

# Panoptic Segmentation of Satellite Image Time Series with Convolutional Temporal Attention Networks

Vivien Sainte Fare Garnot

Loic Landrieu

LASTIG, Univ. Gustave Eiffel, ENSG, IGN, F-94160 Saint-Mande, France

{vivien.sainte-fare-garnot, loic.landrieu}@ign.fr

## Abstract

*Unprecedented access to multi-temporal satellite imagery has opened new perspectives for a variety of Earth observation tasks. Among them, pixel-precise panoptic segmentation of agricultural parcels has major economic and environmental implications. While researchers have explored this problem for single images, we argue that the complex temporal patterns of crop phenology are better addressed with temporal sequences of images. In this paper, we present the first end-to-end, single-stage method for panoptic segmentation of Satellite Image Time Series (SITS). This module can be combined with our novel image sequence encoding network which relies on temporal self-attention to extract rich and adaptive multi-scale spatio-temporal features. We also introduce PASTIS, the first open-access SITS dataset with panoptic annotations. We demonstrate the superiority of our encoder for semantic segmentation against multiple competing architectures, and set up the first state-of-the-art of panoptic segmentation of SITS. Our [implementation](#) and [PASTIS](#) are publicly available.*

## 1. Introduction

The precision and availability of Earth observations have continuously improved thanks to sustained advances in space-based remote sensing, such as the launch of the Planet [5] and the open-access Sentinel constellations [8]. In particular, satellites with high revisit frequency contribute to a better understanding of phenomena with complex temporal dynamics. Crop mapping—the driving application of this paper—relies on exploiting such temporal patterns [38] and entails major financial and environmental stakes. Indeed, remote monitoring of the surface and nature of agricultural parcels is necessary for a fair allocation of agricultural subsidies (50 and 22 billion euros per year in Europe and in the US, respectively) and for ensuring that best crop rotation practices are respected. More generally, the automated analysis of SITS represents a significant interest for a wide

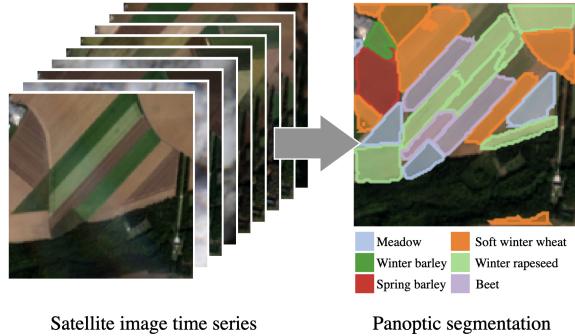


Figure 1: **Overview.** We propose an end-to-end, single-stage model for panoptic segmentation of agricultural parcels from time series of satellite images. Note the difficulty of resolving the parcels’ borders from a single image, highlighting the need for modeling temporal dynamics.

range of applications, such as surveying urban development and deforestation.

The task of monitoring both the content and extent of agricultural parcels can be framed as the panoptic segmentation of an image sequence. Panoptic segmentation consists of assigning to each pixel a class and a unique instance label, and has become a standard visual perception task in computer vision [19, 26]. However, panoptic segmentation is a fundamentally different task for SITS versus sequences of natural images or videos. Indeed, understanding videos requires tracking objects through time and space [44]. In yearly SITS, the targets are static in a geo-referenced frame, which removes the need for spatial tracking. Additionally, SITS share a common temporal frame of reference, which means that the time of acquisition itself contains information useful for modeling the underlying temporal dynamics. In contrast, the frame number in videos is often arbitrary. Finally, while objects on the Earth surface generally do not occlude one another, as is commonly the case for objects in natural images, varying cloud cover can make the analysis of SITS arduous. For the specific problem addressed in this

paper, individualizing agricultural parcels requires learning complex and specific temporal, spatial, and spectral patterns not commonly encountered in video processing, such as differences in plant phenological profiles, subpixel border information, and swift human interventions such as harvests.

