SatMAE: Pre-training Transformers for Temporal and Multi-Spectral Satellite Imagery

Erik Rozi Yutong He Marshall Burke David B. Lobell Stefano Ermon

Stanford University

Abstract

Unsupervised pre-training methods for large vision models have shown to enhance performance on downstream supervised tasks. Developing similar techniques for satellite imagery presents significant opportunities as unlabelled data is plentiful and the inherent temporal and multi-spectral structure provides avenues to further improve existing pre-training strategies. In this paper, we present SatMAE, a pre-training framework for temporal or multi-spectral satellite imagery based on Masked Autoencoder (MAE). To leverage temporal information, we include a temporal embedding along with independently masking image patches across time. In addition, we demonstrate that encoding multi-spectral data as groups of bands with distinct spectral positional encodings is beneficial. Our approach yields strong improvements over previous state-of-the-art techniques, both in terms of supervised learning performance on benchmark datasets (up to \uparrow 7%), and transfer learning performance on downstream remote sensing tasks, including land cover classification (up to \uparrow 14%) and semantic segmentation. Code and data are available on the project website: https://sustainlab-group.github.io/SatMAE/

1 Introduction

In recent years, self-supervised learning techniques have quickly become the norm for pre-training models on large-scale natural image datasets [1, 2, 3, 4, 5, 6, 7, 8], and have demonstrated strong performance on downstream tasks including image classification [3, 4, 9, 10], image segmentation [3, 11], representation learning [12, 13, 14], image compression [12, 15], image reconstruction [1], and image generation [16]. Unlike supervised learning approaches, self-supervised learning techniques do not require human labeling, making them appealing in settings where unlabeled data are abundant but labeled data are scarce, such as remote sensing data (e.g., satellite imagery). While several large-scale satellite image datasets have been carefully curated in the past few years, including Functional Map of the World (fMoW) [17], BigEarthNet [18], xView [19], SpaceNet [20], annotating these datasets requires specialized skills and is more expensive than traditional computer vision datasets. Moreover, automatic analysis of satellite imagery is often needed for tasks with large societal impact such as poverty or crop yield prediction [21, 22, 23, 24, 25, 26, 27, 28, 29, 30], where acquiring large amounts of labeled data through surveys is impossible or prohibitively expensive. This suggests that self-supervised learning approaches for satellite imagery could be especially valuable.

However, existing self-supervised learning approaches [1, 2, 3, 4, 5, 6] are mainly designed for natural images. As opposed to natural images such as ImageNet [31], satellite imagery is usually associated

^{*}Equal contribution. Order determined via coin flip.

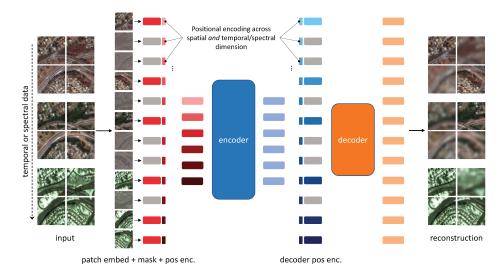


Figure 1: With carefully-designed masking strategies across mutli-spectral and temporal images, and temporal and spectral positional encodings, our SatMAE serves as a powerful SSL vision learner for remote sensing tasks.

with meaningful geographical and temporal information, and can consist of multiple spectral bands representing sensor readings besides visible light (i.e., RGB channels typical in natural images). Depending on the data source, satellite imagery can also vary significantly in resolution [32] [33]. While self-supervised learning methods for satellite imagery exist [34] [35], these approaches cannot learn general representations for both temporal and multi-spectral remote sensing data.

To address this issue we propose **SatMAE**, a self-supervised learning framework based on masked autoencoders (MAEs) which naturally handles temporal and multi-spectral input data. We show that introducing a positional encoding for the temporal/spectral dimension and *independently* masking patches across the temporal/spectral dimension benefits pre-training, allowing the model to learn representations of the data that are more conducive to finetuning. Specifically, our contributions are:

- 1. We propose a novel method to leverage temporal or multi-spectral information in satellite imagery to improve self-supervised pre-training with masked autoencoders (see 4).
- 2. We introduce fMoW-Sentinel, a new Sentinel-2 dataset cross-referenced with fMoW, as a benchmark for training models on multi-spectral satellite imagery (see [5.1]).
- 3. We demonstrate the effectiveness of pre-training transformers [36] on satellite imagery, achieving significant improvement over previous state-of-the-art methods on benchmark datasets as well as downstream remote sensing tasks (see [5])

2 Related Work

ML for SITS Deep learning has been used for many Satellite Image Time Series (SITS) supervised-learning tasks such as crop-type mapping [29, 28, 37, 38], yield prediction [39, 40], understanding the economy [41, 42, 43, 44], precipitation forcasting [45], and land-cover classification [46, 47, 48, 27]. These works establish the usefulness of tailoring architectures such as LSTMs, self-attention, and transformers to temporal data. However, outside of their specific task, they are often not directly applicable to other remote-sensing datasets.

SSL for Satellite Imagery Self-supervised learning [2, 3, 4, 5, 6] has emerged as a promising approach in remote sensing domains. For instance, [34] and [35] propose incorporating spatially aligned images over time for contrastive self-supervised learning. Despite promising results, these two contrastive learning approaches rely heavily on the quality of positive pairs, which is often hard to control. [49] combines different sensor channels to generate co-located images that serve as positive pairs. [50, 51, 52] apply off-the-shelf contrastive learning algorithms to satellite images. [52] utilizes image inpainting and transformation prediction as additional pretext tasks. [53] leverages geographical knowledge to aid SSL, which, however, can be difficult to obtain as annotations.

Masked Autoencoder MAE [1] is a recent powerful self-supervised learning method. Instead of constructing a contrastive objective, it proposes the pretext task of reconstructing masked patches of the input, and largely avoids the need for designing specific data augmentation. Inspired by MAE's state-of-the-art performance on a wide collection of vision benchmarks [1], many follow-up works extend MAE to different data modalities. VideoMAE [54] proposes video tube masking and reconstruction as a pretext task for video analysis. GMAE [55] adapts MAE to the domain of graphs. MultiMAE [56] takes optional inputs of different modalities and accordingly includes other training objectives to facilitate multi-modality learning. However, these works fail to optimally handle temporal and multi-spectral input. VideoMAE requires equally-spaced image frames in the temporal dimension, which is not the case for satellite data given the temporal irregularity and discontinuity in sampling images of a location. In this work, we incorporate temporal and spectral information into a masked autoencoder architecture, and propose a novel self-supervised framework for satellite data.

