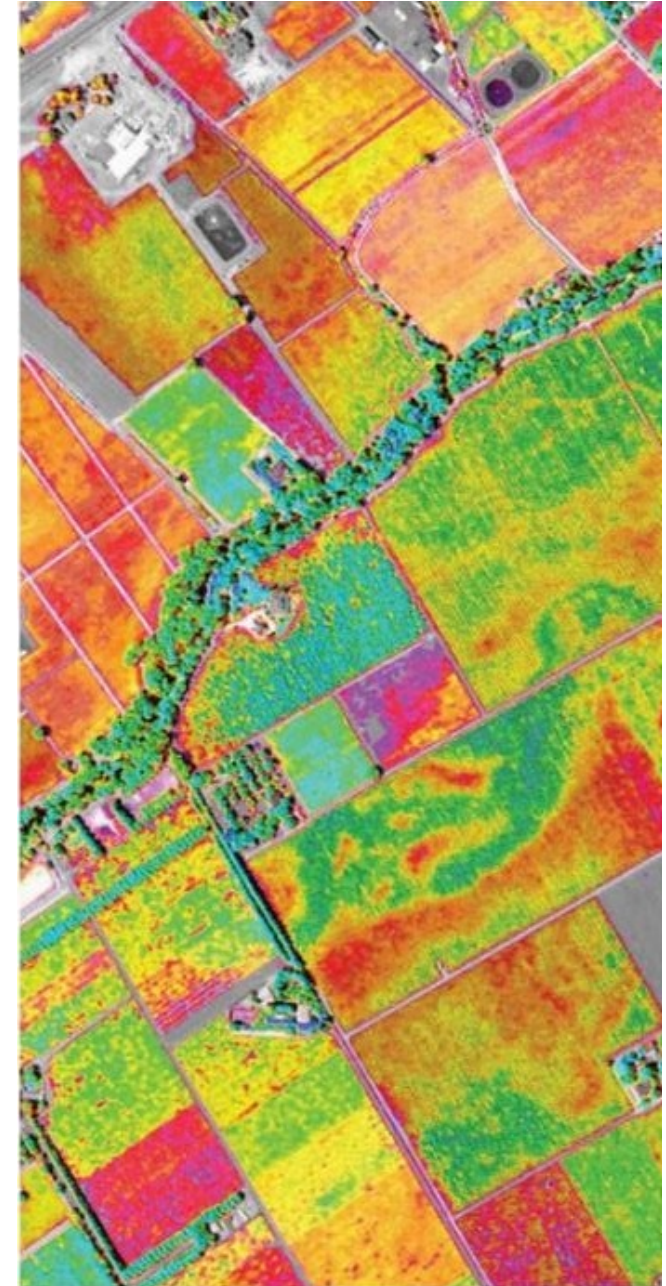
Satellite data is a distinct modality

- Logarithmic feature scales
- Varying T+S resolutions
- Non-standard image channels
- Vast amounts of data
- Data annotations
 - Sparse
 - Uncertain
 - Often difficult to collect
 - More common in e.g., the Global North