While deep networks have proven efficient for learning such complex patterns for pixel classification [16, 12, 1], there is no dedicated approach for detecting individual objects in SITS. Existing work on instance segmentation has been restricted to analysing a single satellite image [33]. In summary, specialized remote sensing methods are limited to semantic segmentation or single-image instance segmentation, while computer vision’s panoptic-ready networks require significant adaptation to be applied to SITS.

In this paper, we introduce U-TAE (U-net with Temporal Attention Encoder), a novel spatio-temporal encoder combining multi-scale spatial convolutions [34] and a temporal self-attention mechanism [38] which learns to focus on the most salient acquisitions. While convolutional-recurrent methods are limited to extracting temporal features at the highest [35] or lowest [37] spatial resolutions, our proposed method can use the predicted temporal masks to extract specialized and adaptive spatio-temporal features at different resolutions simultaneously. We also propose Parcels-as-Points (PaPs), the first end-to-end deep learning method for panoptic segmentation of SITS. Our approach is built upon the efficient CenterMask network [49], which we modify to fit our problem. Lastly, we present Panoptic Agricultural Satellite TIIme-Series (PASTIS), the first open-access dataset for training and evaluating panoptic segmentation models on SITS, with over 2 billion annotated pixels covering over 4000km<sup>2</sup>. Evaluated on this dataset, our approach outperforms all reimplemented competing methods for semantic segmentation, and defines the first state-of-the-art of SITS panoptic segmentation.

## 2. Related Work

To the best of our knowledge, no instance or panoptic segmentation method operating on SITS has been proposed to date. However, there is a large body of work on both the encoding of satellite sequences, and the panoptic segmentation of videos and single satellite images.

**Encoding Satellite Image Sequences.** While the first automated tools for SITS analysis relied on traditional machine learning [13, 46], deep convolutional networks allow for the extraction of richer spatial descriptors [20, 12, 1, 16]. The temporal dimension was initially dealt via handcrafted temporal descriptors [2, 43, 52] or probabilistic models [3], which have been advantageously replaced by recurrent [35, 38, 28], convolutional [30, 37, 15], or differential [25] architectures. Recently, attention-based approaches have

been adapted to encode sequences of remote sensing images and have led to significant progress for pixel-wise and parcel-wise classification [39, 36, 54]. In parallel, hybrid architectures [42, 37, 29] relying on U-Net-type architectures [34] for encoding the spatial dimension and recurrent networks for the temporal dimension have shown to be well suited for the semantic segmentation of SITS. In this paper, we propose to combine this hybrid architecture with the promising temporal attention mechanism.

**Instance Segmentation of Satellite Images.** The first step of panoptic segmentation is to delineate all individual instances, *i.e.* instance segmentation. Most remote sensing instantiation approaches operate on a single acquisition. For example, several methods have been proposed to detect individual instances of trees [32, 55], buildings [47], or fields [33]. Several algorithms start with a delineation step (border detection) [9, 24, 48], and require postprocessing to obtain individual instances. Other methods use segmentation as a preprocessing step and compute cluster-based features [6, 7], but do not produce explicit cluster-to-object mappings. Petitjean *et al.* [31] propose a segmentation-aided classification method operating on image time series. However, their approach partitions each image separately and does not attempt to retrieve individual objects consistently across the entire sequence. In this paper, we propose the first end-to-end framework for directly performing joint semantic and instance segmentation on SITS.