3 Background

Masked Autoencoder The MAE is an autoencoder with asymmetrical encoding and decoding stages [1]. It operates on images $I \in \mathbb{R}^{C \times H \times W}$, where H,W are the height and width of the image, respectively, and C is the number of channels. The input image I is resized to a sequence of non-overlapping patches, $S \in \mathbb{R}^{L \times P^2C}$, where P is the height and width of the patch, and $L = (H/P) \cdot (W/P)$ is the number of patches. Each patch is passed through a patch embedding $f_p : \mathbb{R}^{P^2C} \mapsto \mathbb{R}^D$ to create a sequence $S' \in \mathbb{R}^{L \times D}$ of embedded patch "tokens". A fraction p_m of the L tokens are masked and only the remaining $(1-p_m)L$ "visible" patch tokens are fed to the encoder, a Vision Transformer (ViT) [36] with positional embeddings to capture the spatial location of the patch in the image. The decoder is a series of transformer blocks that operates on all L tokens (with positional embeddings added), where the p_mL encoded visible patches are placed in their original sequence position among $(1-p_m)L$ masked patches represented by a learnable mask token. The decoder outputs a reconstructed image $\hat{I} \in \mathbb{R}^{C \times H \times W}$, which is compared to the original image using the mean-squared error (MSE) loss, computed per-pixel only on the masked patches [1].

Positional encoding Positional encoding allows transformers to make their learned representations position-aware. In MAE [1] and in many transformers [57, 58], the positional encoding is:

$$\operatorname{Encode}(k,2i) = \sin\frac{k}{\Omega^{\frac{2i}{d}}}, \ \operatorname{Encode}(k,2i+1) = \cos\frac{k}{\Omega^{\frac{2i}{d}}} \tag{1}$$

Here, k is the position, i is the index of feature dimension in the encoding, d is the number of possible positions, and Ω is a large constant (normally set to 10000). In MAE, position is defined as the index of the patch along the x or y axes. Therefore, k ranges from 0 to H/P (or W/P). The final encoding is generated by concatenating the encodings of the x and y coordinates.

4 Method

In this section, we describe SatMAE with temporal (4.1) and multi-spectral (4.2) satellite images.

4.1 Temporal SatMAE

We now consider input tensors $I_T \in \mathbb{R}^{T \times C \times H \times W}$, where T denotes the number of images in a temporal sequence. In video data, T frames are usually equally spaced. However, temporal satellite imagery rarely has images at regular intervals. More commonly, several snapshots, or versions, of a given location are taken at irregular times. The length and sample frequency of these sequences of satellite images vary drastically over years and across different regions.

Naïvely, one could reshape I_T to $I_T' \in \mathbb{R}^{TC \times H \times W}$, effectively concatenating the temporal sequence of images along the spectral (i.e. channel) dimension, and then apply the MAE machinery verbatim. This method poses a few difficulties: (i) the model may be unable to generalise to a temporal ordering different to the one used in pre-training, since it can only understand order through the position of images in the stacked-timeseries (ii) the model cannot reason about the length of time separating two consecutive images in a time sequence, which may be variable when images of a location are

sampled at irregular intervals (iii) the model loses access to temporal fine-grained information in deeper layers, as its only direct exposure to encode temporal information is through the initial patch embedding f_p (iv) the model is not temporally-shift invariant (i.e. the model would need to separately learn to detect the same event in two different segments of a temporal sequence).

To address these challenges and to avoid losing temporal information, we resize the temporal sequence I_T to $S_T \in \mathbb{R}^{L_T \times P_T P^2 C}$, where $L_T = L \cdot (T/P_T) = (H/P) \cdot (W/P) \cdot (T/P_T)$, P_T is the "patch size" in the temporal dimension, and L and P are defined in 3. Prior works using transformers for video data suggest using $P_T = 2$, where each "patch" is a cube of shape $2 \times 16 \times 16$ [54, 59, 60]. Since our data has much shorter temporal sequence lengths [17], we let $P_T = 1$ such that $L_T = L \cdot T$. In order to operate on inputs of any temporal order, we re-use the same patch embedding $f_p : \mathbb{R}^{P^2C} \mapsto \mathbb{R}^D$ for each image in the time series, giving us an embedded sequence of tokens $S_T' \in \mathbb{R}^{L_T \times D}$.

4.1.1 Temporal Encoding

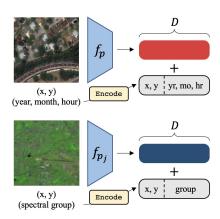


Figure 2: Top: Encoding each temporal patch with a shared patch embedding f_p . Bottom: Encoding each spectral patch with a different patch embedding f_{p_j} for each group j.

For each embedded token in the L_T length sequence, we need to ensure the model retains information about its spatial and temporal position. As shown in many prior works [34, 35], the timestamp of a satellite image is useful for many pre-training or downstream vision tasks. We propose a temporal encoding scheme compatible with the masked autoencoder architecture by treating the temporal dimension similarly to the positional dimensions (see 3).

The timestamp of a satellite image is represented as "year-month-day-hour-minute-second". Instead of passing the entire numerized timestamp into a feature encoder, we propose only keeping the useful parts. Intuitively, the day, minute, and second should be unrelated to the visual appearance of a region. Thus, including these components in the temporal encoding may not be beneficial, and can even be detrimental. In contrast, a landscape may evolve over years due to weather, geology, and human activity. The month reflects season and climate, and the hour reflects daylight and temperature.

Then, the temporal encoding is formulated as:

$$t_{k,i} = \text{CONCAT}[\text{Encode}(k_{\text{year}}, i), \text{Encode}(k_{\text{month}}, i), \text{Encode}(k_{\text{hour}}, i)]$$
(2)

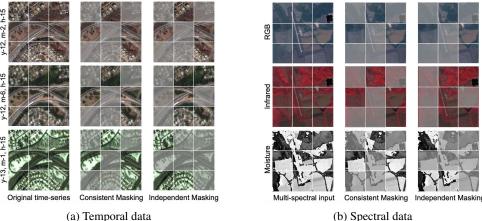
And the final encoding is generated by concatenating the temporal encoding to the positional encoding defined in 3 such that the total length of the encoding is D.