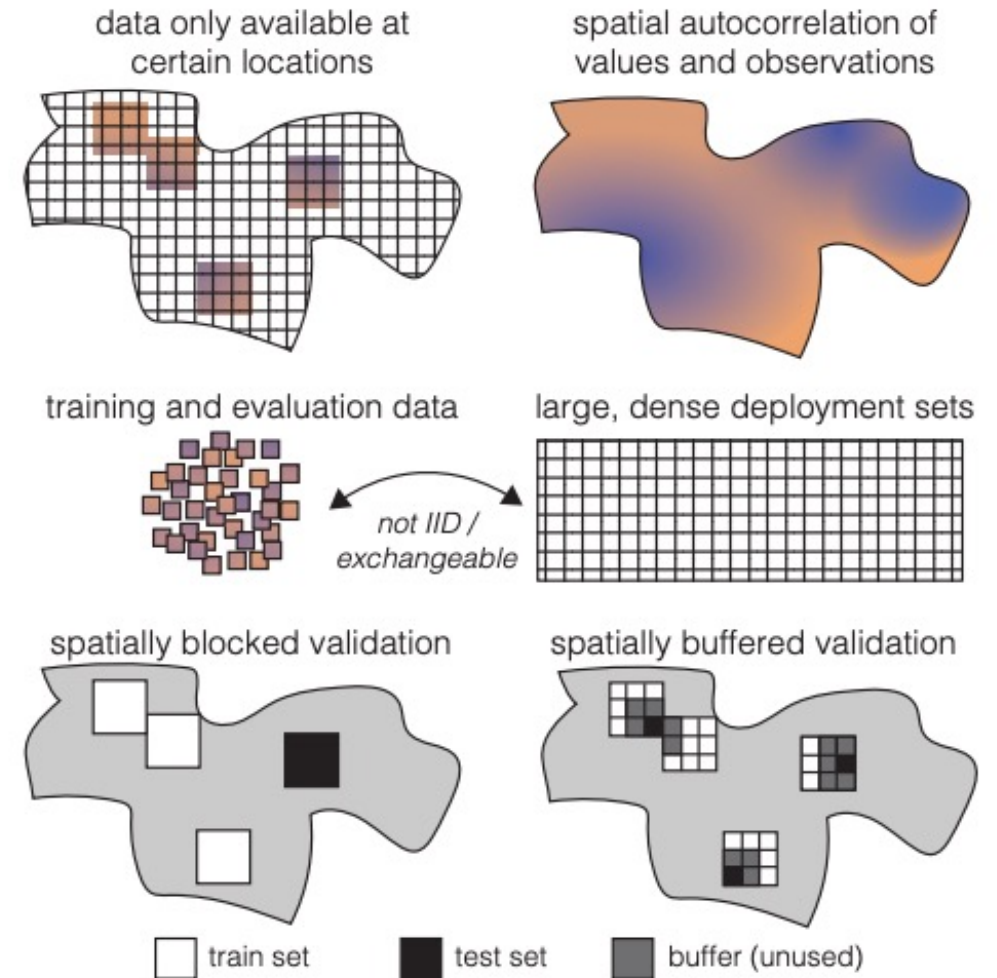
Deployment challenges

- Desired output is often a set of dense predictions
 - E.g., crop type prediction and many others requires pixel-level classification
 - Some tasks may need to be updated daily, e.g., deforestation monitoring
 - Extremely computationally expensive - makes model efficiency important