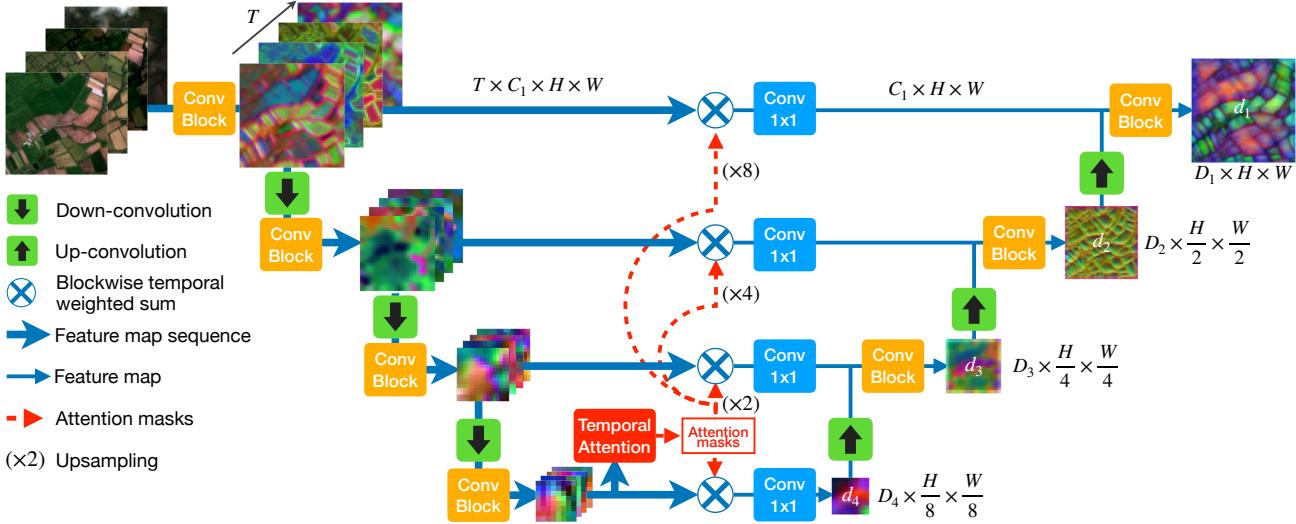
**Panoptic Segmentation of Videos.** Among the vast literature on instance segmentation, Mask-RCNN [11] is the leading method for natural images. Recently, Wang *et al.* proposed CenterMask [49], a lighter and more efficient single-stage method which we use as a starting point in this paper. Several approaches propose extending instance or panoptic segmentation methods from image to video [51, 44, 17]. However, as explained in the introduction, SITS differs from natural video in several key ways which require specific algorithmic and architectural adaptations.

## 3. Method

We consider an image time sequence  $X$ , organized into a four-dimensional tensor of shape  $T \times C \times H \times W$ , with  $T$  the length of the sequence,  $C$  the number of channels, and  $H \times W$  the spatial extent.

### 3.1. Spatio-Temporal Encoding

Our model, dubbed U-TAE (U-Net with Temporal Attention Encoder), encodes a sequence  $X$  in three steps: **(a)** each image in the sequence is embedded simultaneously and independently by a shared multi-level spatial convolutional encoder, **(b)** a temporal attention encoder collapses



**Figure 2: Spatio-temporal Encoding.** A sequence of images is processed in parallel by a shared convolutional encoder. At the lowest resolution, an attention-based temporal encoder produces a set of temporal attention masks for each pixel, which are then spatially interpolated at all resolutions. These masks are used to collapse the temporal dimension of the feature map sequences into a single map per resolution. A convolutional decoder then computes features at all resolution levels. All convolutions operate purely on the spatial and channel dimensions, and we use strided convolutions for both spatial up and down-sampling. The feature maps are projected in RGB space to help visual interpretation.

the temporal dimension of the resulting sequence of feature maps into a single map for each level, (c) a spatial convolutional decoder produces a single feature map with the same resolution as the input images, see Figure 2.

**a) Spatial Encoding.** We consider a convolutional encoder  $\mathcal{E}$  with  $L$  levels  $1, \dots, L$ . Each level is composed of a sequence of convolutions, Rectified Linear Unit (ReLU) activations, and normalizations. Except for the first level, each block starts with a strided convolution, dividing the resolution of the feature maps by a factor 2.