4.1.2 Masking Strategies

With an additional temporal dimension, masking a subset of the L_T tokens needs to be treated with care. As seen in figure 3, there are different ways to mask a temporal series of satellite images.

Consistent Masking Each image is "patchified" separately, but the masked regions are consistent across all images (fig. 3a). This approach is also used in VideoMAE [54], with video input.

Independent Masking Each image is "patchified" separately, and masked regions may not be the same across every image. Instead, a fraction p_m of the full sequence of all patch tokens are masked. Another variant is to independently mask the regions of each image, but keep the ratio p_m of masked regions fixed per image. Both variants are equivalent in expectation. Effectively, the model may look at unmasked values of a region that is masked in one image but not in others. This setting may lead to an easier task for video data since the model can "cheat" and exploit temporal redundancy in videos with high framerates [54]. However, we argue that this form of "cheating" is less feasible in temporal satellite imagery, given the strong impact of seasonal variation and changing human activity over periods of time and the much larger time deltas between temporally consecutive images (see fig. 3a).



(a) Temporal data

Figure 3: 3a Temporal masking: For images in a timeseries, we can choose to keep a patch fully visible or fully masked across time (consistent masking), or independently mask all patches (independent masking). In both cases, a fraction p_m patches are masked. Here, T=3, and the leftmost column orders the temporal sequence according to the timestamp features. For example, "y-12, m-12, h-15" is 12 years from the minimum year (2002), the zero-indexed month 2, and the 15th hour of the day; i.e., roughly 2014, March, 15:00. 3b Spectral Masking: The same masking strategies are adapted to groups of the 13 spectral bands in Sentinel-2 images.

Independent Masking + Inconsistent Cropping During data pre-processing, we can crop square regions for input inconsistently so that images in the same temporal sequence may be spatiallyunaligned. This strategy may help the model learn better representations as it may learn to align images in the sequence across the spatial and temporal dimensions.

4.2 Multi-spectral SatMAE

While MAE does operate on images $I \in \mathbb{R}^{C \times H \times W}$, usually C = 3 for RGB images. Satellite data, on the other hand, can often have multiple spectral bands. For example, Sentinel-2 imagery has C=13 bands of 10m, 20m and 60m spatial resolution, each of different wavelengths (see A.2.2). Below, we discuss and later experimentally compare various ways to encode spectral information.

Stack Channels The sequence of patches $S \in \mathbb{R}^{L \times P^2C}$ is embedded to a sequence of tokens $S' \in \mathbb{R}^{L \times D}$, thus treating the multi-band image as is. We denote this method **SatMAE+Stack**.

Group Channels There are limitations to naively stacking the spectral information, especially that a single convolutional patch embedding may be insufficient to fully capture fine-grained information present in multiple bands of different wavelengths and spatial resolution. We would like the model to preserve information about the different bands through the encoding and decoding stages.

To address this limitation, we propose grouping subsets of spectral bands. Given C channels, we To address this limitation, we propose grouping subsets of spectral bands. Given C channels, we form G groups g_1,g_2,\ldots,g_G such that $g_1+g_2+\cdots+g_G=C$. This is analogous to slicing the image I in the channel dimension, creating images I_1,\ldots,I_G , where $I_j\in\mathbb{R}^{g_j\times H\times W}$. We use a separate patch embedding $f_{p_j}:\mathbb{R}^{P^2g_j}\mapsto\mathbb{R}^D$ for each group j, thus allowing the model to best represent each possibly different group of channels as token embeddings. Therefore, each group j is first resized from $I_j\in\mathbb{R}^{g_j\times H\times W}$ to $S_j\in\mathbb{R}^{L\times P^2g_j}$, and then each patch is embedded with f_{p_j} to produce a sequence of embedded tokens $S_j'\in\mathbb{R}^{L\times D}$. The sequences S_1',\ldots,S_G' are concatenated to produce the final set of tokens $S'\in\mathbb{R}^{GL\times D}$.

Spectral Encoding Since the tokens in S' correspond to a patch location (m, n) in the input image and a group of channels g_j , we include an encoding for the group index k_g similar to 4.1.1

$$g_{k_g,i} = \text{Encode}(k_g, i) \tag{3}$$

Note that this encoding simply depends on a user-devised channel grouping, and differs from eq. (2) since additional metadata for the imagery, like its date, is not needed. The final encoding is a

concatenation of the positional $x_{k,i}$, $y_{k,i}$ and the spectral encoding $g_{k,i}$ such that the total dimension is D (see fig. 2). This positional encoding is added to S' before inputting it to the encoder. We denote the combined setting of grouping channels and using a group encoding as **SatMAE+Group**.

Masking Strategies We consider *consistent masking* (denoted SatMAE+Group+CM) and *independent masking* (SatMAE+Group+IM) as defined in section 4.1.2 and as visualized in fig. 3b.

5 Experiments

In this section, we first introduce the datasets we considered, including a new multi-spectral remote sensing image dataset for downstream task evaluation (5.1). We then present our results on benchmark datasets (5.2, 5.3, 5.4) and various remote sensing transfer-learning and downstream tasks 5.5. For all experiments, we compare with the current state-of-the-art methods [34, 35] and with supervised learning from scratch using the ViT backbone of SatMAE. In summary, our approach demonstrates strong performance on all the tasks we considered, yielding improvements over previous state-of-the-art techniques by up to 6% on supervised learning benchmarks, and up to 14% on remote sensing transfer-learning downstream remote sensing tasks.

5.1 Datasets for Pre-training

fMoW RGB Functional Map of the World (fMoW) [17] is a dataset of high-resolution satellite image time series across the world, with a task of classification among 62 categories.

fMoW Sentinel We create a new dataset based on the fMoW RGB dataset. We collect all 13 frequency bands provided by Sentinel-2 (B1-12 and B8A) for the original fMoW locations, at some of the same times as fMoW images plus some extra times, for a total of 712,874 training images, 84,939 validation images, and 84,966 test images. More details are included in appendix A.1.

5.2 fMoW RGB (non-temporal)

Method	Backbone	Frozen/Finetune
Sup.*	ResNet50	-/69.05
Sup.†	ResNet50	-/69.07
GASSL [34]	ResNet50	68.32/71.55
Sup.*	ViT-Large	-/62.48
Sup.†	ViT-Large	-/75.70
Sup.‡	ViT-Large	-/76.91
SatMAE	ViT-Large	65.94/ 77.84

Table 1: Top 1 Accuracy on fMoW classification. **Frozen**: only performing linear classification on frozen features of the pre-trained model. **Finetune**: end-to-end finetuning the whole model. * is training from scratch, and † is using supervised-learning ImageNet weights, and ‡ is SSL MAE ImageNet weights.