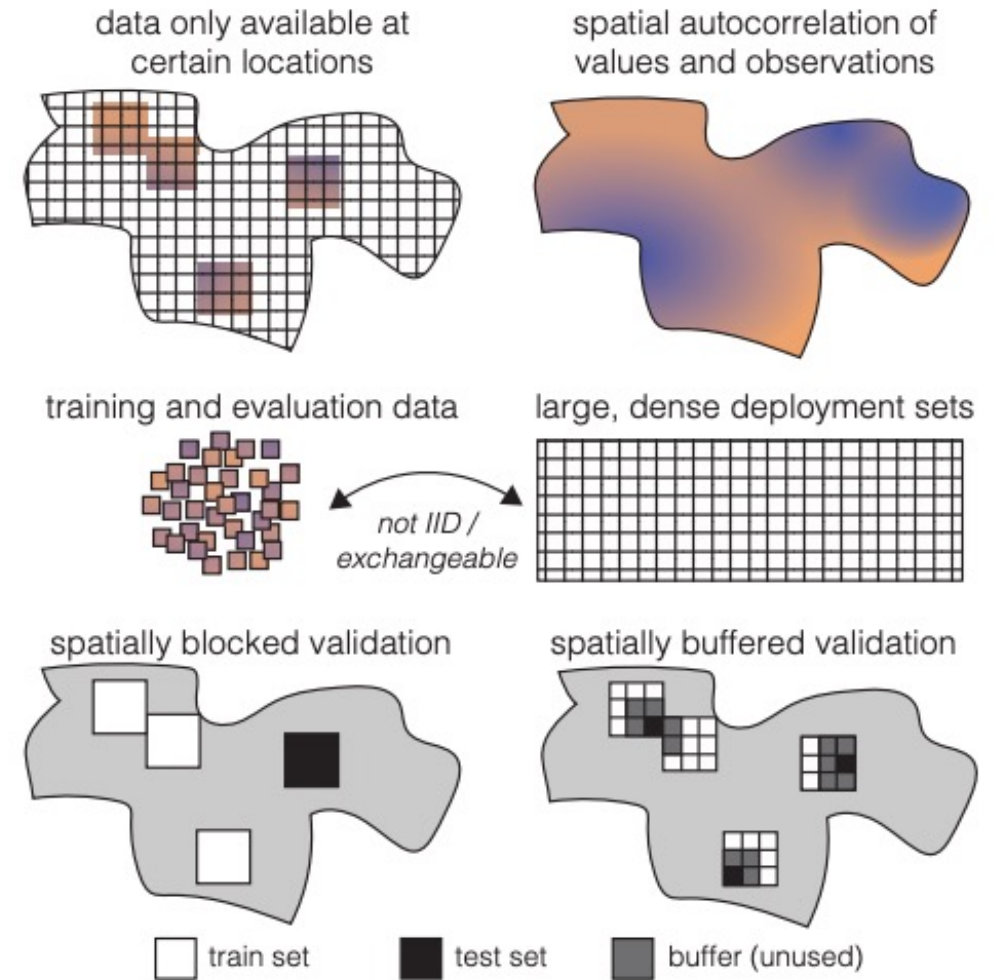
Evaluation challenges

- Uniformly at random sampled train-test splits are generally not appropriate for SatML data
 - SatML data is often temporally or spatially correlated, and your splits should account for this



Evaluation challenges

- Uniformly at random sampled train-test splits are generally not appropriate for SatML data
- Reported metrics often aren't informative for real-world applications
 - Require reliable estimation of metrics from cluster-sampled data
 - Require notions of 'regions of applicability'



Ethical concerns

- Like many large-scale datasets, there is no clear way to opt-out of satellite data that covers you or your community
- Satellite data is excluded from laws that would prevent otherwise invasive images of your land or self being taken/used
- Irresponsible/inaccurate outcomes from SatML models can worsen outcomes for product users



Specialising ML for satellite data

- Learning strategies
 - Pre-training models on satellite data instead of ImageNet can improve performance by up to 18% (Bastani et al., 2023)
 - Including information like geographic domain, scale, temporal dimension, and availability of sensor in the learning strategy also increases performance.
 - Ayush et al. (2021) add auxiliary geo task.
- Model architectures
- Domain context

Method	Average	
	50	All
Random Initialization	0.33	0.65
ImageNet [25]	0.38	0.75
BigEarthNet [51]	0.40	0.76
MillionAID [39]	0.49	0.78
DOTA [55]	0.50	0.79
iSAID [58]	0.50	0.79
MoCo [17]	0.21	0.27
SeCo [40]	0.38	0.71
SatlasPretrain	0.56	0.83