For each time stamp  $t$  simultaneously, the encoder  $\mathcal{E}_l$  at level  $l$  takes as input the feature map of the previous level  $e_t^{l-1}$ , and outputs a feature map  $e_t^l$  of size  $C_l \times H_l \times W_l$  with  $H_l = H/2^{l-1}$  and  $W_l = W/2^{l-1}$ . The resulting feature maps are then temporally stacked into a feature map sequence  $e^l$  of size  $T \times C_l \times H_l \times W_l$ :

$$e^l = [\mathcal{E}_l(e_t^{l-1})]_{t=0}^T \text{ for } l \in [1, L], \quad (1)$$

with  $e^0 = X$  and  $[\cdot]$  the concatenation operator along the temporal dimension. When constituting batches, we flatten the temporal and batch dimensions. Since each sequence comprises images acquired at different times, the batches' samples are not identically distributed. To address this issue, we use Group Normalization [50] with 4 groups instead of Batch Normalization [14] in the encoder.

**b) Temporal Encoding.** In order to obtain a single representation per sequence, we need to collapse the temporal dimension of each feature map sequence  $e^l$  before using them as *skip connections*. Convolutional-recurrent U-Net networks [42, 37, 29] only process the temporal dimension of the lowest resolution feature map with a temporal encoder. The rest of the skip connections are collapsed with a simple temporal average. This prevents the extraction of spatially adaptive and parcel-specific temporal patterns at higher resolutions. Conversely, processing the highest resolution would result in small spatial receptive fields for the temporal encoder, and an increased memory requirement. Instead, we propose an attention-based scheme which only processes the temporal dimension at the lowest feature map resolution, but is able to utilize the predicted temporal attention masks at all resolutions simultaneously.

Based on its performance and computational efficiency, we choose the Lightweight-Temporal Attention Encoder (L-TAE) [10] to handle the temporal dimension. The L-TAE is a simplified multi-head self-attention network [45] in which the attention masks are directly applied to the input sequence of vectors instead of predicted *values*. Additionally, the L-TAE implements a channel grouping strategy similar to Group Normalization [50].

We apply a shared L-TAE with  $G$  heads independently at each pixel of  $e^L$ , the feature map sequence at the lowest level resolution  $L$ . This generates  $G$  temporal attention masks for each pixel, which can be arranged into  $G$  tensors

$a^{L,g}$  with values in  $[0, 1]$  and of shape  $T \times H_L \times W_L$ :

$$a^{L,1}, \dots, a^{L,G} = \text{LTAE}(e^L), \text{ applied pixelwise.} \quad (2)$$

In order to use these attention masks at all scale levels  $l$  of the encoder, we compute spatially-interpolated masks  $a^{l,g}$  of shape  $T \times H_l \times W_l$  for all  $l$  in  $[1, L - 1]$  and  $g$  in  $[1, G]$  with bilinear interpolation:

$$a^{l,g} = \text{resize } a^{L,g} \text{ to } H_l \times W_l. \quad (3)$$

The interpolated masks  $a^{l,g}$  at level  $l$  of the encoder are then used as if they were generated by a temporal attention module operating at this resolution. We apply the L-TAE channel-grouping strategy at all resolution levels: the channels of each feature map sequence  $e^l$  are split into  $G$  contiguous groups  $e^{l,1}, \dots, e^{l,G}$  of identical shape  $T \times C_l/G \times W_l \times H_l$ . For each group  $g$ , the feature map sequence  $e^{l,g}$  is averaged on the temporal dimension using  $a^{l,g}$  as weights. The resulting maps are concatenated along the channel dimension, and processed by a shared  $1 \times 1$  convolution layer  $\text{Conv}_{1 \times 1}^l$  of width  $C_l$ . We denote by  $f^l$  the resulting map of size  $C_l \times W_l \times H_l$  by :

$$f^l = \text{Conv}_{1 \times 1}^l \left( \left[ \sum_{t=1}^T a_t^{l,g} \odot e_t^{l,g} \right]_{g=1}^G \right), \quad (4)$$

with  $[\cdot]$  the concatenation along the channel dimension and  $\odot$  the term-wise multiplication with channel broadcasting.