In this section, we perform experiments on fMoW single image classification task. Following [34], we report both the performance of linear probing and finetuning setting. Table 1 shows that compared to the previous stateof-the-art self-supervised method using a contrastive momentum encoding approach [34, 3], our SatMAE achieved a 6.29% improvement in top 1 classification accuracy. Interestingly, without SatMAE pre-training the ViT-large model could only reach 62.48% at convergence after 50 epochs of finetuning compared to 69.05% achieved by training a ResNet-50 model from scratch. This is likely because the ViT [36] backbone is harder to finetune from scratch than ResNet50 [61], which makes the pre-trained model more valuable.

5.3 fMoW RGB (temporal)

Main experiments We perform image-sequence classification on the temporal version of fMoW RGB to evaluate our temporal SatMAE. The temporal fMoW consists of co-located image sequences with a length of 3. As seen in table 2, SatMAE surpasses the previous state-of-the-art by 4.48% and improves the non-temporal result by 2.06% in top 1 classification accuracy. We also outperform UTAE [48], a SITS state-of-the-art, by 18%. We can observe from rows 5-8 that this gain is not from the larger model to handle sequences of data. Naively stacking the image sequences in the channel dimension performs even worse than the non-temporal SatMAE. Again, SatMAE pre-training is crucial for ViT to outperform ResNet50. Training details are in appendix A.3.2.

Method	Backbone	Top Acc. (1/5)
Sup.*	ResNet50	73.24/-
SeCo [35]	ResNet50	66.80/-
GASSL [34]	ResNet50	74.11/-
UTAE [48]	U-Net	61.59/86.45
Sup.*	ViT-Large	61.89/84.23
SatMAE+Stack	ViT-Large	75.85/88.68
MAE+Test Aug.	ViT-Large	78.90/93.31
MAE	ViT-Large	76.78/92.01
SatMAE	ViT-Large	81.49/93.26

Table 2: Classification results on the temporal fMoW RGB dataset. * means finetuning from scratch. || means copying the input image 3 times instead of using temporal sequences as input. SatMAE+Stack here means stacking the image sequence along the channel space.

Temp. Enc.	Indep. Mask.	Cons. Crop.	Test Aug.	Top 1 Acc.
	√	√		78.07
$\overline{\hspace{1cm}}$		√		78.45
$\overline{\hspace{1cm}}$	√			79.90
$\overline{\hspace{1cm}}$	√	√		79.69
$\overline{\hspace{1cm}}$	√	√	√	81.49

Table 4: Ablation studies on different components of temporal SatMAE on the temporal fMoW classification task. The first column is whether using temporal encoding, the second is whether using independent masking, the third is whether cropping consistently, and the last one is whether applying test-time augmentation.

Method	Backbone	Top Acc. (1/5)
Sup. Learning*	ResNet152	49.12/75.73
Sup. Learning‡	ResNet152	54.46/78.99
MoCo-v3	ViT-Base	50.45/76.37
MoCo-v3+Group	ViT-Base	51.33/75.68
SatMAE+Group*	ViT-Large	53.03/77.14
SatMAE+Group†	ViT-Large	51.61/77.26
SatMAE+Group‡	ViT-Large	47.57/72.26
SatMAE+Group§	ViT-Large	49.49/76.30
SatMAE+Stack	ViT-Large	57.37/81.63
SatMAE+Group+IM	ViT-Large	59.30/82.81
SatMAE+Group+IM	ViT-Large	61.48/85.17

Table 3: Top 1 & Top 5 Accuracy on the fMoW Sentinel validation set. The different initializations are: * from scratch, † MAE ImageNet weights, ‡ supervised ImageNet weights, § SatMAE fMoW RGB weights. Other rows use fMoW Sentinel for pre-training. The last row includes additional data augmentations (5.4).

Back.	Group Strat.	Indp. Mask.	Spec. Enc.	Top 1 Acc.
Base	X	√	√	59.11
Large	X	√		58.87
Large	X		√	57.76
Large	Н	√	√	57.78
Large	R	√	√	58.76
Large	X	√	√	59.30

Table 5: Ablation studies on spectral SatMAE on fMoW-Sentinel. The first column denotes using ViT-Base or ViT-Large. The second column is the grouping strategy (see 5.4). The third column denotes independent or consistent masking. The last column is whether the spectral group encoding 3 is used.

Ablation studies Table 4 provides a comprehensive ablation study on the components of temporal SatMAE. We see that improved performance is mainly due to the temporal encoding and adopting independent masking rather than the consistent masking strategy suggested in VideoMAE [54]. Interestingly, consistent cropping slightly decreases performance, indicating that the model does not rely on perfectly spatially-aligned image sequences. In addition, using test-time augmentations similar to [34] is beneficial. Further ablations on mask ratio p_m and patch size P are in appendix A.4.

5.4 fMoW Sentinel (Multi-spectral)

In this section, we pre-train and finetune SatMAE on the image classification task of the fMoW-Sentinel dataset. We pre-train SatMAE+Stack 4.2 and investigate SatMAE+Group+CM 4.1.2 and SatMAE+Group+IM 4.1.2, (see 4.2, 4.2). The full models are then finetuned on the fMoW-Sentinel image classification task. For comparison, we also finetune the ResNet-152 model [61] from scratch and from a supervised ImageNet initialization. We pick the largest model, ResNet-152, for fairer comparison with ViTs. We also include MoCo-v3 [62, 3], a popular SSL method. Given the differences in applying RGB-image augmentations to satellite imagery, we implement two versions: (i) MoCo-v3: we apply all of the same augmentations, except random grayscale and solarize, to create 2 views of the 10-channel image. (ii) MoCo-v3+Group: we split the 10 bands into two groups suggested by [2], and apply augmentations to each to create a positive pair of two 5-channel images.

Model configuration As not all of the 13 Sentinel-2 bands may be useful, in our experiments we drop bands B1, B9 and B10, which correspond to a spatial resolution of 60m. Of the remaining 10 bands, we form three groups: (i) RGB+NIR: B2, B3, B4, B8 (ii) Red Edge: B5, B6, B7, B8A (iii) SWIR: B11, B12. We choose this grouping to ensure each group has bands of the same spatial resolution and similar wavelength (see A.2.2, A.6). Only the last row of table 3 includes additional data augmentations used during finetuning as in [1]. See A.3.3 for pre-training and finetuning details.