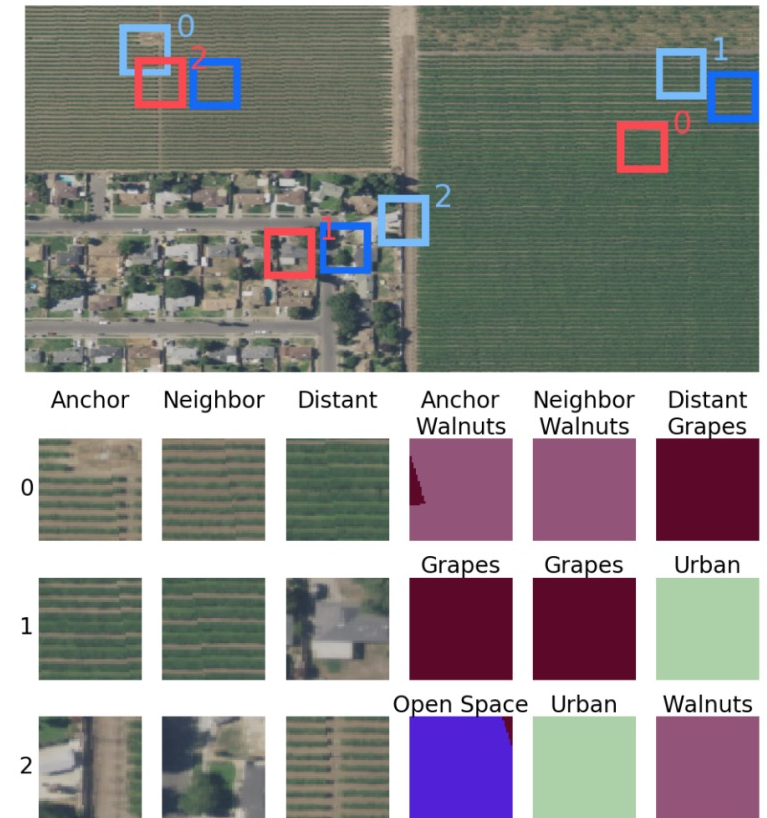
Specialising ML for satellite data

- Learning strategies
- Model architectures
 - SatML specific architectures add rotation invariance, which isn't typically required
 - Architectures proposed that focus on small objects in very high-resolution imagery
 - Architectures which are designed to effectively leverage data from similar sensors
- Domain context



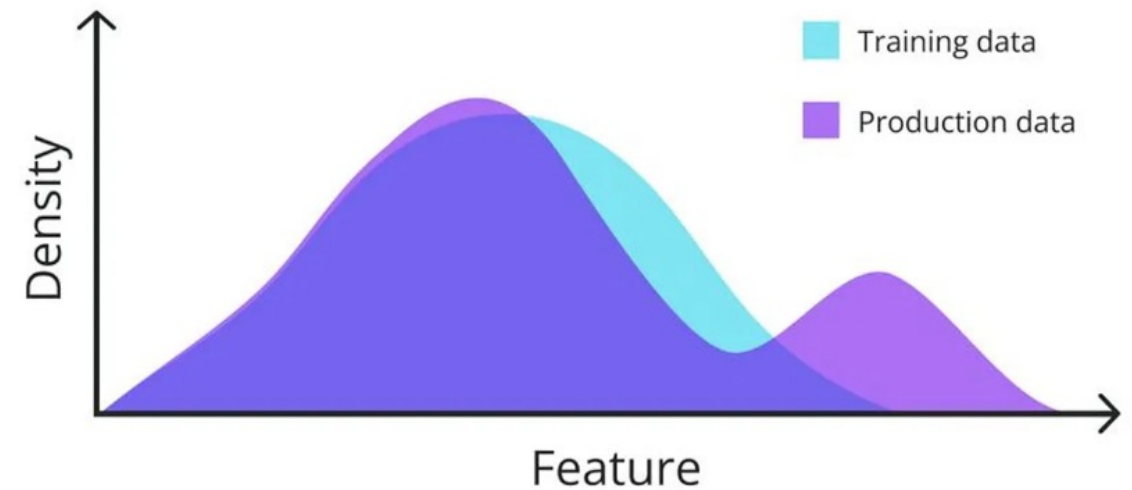
Specialising ML for satellite data

- Learning strategies
- Model architectures
- Domain context
 - Tobler's first law of geography: *Everything is related to everything else, but nearby things are related more*
 - Incorporating autocorrelation measures into the model loss (including a triplet loss)
 - Man-made objects remain fixed temporally, used for change point detection