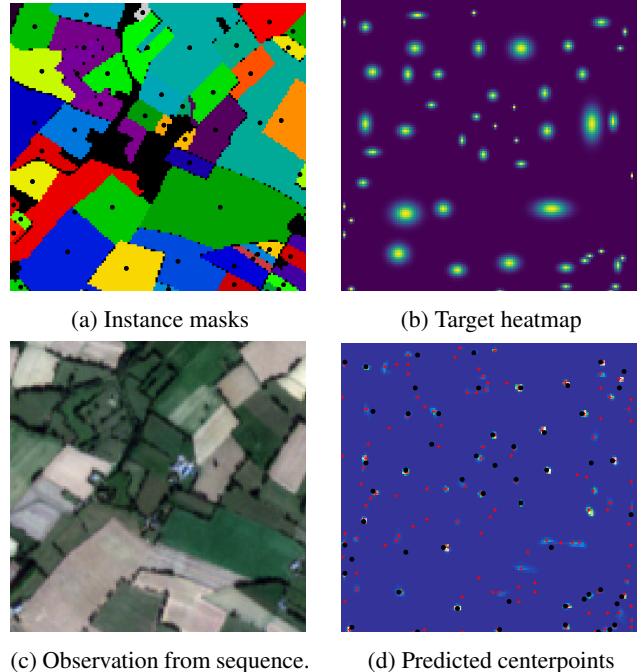
**c) Spatial Decoding.** We combine the feature maps  $f^l$  learned at the previous step with a convolutional decoder to obtain spatio-temporal features at all resolutions. The decoder is composed of  $L - 1$  blocks  $\mathcal{D}_l$  for  $1 \leq l < L$ , with convolutions, ReLu activations, and BatchNorms [14]. Each decoder block uses a strided transposed convolution  $\mathcal{D}_l^{\text{up}}$  to up-sample the previous feature map. The decoder at level  $l$  produces a feature map  $d^l$  of size  $D_l \times H_l \times W_l$ . In a U-Net fashion, the encoder’s map at level  $l$  is concatenated with the output of the decoder block at level  $l - 1$ :

$$d^l = \mathcal{D}_l([\mathcal{D}_l^{\text{up}}(d^{l+1}), f^l]) \text{ for } l \in [1, L - 1], \quad (5)$$

with  $d^L = f^L$  and  $[\cdot]$  is the channelwise concatenation.

### 3.2. Panoptic Segmentation

Our goal is to use the multi-scale feature maps  $\{d^l\}_{l=1}^L$  learned by the spatio-temporal encoder to perform panoptic segmentation of a sequence of satellite images over an area of interest. The first stage of panoptic segmentation is to produce instance proposals, which are then combined into a single panoptic instance map. Since an entire sequence of images (often over 50) must be encoded to compute  $\{d^l\}_{l=1}^L$ , we favor an efficient design for our panoptic segmentation module. Furthermore, given the relative