Method	Backbone	Top 1 Acc.
Sup. (Scratch)	ResNet50	54.46
GASSL [34]	ResNet50	57.63
Sup. (Scratch)	ViT-Large	69.65
SatMAE	ViT-Large	71.77

Table 6: NAIP land cover classification results.

Method	Backbone	Top 1 Acc.
Sup. (Scratch)	ResNet18	63.21
Sup. (IN init.)	ResNet18	86.44
GASSL [34]	ResNet18	89.51
SeCo [35]	ResNet18	93.14
SatMAE*	ViT-Large	95.74
SatMAE	ViT-Large	98.94
SatMAE+Group+IM	ViT-Large	98.98

Table 8: EuroSAT land cover classification results. * means we only use the RGB channels of the data.

Method	Backbone	mIoU
Sup. (Scratch)	ResNet50	75.57
GASSL [34]	ResNet50	78.51
Sup. (Scratch)	ViT-Large	74.71
SatMAE	ViT-Large	78.07

Table 7: SpaceNet v1 building segmentation results.

Method	Backbone	mAP
Sup. (Scratch)	ResNet50	69.49
Sup. (IN init.)	ResNet50	80.04
GASSL [34]	ResNet50	80.20
SeCo [35]	ResNet50	82.62
Sup. (Scratch)	ViT-Large	80.07
SatMAE	ViT-Large	82.13

Table 9: BigEarthNet multi-label classification results. Following [35], we use mean Average Precision (mAP) as the metric, and use a newer set of class labels.

Results We present results in table 3. Our method SatMAE+Group+IM achieves the highest accuracy, outperforming supervised training from scratch (\uparrow 6.27%) and ImageNet-initialized backbones (\uparrow 4.84%). ImageNet initializations may be less useful than in fMoW-RGB given the larger distributional shift to multi-spectral input data. We also note the effectiveness of grouping channels over processing all bands only at the patch embedding level (i.e. SatMAE+Stack).

Ablation Studies We investigate the design of SatMAE for multi-spectral data in table 5. For grouping strategy, we implement alternate band groups to test the hypothesis that grouping bands based on wavelength and resolution is beneficial. X represents the band groups in 5.4. H represents splitting the 10 bands into two halves, {(2,3,4,5,6), (7,8,8A,11,12)}. R represents a random split into three groups {(6,5,11,12), (8A,4,8,3), (7,2)}, reflecting the same group sizes as X. As seen, the choice of band groups does influence performance, yielding a gain of about 0.6%. Moreover, ViT-Base performs strongly, suggesting that SatMAE is the reason for improved performance rather than the number of parameters in ViT. Interestingly, *independent masking* performs the best, which prompts the model to "peek" at unmasked band groups to reconstruct the same region in a masked band group.

We also include further experiments on the length of pre-training (see A.3.3), the impact of mask ratio p_m and patch size P (see A.5), and the usefulness of the 13 Sentinel-2 spectral bands (see A.6).

5.5 Transfer Learning Experiments

Now, we finetune our pre-trained SatMAE on downstream tasks on remote-sensing datasets, including land cover classification (5.5), multi-label classification (5.5), and building segmentation (5.5). Finetuning details are included in A.7, A.8, A.9, A.10.

Land Cover Classification We perform transfer learning experiments on land cover classification using the NAIP and EuroSAT [63] dataset. NAIP consists of RGB+CIR images of 66 land cover classes obtained by the USDA's National Agricultural Imagery Program, which are split into 244,471 training and 55,529 validation images. EuroSAT is a small dataset containing 27,000 13-band satellite images of 10 classes based on Sentinel-2. We follow [35, 64] for the train/val splits on EuroSAT.

Table 6 and table 8 shows the remarkable improvement of our SatMAE over the state-of-the-arts. Although using the ViT-Large backbone already achieved good results, initializing the model with SAT-MAE pre-trained weights further increased the accuracy by 2%-3%.

Multi-label Classification We also use the BigEarthNet [18] dataset for multi-label classification, which consists of 13-band Sentinel-2 images of 19 classes in total. There are 354,196 images for training and 118,065 images for validation. Following [35], we use a 10% subset of the train set.

Table 9 shows SatMAE pre-training improves upon the model trained from scratch by over 2%, and achieves comparable results to the state-of-the-art. GASSL and SeCo were actually trained on a larger pre-train dataset (1M Sentinel-2 images v.s. 713k) and with all 13 bands than our fMoW Sentinel. Therefore we expect further improvement when we pre-train SatMAE with more data and for longer.

Building Segmentation In this section, we evaluate SatMAE on the semantic segmentation down-stream task of the SpaceNet v1 dataset [20]. The SpaceNet v1 dataset consists of 6940 high resolution satellite images with segmentation masks for buildings, which are divided into train and test sets of 5000 and 1940 images, respectively.

The results in table 7 show that our method achieves a larger performance gain from supervised learning from scratch compared to [34]. The incompatibility of the ViT backbone with PSANet could explain why the baseline performance is not as strong as that of using a ResNet50 backbone.

5.6 Visualizing reconstruction quality for SatMAE

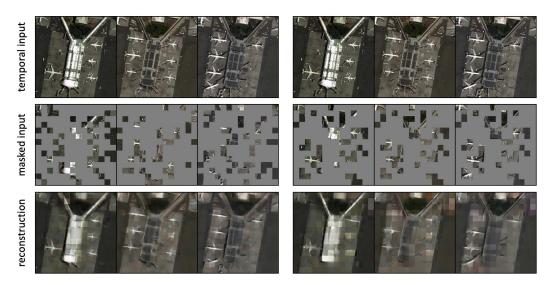


Figure 4: Reconstruction quality of SatMAE+IM (left) vs. SatMAE+CM (right). Further results in appendix C.

We show the visualization of the reconstruction quality of two different SatMAE masking strategies in fig. 4. SatMAE+IM successfully reconstructs all the airplanes even though their number varies across time. In contrast, the SatMAE with Consistent Masking missed some airplanes in the reconstruction.