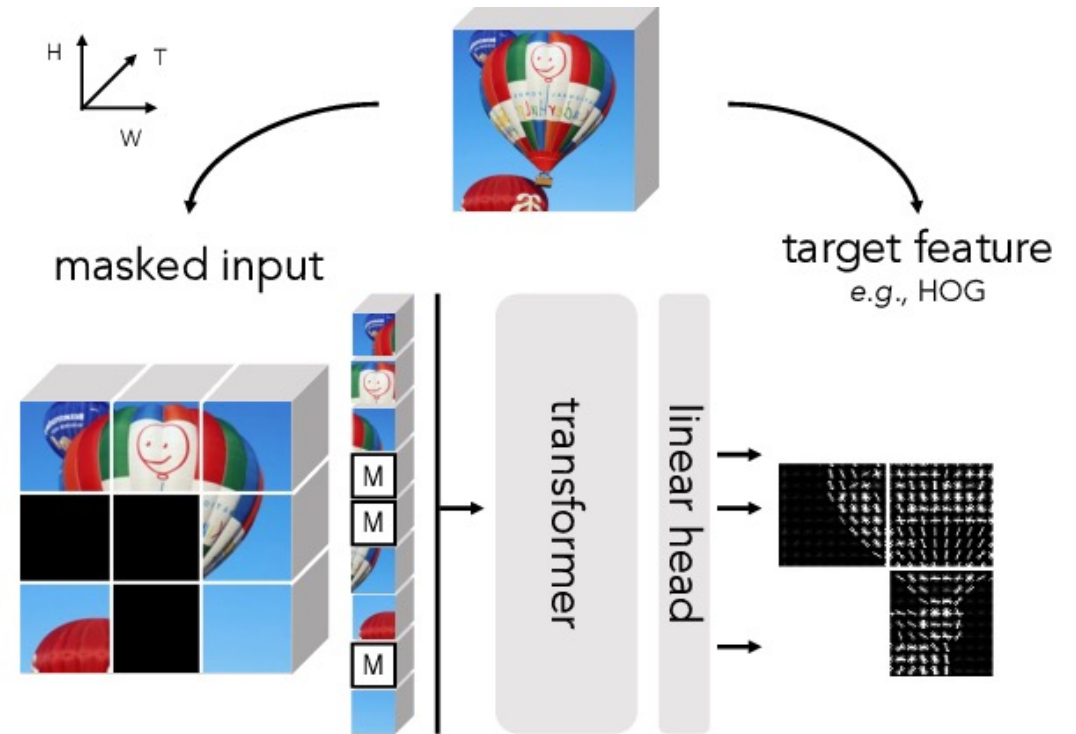
Learnings from SatML for 'core' ML research

- Distribution shift
 - Pervasive in SatML due to temporal shifts



Learnings from SatML for 'core' ML research

- Distribution shift
- Self-supervised learning
 - Vast amounts of unlabelled data makes this a key SatML problem



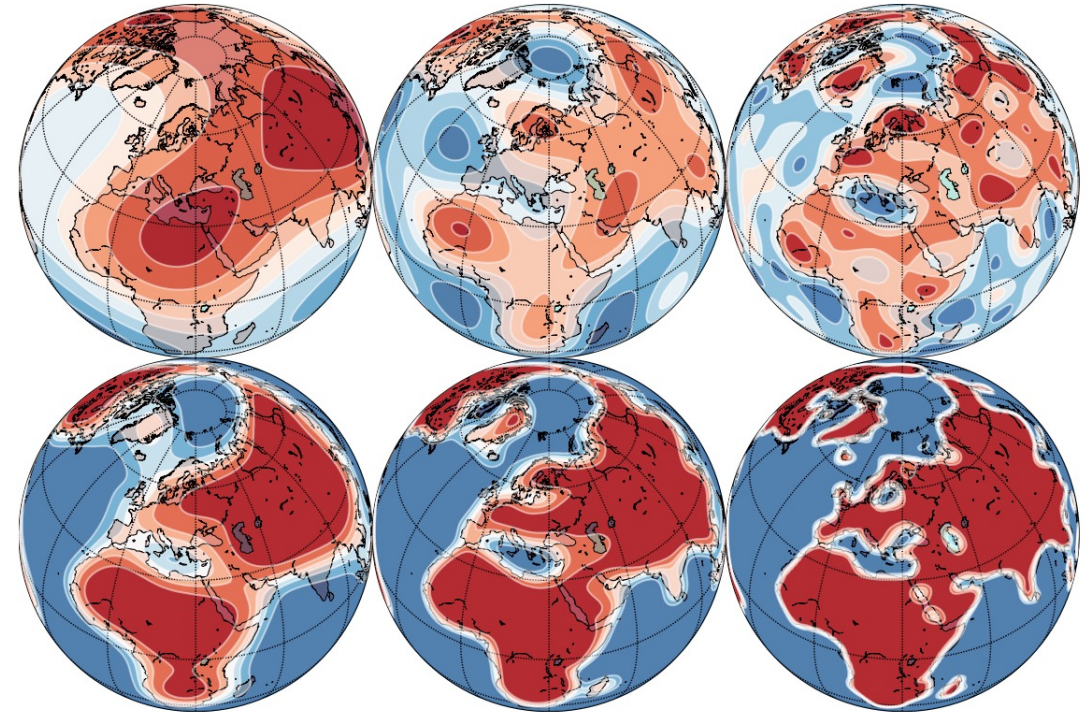
Learnings from SatML for ‘core’ ML research