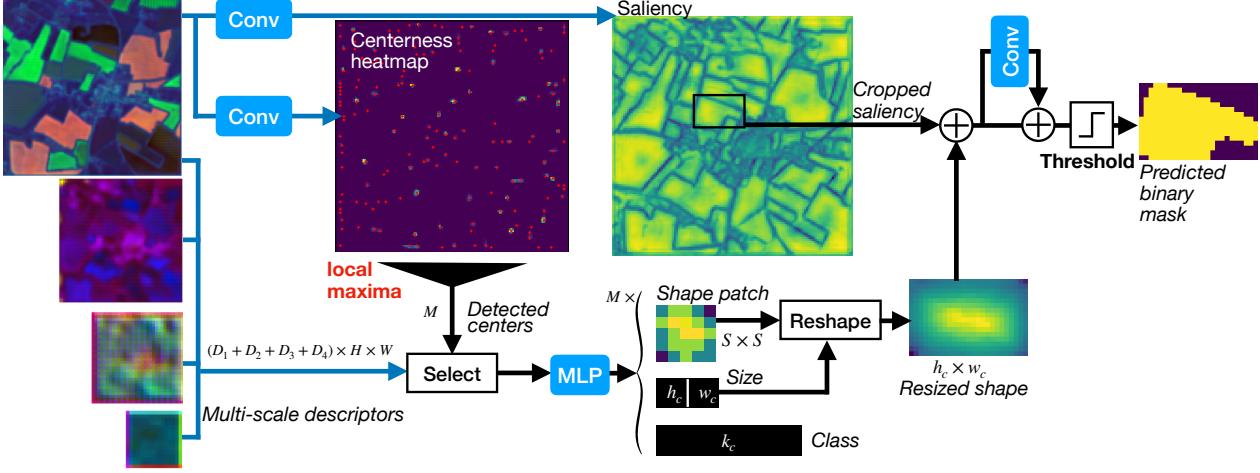


**Figure 3: Centerpoint Detection.** The ground truth instance masks (a) is used to construct a target heatmap (b). Our parcel detection module maps the raw sequence of observation (c) to a predicted heatmap (d). The predicted centerpoints (red crosses) are the local maxima of the predicted heatmap (d). The black dots are the true parcels centers.

simplicity of parcels’ borders, we avoid complex region proposal networks such as Mask-RCNN. Instead, we adapt the single-stage CenterMask instance segmentation network [49], and detail our modifications in the following paragraphs. We name our approach *Parcels-as-Points* (PaPs) to highlight our inspiration from CenterNet/Mask [56, 49].

We denote by  $P$  the set of ground truth parcels in the image sequence  $X$ . Note that the position of these parcels is time-invariant and hence only defined by their spatial extent. Each parcel  $p$  is associated with (i) a centerpoint  $\hat{i}_p, \hat{j}_p$  with integer coordinates, (ii) a bounding box of size  $\hat{h}_p, \hat{w}_p$ , (iii) a binary instance mask  $\hat{s}_p \in \{0, 1\}^{H \times W}$ , (iv) a class  $\hat{k}_p \in [1, K]$  with  $K$  the total number of classes.

**Centerpoint Detection.** Following CenterMask, we perform parcel detection by predicting *centerness heatmaps* supervised by the ground truth parcels’ bounding boxes. In the original approach [56], each class has its own heatmap: detection doubles as classification. This is a sensible choice for natural images, since the tasks of detecting an object’s nature, location, and shape are intrinsically related. In our setting however, the parcels’ shapes and border characteristics are mostly independent of the cultivated crop. For this



**Figure 4: Panoptic Segmentation.** The local maxima of the predicted centerness heatmap defines  $M$  tentative parcels. For each one, the pixel features at all levels are concatenated and used to predict a bounding box size, a semantic class, and an  $S \times S$  shape patch. The latter is combined with a global saliency map for predicting pixel-precise masks. The instance predictions are combined into a panoptic segmentation using the centerness as quality.

reason, we use a single centerness heatmap and postpone class identification to a subsequent specialized module. See Figure 3 for an illustration of our parcel detection method.