6 Conclusion

In this paper, we propose a new SSL framework based on the MAE architecture [1] tailored to remote-sensing data (satellite imagery). Our novel masking strategy in a joint positional, temporal/spectral space, along with the temporal and spectral encoding, enables our model to handle temporal and multi-spectral satellite images as input and learn useful representations. Experiments on the datasets for pre-training and multiple downstream datasets demonstrate the effectiveness of our pre-trained SatMAE model, outperforming previous state-of-the-art results by large margins.

In the future, it would be useful to design more efficient transformer architectures. While SatMAE has a similar number of parameters for both the temporal and multi-spectral setting as a regular ViT, the increased length of token sequences can strain computational resources. Moreover, it is also worth exploring optimal positional encodings for spectral and temporal data, as well as optimal groups of spectral bands, either by neural-based search methods, or using prior knowledge. Lastly, investigating better architectures for object detection and semantic segmentation using ViTs will be important in generalising SatMAE to further downstream tasks.

Broader Impact

Accurate measurements of economic, social, and environmental phenomena are key inputs into policy decisions made around the world, but the sparsity of labelled data on many outcomes means that such decisions are often not guided by timely or accurate data. We demonstrate how a pre-training framework could relieve the dependence on labelled data for many downstream tasks that use satellite imagery as input. We hope our SatMAE method will help close the gap between SSL performance on natural imagery and on the more challenging satellite imagery, and prompt further attention from the ML community on the usefulness of SSL in satellite-imagery-related tasks.

Better extraction of information from satellite imagery has profound implications for our ability to measure and understand a broad array of social, economic and environmental phenomena that are critical for decision making. Our approach further amplifies the usefulness of the sparse amount of labelled data that exist on key human outcomes, and could enable rapid and accurate extraction of imagery features relevant for critical downstream tasks, including poverty prediction, infrastructure development, and population estimation. Such information could aid governments in more rapid and data-informed decision making and ultimately bring large societal benefits.

7 Acknowledgements

This research is based upon work supported in part by the Office of the Director of National Intelligence (ODNI), Intelligence Advanced Research Projects Activity (IARPA), via 2021-2011000004, HAI, NSF(#1651565), AFOSR (FA95501910024), ARO (W911NF-21-1-0125) and Sloan Fellowship. The views and conclusions contained herein are those of the authors and should not be interpreted as necessarily representing the official policies, either expressed or implied, of ODNI, IARPA, or the U.S. Government. The U.S. Government is authorized to reproduce and distribute reprints for governmental purposes not-withstanding any copyright annotation therein.

Checklist

- 1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper's contributions and scope? [Yes] We tried our best to be precise.
 - (b) Did you describe the limitations of your work? [Yes] We described the limitations in the conclusion section.
 - (c) Did you discuss any potential negative societal impacts of your work? [Yes] We discussed in Appendix B
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [Yes]
- 2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [N/A]
 - (b) Did you include complete proofs of all theoretical results? [N/A]
- 3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [Yes] We list the url to our project website in the abstract. The website will contain links to the code and data on Github.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [Yes] We list part of them in the experiments section and part of them in Appendix A.
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [No] We trained on large datasets with small variation across different runs expected. Limited by computation resources, we only ran each experiment once.
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [Yes] We list that in Appendix A.
- 4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [Yes] Yes, we do that by citing creators of datasets and authors of prior works.
 - (b) Did you mention the license of the assets? [Yes] We do that in Appendix A
 - (c) Did you include any new assets either in the supplemental material or as a URL? [Yes] We are releasing a new dataset. The instructions and links will be released in our codebase mentioned above.
 - (d) Did you discuss whether and how consent was obtained from people whose data you're using/curating? [Yes] All data we used are released publicly.
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [Yes] We concluded that no data we are using has such concerns.
- 5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [N/A]
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [N/A]
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [N/A]

References

[1] Kaiming He, Xinlei Chen, Saining Xie, Yanghao Li, Piotr Dollár, and Ross Girshick. Masked autoencoders are scalable vision learners. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 16000–16009, 2022.

- [2] Yonglong Tian, Dilip Krishnan, and Phillip Isola. Contrastive multiview coding. In *European conference on computer vision*, pages 776–794. Springer, 2020.
- [3] Kaiming He, Haoqi Fan, Yuxin Wu, Saining Xie, and Ross Girshick. Momentum contrast for unsupervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 9729–9738, 2020.
- [4] Ting Chen, Simon Kornblith, Mohammad Norouzi, and Geoffrey Hinton. A simple framework for contrastive learning of visual representations. In *International conference on machine learning*, pages 1597–1607. PMLR, 2020.
- [5] Jean-Bastien Grill, Florian Strub, Florent Altché, Corentin Tallec, Pierre Richemond, Elena Buchatskaya, Carl Doersch, Bernardo Avila Pires, Zhaohan Guo, Mohammad Gheshlaghi Azar, et al. Bootstrap your own latent-a new approach to self-supervised learning. *Advances in Neural Information Processing Systems*, 33:21271–21284, 2020.
- [6] Xinlei Chen and Kaiming He. Exploring simple siamese representation learning. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 15750–15758, 2021.
- [7] Ashish Jaiswal, Ashwin Ramesh Babu, Mohammad Zaki Zadeh, Debapriya Banerjee, and Fillia Makedon. A survey on contrastive self-supervised learning. *CoRR*, abs/2011.00362, 2020.
- [8] Longlong Jing and Yingli Tian. Self-supervised visual feature learning with deep neural networks: A survey. IEEE transactions on pattern analysis and machine intelligence, 43(11):4037– 4058, 2020.
- [9] Ishan Misra and Laurens van der Maaten. Self-supervised learning of pretext-invariant representations. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 6707–6717, 2020.
- [10] Mathilde Caron, Ishan Misra, Julien Mairal, Priya Goyal, Piotr Bojanowski, and Armand Joulin. Unsupervised learning of visual features by contrasting cluster assignments. *Advances in Neural Information Processing Systems*, 33:9912–9924, 2020.
- [11] Yude Wang, Jie Zhang, Meina Kan, Shiguang Shan, and Xilin Chen. Self-supervised equivariant attention mechanism for weakly supervised semantic segmentation. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 12275–12284, 2020.
- [12] Yann Dubois, Benjamin Bloem-Reddy, Karen Ullrich, and Chris J Maddison. Lossy compression for lossless prediction. *Advances in Neural Information Processing Systems*, 34, 2021.
- [13] Alexander Kolesnikov, Xiaohua Zhai, and Lucas Beyer. Revisiting self-supervised visual representation learning. In *Proceedings of the IEEE/CVF conference on computer vision and pattern recognition*, pages 1920–1929, 2019.
- [14] Zeyu Feng, Chang Xu, and Dacheng Tao. Self-supervised representation learning by rotation feature decoupling. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 10364–10374, 2019.
- [15] Soroush Abbasi Koohpayegani, Ajinkya Tejankar, and Hamed Pirsiavash. Compress: Self-supervised learning by compressing representations. Advances in Neural Information Processing Systems, 33:12980–12992, 2020.
- [16] Abhishek Sinha, Jiaming Song, Chenlin Meng, and Stefano Ermon. D2c: Diffusion-decoding models for few-shot conditional generation. Advances in Neural Information Processing Systems, 34:12533–12548, 2021.
- [17] Gordon Christie, Neil Fendley, James Wilson, and Ryan Mukherjee. Functional map of the world. In CVPR, 2018.
- [18] Gencer Sumbul, Marcela Charfuelan, Begüm Demir, and Volker Markl. Bigearthnet: A large-scale benchmark archive for remote sensing image understanding. In *IGARSS* 2019-2019 IEEE International Geoscience and Remote Sensing Symposium, pages 5901–5904. IEEE, 2019.