- Distribution shift
- Self-supervised learning
- Multi-modal learning
 - SatML is fundamentally multi-modal, covering many bands and instruments



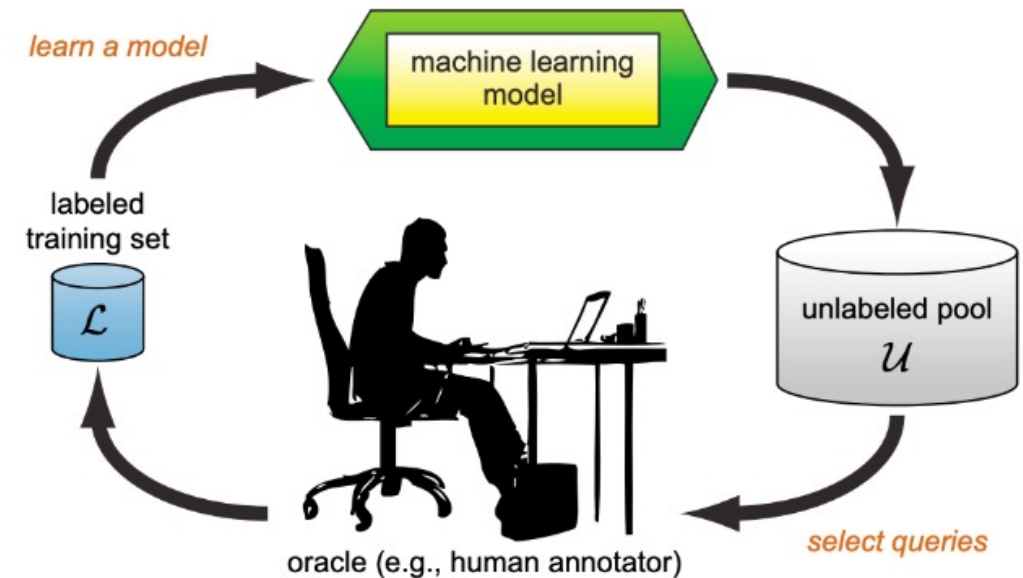
Learnings from SatML for ‘core’ ML research

- Distribution shift
- Self-supervised learning
- Multi-modal learning
- Positional encoding
 - SatML has already led to more advanced positional encodings



Learnings from SatML for ‘core’ ML research

- Distribution shift
- Self-supervised learning
- Multi-modal learning
- Positional encoding
- Human-in-the-loop and active learning
 - Model outputs should be HITL and active labelling for the sparse datasets



Discussion Topics

- How do we prioritise and collaborate on the most important SatML challenges?
- How will we align SatML research progress and real-world impact?
- How will we ensure SatML progress benefits global and local communities?