We associate each parcel  $p$  with a Gaussian kernel of deviations  $\sigma_p^{\text{ver}}$  and  $\sigma_p^{\text{hor}}$  taken respectively as 1/20 of the height and width of the parcels’ bounding box. Unlike Law and Deng [21], we use heteroschedastic kernels to reflect the potential narrowness of parcels. We then define the target centerness heatmap  $\hat{m} \in [0, 1]^{H \times W}$  as the maximum value of all parcel kernels at each pixel  $(i, j)$  in  $H \times W$ :

$$\hat{m}_{i,j} = \max_{p \in P} \exp \left( - \left[ \frac{(i - \hat{i}_p)^2}{2(\sigma_p^{\text{ver}})^2} + \frac{(j - \hat{j}_p)^2}{2(\sigma_p^{\text{hor}})^2} \right] \right) \quad (6)$$

A convolutional layer takes the highest-resolution feature map  $d^1$  as input and predicts a centerness heatmap  $m \in [0, 1]^{H \times W}$ . The predicted heatmap is supervised using the loss defined in Equation 7 with  $\beta = 4$ :

$$\mathcal{L}_{\text{center}} = \frac{-1}{|P|} \sum_{i=1 \dots H} \sum_{j=1 \dots W} \begin{cases} \log(m_{i,j}) & \text{if } \hat{m}_{i,j} = 1 \\ ((1 - \hat{m}_{i,j})^\beta \log(1 - m_{i,j})) & \text{else.} \end{cases} \quad (7)$$

We define the predicted centerpoints as the local maxima of  $m$ , *i.e.* pixels with larger values than their 8 adjacent neighbors. This set can be efficiently computed with a single max-pooling operation. Replacing the max operator by argmax in Equation 6 defines a mapping  $H \times W \mapsto P$  between pixels and parcels. During training, we associate each true parcel  $p$  with the predicted centerpoint  $\mathbf{c}(p)$  with highest predicted centerness  $m$  among the set of centerpoints which coordinates are mapped to  $p$ . If this set is

empty, then  $\mathbf{c}(p)$  is undefined: the parcel  $p$  is not detected. We denote by  $P'$  the subset of detected parcels, *i.e.* for which  $\mathbf{c}(p)$  is well defined.

**Size and Class Prediction.** We associate to a predicted centerpoint  $c$  of coordinate  $(i_c, j_c)$  the multi-scale feature vector  $\tilde{d}_c$  of size  $D_1 + \dots + D_L$  by concatenating channelwise the pixel features at location  $(i_c, j_c)$  in all maps  $d^l$ :

$$\tilde{d}_c = [d^l ([i_c/2^{l-1}], [j_c/2^{l-1}])]_{l=1}^L, \quad (8)$$

with  $[\cdot]$  the channelwise concatenation. This vector  $\tilde{d}_c$  is then processed by four different multilayer perceptrons (MLP) to obtain three vectors of sizes 2,  $K$ , and  $S^2$  representing respectively: (i) a bounding box size  $h_c, w_c$ , (ii) a vector of class probabilities  $k_c$  of size  $K$ , and (iii) a shape patch  $s_c$  of fixed size  $S \times S$ . The latter is described in the next paragraph.

The class prediction  $k_{\mathbf{c}(p)}$  associated to the true parcel  $p$  is supervised with the cross-entropy loss, and the size prediction with a normalized L1 loss. For all  $p$  in  $P'$ , we have:

$$\mathcal{L}_{\text{class}}^p = -\log(k_{\mathbf{c}(p)}[\hat{k}_p]) \quad (9)$$

$$\mathcal{L}_{\text{size}}^p = \frac{|h_{\mathbf{c}(p)} - \hat{h}_p|}{\hat{h}_p} + \frac{|w_{\mathbf{c}(p)} - \hat{w}_p|}{\hat{w}_p}. \quad (10)$$