- [19] Darius Lam, Richard Kuzma, Kevin McGee, Samuel Dooley, Michael Laielli, Matthew Klaric, Yaroslav Bulatov, and Brendan McCord. xview: Objects in context in overhead imagery. *arXiv* preprint arXiv:1802.07856, 2018.
- [20] Adam Van Etten, Dave Lindenbaum, and Todd M Bacastow. Spacenet: A remote sensing dataset and challenge series. *arXiv preprint arXiv:1807.01232*, 2018.
- [21] Marshall Burke, Anne Driscoll, David B Lobell, and Stefano Ermon. Using satellite imagery to understand and promote sustainable development. *Science*, 371(6535):eabe8628, 2021.
- [22] Kumar Ayush, Burak Uzkent, Kumar Tanmay, Marshall Burke, David Lobell, and Stefano Ermon. Efficient poverty mapping from high resolution remote sensing images. In *Proc. AAAI Conf. Artif. Intell*, volume 35, pages 12–20, 2021.
- [23] Kumar Ayush, Burak Uzkent, Marshall Burke, David Lobell, and Stefano Ermon. Generating interpretable poverty maps using object detection in satellite images. *arXiv* preprint *arXiv*:2002.01612, 2020.
- [24] Neal Jean, Marshall Burke, Michael Xie, W Matthew Davis, David B Lobell, and Stefano Ermon. Combining satellite imagery and machine learning to predict poverty. *Science*, 353(6301):790–794, 2016.
- [25] Jiaxuan You, Xiaocheng Li, Melvin Low, David Lobell, and Stefano Ermon. Deep gaussian process for crop yield prediction based on remote sensing data. In *Thirty-First AAAI conference* on artificial intelligence, 2017.
- [26] Anna X Wang, Caelin Tran, Nikhil Desai, David Lobell, and Stefano Ermon. Deep transfer learning for crop yield prediction with remote sensing data. In *Proceedings of the 1st ACM SIGCAS Conference on Computing and Sustainable Societies*, pages 1–5, 2018.
- [27] Marc Rußwurm and Marco Körner. Self-attention for raw optical satellite time series classification. *ISPRS journal of photogrammetry and remote sensing*, 169:421–435, 2020.
- [28] Jorge Andres Chamorro Martinez, Laura Elena Cué La Rosa, Raul Queiroz Feitosa, Ieda Del'Arco Sanches, and Patrick Nigri Happ. Fully convolutional recurrent networks for multidate crop recognition from multitemporal image sequences. *ISPRS Journal of Photogrammetry and Remote Sensing*, 171:188–201, 2021.
- [29] Rose M Rustowicz, Robin Cheong, Lijing Wang, Stefano Ermon, Marshall Burke, and David Lobell. Semantic segmentation of crop type in africa: A novel dataset and analysis of deep learning methods. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition Workshops*, pages 75–82, 2019.
- [30] Christopher Yeh, Chenlin Meng, Sherrie Wang, Anne Driscoll, Erik Rozi, Patrick Liu, Jihyeon Lee, Marshall Burke, David B Lobell, and Stefano Ermon. Sustainbench: Benchmarks for monitoring the sustainable development goals with machine learning. *arXiv preprint arXiv:2111.04724*, 2021.
- [31] Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. Imagenet: A large-scale hierarchical image database. In 2009 IEEE Conference on Computer Vision and Pattern Recognition, pages 248–255, 2009.
- [32] Yutong He, Dingjie Wang, Nicholas Lai, William Zhang, Chenlin Meng, Marshall Burke, David Lobell, and Stefano Ermon. Spatial-temporal super-resolution of satellite imagery via conditional pixel synthesis. *Advances in Neural Information Processing Systems*, 34:27903–27915, 2021.
- [33] Yutong He, William Zhang, Chenlin Meng, Marshall Burke, David B Lobell, and Stefano Ermon. Tracking urbanization in developing regions with remote sensing spatial-temporal super-resolution. *arXiv* preprint arXiv:2204.01736, 2022.
- [34] Kumar Ayush, Burak Uzkent, Chenlin Meng, Kumar Tanmay, Marshall Burke, David Lobell, and Stefano Ermon. Geography-aware self-supervised learning. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 10181–10190, 2021.

- [35] Oscar Mañas, Alexandre Lacoste, Xavier Giro-i Nieto, David Vazquez, and Pau Rodriguez. Seasonal contrast: Unsupervised pre-training from uncurated remote sensing data. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 9414–9423, 2021.
- [36] Alexey Dosovitskiy, Lucas Beyer, Alexander Kolesnikov, Dirk Weissenborn, Xiaohua Zhai, Thomas Unterthiner, Mostafa Dehghani, Matthias Minderer, Georg Heigold, Sylvain Gelly, Jakob Uszkoreit, and Neil Houlsby. An image is worth 16x16 words: Transformers for image recognition at scale. *ICLR*, 2021.
- [37] Marc Rußwurm and Marco Körner. Convolutional lstms for cloud-robust segmentation of remote sensing imagery. *arXiv preprint arXiv:1811.02471*, 2018.
- [38] Charlotte Pelletier, Geoffrey I Webb, and François Petitjean. Temporal convolutional neural network for the classification of satellite image time series. *Remote Sensing*, 11(5):523, 2019.
- [39] Elisa Kamir, François Waldner, and Zvi Hochman. Estimating wheat yields in australia using climate records, satellite image time series and machine learning methods. ISPRS Journal of Photogrammetry and Remote Sensing, 160:124–135, 2020.
- [40] Raí A Schwalbert, Telmo Amado, Geomar Corassa, Luan Pierre Pott, PV Vara Prasad, and Ignacio A Ciampitti. Satellite-based soybean yield forecast: Integrating machine learning and weather data for improving crop yield prediction in southern brazil. Agricultural and Forest Meteorology, 284:107886, 2020.
- [41] Chenlin Meng, Enci Liu, Willie Neiswanger, Jiaming Song, Marshall Burke, David Lobell, and Stefano Ermon. Is-count: Large-scale object counting from satellite images with covariate-based importance sampling. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 36, pages 12034–12042, 2022.
- [42] Evan Sheehan, Chenlin Meng, Matthew Tan, Burak Uzkent, Neal Jean, Marshall Burke, David Lobell, and Stefano Ermon. Predicting economic development using geolocated wikipedia articles. In *Proceedings of the 25th ACM SIGKDD international conference on knowledge discovery & data mining*, pages 2698–2706, 2019.
- [43] Enci Liu, Chenlin Meng, Matthew Kolodner, Eun Jee Sung, Sihang Chen, Marshall Burke, David Lobell, and Stefano Ermon. Building coverage estimation with low-resolution remote sensing imagery. *arXiv preprint arXiv:2301.01449*, 2023.
- [44] Burak Uzkent, Evan Sheehan, Chenlin Meng, Zhongyi Tang, Marshall Burke, David Lobell, and Stefano Ermon. Learning to interpret satellite images using wikipedia. In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence*, 2019.
- [45] Xingjian Shi, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting. *Advances in neural information processing systems*, 28, 2015.
- [46] Andrei Stoian, Vincent Poulain, Jordi Inglada, Victor Poughon, and Dawa Derksen. Land cover maps production with high resolution satellite image time series and convolutional neural networks: Adaptations and limits for operational systems. *Remote Sensing*, 11(17):1986, 2019.
- [47] Xin Yang and CP Lo. Using a time series of satellite imagery to detect land use and land cover changes in the atlanta, georgia metropolitan area. *International Journal of Remote Sensing*, 23(9):1775–1798, 2002.
- [48] Vivien Sainte Fare Garnot and Loic Landrieu. Panoptic segmentation of satellite image time series with convolutional temporal attention networks. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pages 4872–4881, 2021.
- [49] Aidan M Swope, Xander H Rudelis, and Kyle T Story. Representation learning for remote sensing: An unsupervised sensor fusion approach. *arXiv preprint arXiv:2108.05094*, 2021.
- [50] Vladan Stojnic and Vladimir Risojevic. Self-supervised learning of remote sensing scene representations using contrastive multiview coding. In *Proceedings of the IEEE/CVF Conference on Computer Vision and Pattern Recognition*, pages 1182–1191, 2021.

- [51] Pallavi Jain, Bianca Schoen-Phelan, and Robert Ross. Multi-modal self-supervised representation learning for earth observation. In 2021 IEEE International Geoscience and Remote Sensing Symposium IGARSS, pages 3241–3244. IEEE, 2021.
- [52] Wenyuan Li, Hao Chen, and Zhenwei Shi. Semantic segmentation of remote sensing images with self-supervised multitask representation learning. *IEEE Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 14:6438–6450, 2021.
- [53] Wenyuan Li, Keyan Chen, Hao Chen, and Zhenwei Shi. Geographical knowledge-driven representation learning for remote sensing images. *IEEE Transactions on Geoscience and Remote Sensing*, 2021.
- [54] Zhan Tong, Yibing Song, Jue Wang, and Limin Wang. Videomae: Masked autoencoders are data-efficient learners for self-supervised video pre-training. arXiv preprint arXiv:2203.12602, 2022.
- [55] Hongxu Chen, Sixiao Zhang, and Guandong Xu. Graph masked autoencoder. arXiv preprint arXiv:2202.08391, 2022.
- [56] Roman Bachmann, David Mizrahi, Andrei Atanov, and Amir Zamir. Multimae: Multi-modal multi-task masked autoencoders. arXiv preprint arXiv:2204.01678, 2022.
- [57] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *Advances in neural information processing systems*, 30, 2017.
- [58] Hugo Touvron, Matthieu Cord, Matthijs Douze, Francisco Massa, Alexandre Sablayrolles, and Hervé Jégou. Training data-efficient image transformers & distillation through attention. In International Conference on Machine Learning, pages 10347–10357. PMLR, 2021.
- [59] Anurag Arnab, Mostafa Dehghani, Georg Heigold, Chen Sun, Mario Lučić, and Cordelia Schmid. Vivit: A video vision transformer. In *International Conference on Computer Vision* (ICCV), 2021.
- [60] Ze Liu, Jia Ning, Yue Cao, Yixuan Wei, Zheng Zhang, Stephen Lin, and Han Hu. Video swin transformer. *arXiv preprint arXiv:2106.13230*, 2021.
- [61] Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 770–778, 2016.
- [62] Xinlei Chen*, Saining Xie*, and Kaiming He. An empirical study of training self-supervised vision transformers. *arXiv* preprint arXiv:2104.02057, 2021.
- [63] Patrick Helber, Benjamin Bischke, Andreas Dengel, and Damian Borth. Eurosat: A novel dataset and deep learning benchmark for land use and land cover classification. *IEEE Journal* of Selected Topics in Applied Earth Observations and Remote Sensing, 12(7):2217–2226, 2019.
- [64] Maxim Neumann, Andre Susano Pinto, Xiaohua Zhai, and Neil Houlsby. In-domain representation learning for remote sensing. arXiv preprint arXiv:1911.06721, 2019.
- [65] Hengshuang Zhao, Yi Zhang, Shu Liu, Jianping Shi, Chen Change Loy, Dahua Lin, and Jiaya Jia. Psanet: Point-wise spatial attention network for scene parsing. In *Proceedings of the European conference on computer vision (ECCV)*, pages 267–283, 2018.
- [66] Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. Quantifying the carbon emissions of machine learning. CoRR, abs/1910.09700, 2019.