**Shape Prediction.** The idea of this step is to combine for a predicted centerpoint  $c$  a rough shape patch  $s_c$  with a full-resolution global saliency map  $z$  to obtain a pixel-precise

instance mask, see Figure 4. For a centerpoint  $c$  of coordinates  $(i_c, j_c)$ , the predicted shape patch  $s_c$  of size  $S \times S$  is resized to the predicted size  $\lceil h_c \rceil \times \lceil w_c \rceil$  with bilinear interpolation. A convolutional layer maps the outermost feature map  $d^1$  to a saliency map  $z$  of size  $H \times W$ , which is shared by all predicted parcels. This saliency map is then cropped along the predicted bounding box  $(i_c, j_c, \lceil h_c \rceil, \lceil w_c \rceil)$ . The resized shape and the cropped saliency are added (11) to obtain a first local shape  $\tilde{l}_c$ , which is then further refined with a residual convolutional network CNN (12). We denote the resulting predicted shape by  $l_c$ :

$$\tilde{l}_c = \text{resize}_c(s_c) + \text{crop}_c(z) \quad (11)$$

$$l_c = \text{sigmoid}(\tilde{l}_c + \text{CNN}(\tilde{l}_c)), \quad (12)$$

with  $\text{resize}_c$  and  $\text{crop}_c$  defined by the coordinates  $(i_c, j_c)$  and predicted bounding box size  $(\lceil h_c \rceil, \lceil w_c \rceil)$ . The shape and saliency predictions are supervised for each parcel  $p$  in  $P'$  by computing the pixelwise binary cross-entropy (BCE) between the predicted shape  $l_{c(p)}$  and the corresponding true binary instance mask  $\hat{s}_p$  cropped along the predicted bounding box  $(i_{c(p)}, j_{c(p)}, \lceil h_{c(p)} \rceil, \lceil w_{c(p)} \rceil)$ :

$$\mathcal{L}_{\text{shape}}^p = \text{BCE}(l_{c(p)}, \text{crop}_{c(p)}(\hat{s}_p)). \quad (13)$$

For inference, we associate a binary mask with a predicted centerpoint  $c$  by thresholding  $l_c$  with the value 0.4.

**Loss Function** : These four losses are combined into a single loss with no weight and optimized end-to-end:

$$\mathcal{L} = \mathcal{L}_{\text{center}} + \frac{1}{|P'|} \sum_{p \in P'} \left( \mathcal{L}_{\text{class}}^p + \mathcal{L}_{\text{size}}^p + \mathcal{L}_{\text{shape}}^p \right). \quad (14)$$

**Differences with CenterMask.** Our approach differs from CenterMask in several key ways: (i) We compute a single saliency map and heatmap instead of  $K$  different ones. This represents the absence of parcel occlusion and the similarity of their shapes. (ii) Accounting for the lower resolution of satellite images, centerpoints are computed at full resolution to detect potentially small parcels, thus dispensing us from predicting offsets. (iii) The class prediction is handled centerpoint-wise instead of pixel-wise for efficiency. (iv) Only the selected centerpoints predict shape, class, and size vectors, saving computation and memory. (v) We use simple feature concatenation to compute multi-scale descriptors instead of deep layer aggregation [53] or stacked Hourglass-Networks [27]. (vi) A convolutional network learns to combine the saliency and the mask instead of a simple term-wise product.

**Converting to Panoptic Segmentation.** Panoptic segmentation consists of associating to each pixel a semantic

label and, for non-background pixels (our only *stuff* class), an instance label [19]. Our predicted binary instance masks can have overlaps, which we resolve by associating to each predicted parcel a quality measure equal to the predicted centerness  $m$  at its associated centerpoint. Masks with higher quality overtake the pixels of overlapping masks with lesser predicted quality. If a mask loses more than 50% of its pixels through this process, it is removed altogether from the predicted instances. Predicted parcels with a quality under a given threshold are dropped. This threshold can be tuned on a validation set to maximize the parcel detection F-score. All pixels not associated with a parcel mask are labelled as background.

**Implementation Details.** Our implementation of U-TAE allows for batch training on sequences of variable length thanks to a simple padding strategy. The complete configuration and training details can be found in the Appendix. A Pytorch implementation is available at <https://github.com/VSainteuf/utae-paps>.

## 4. Experiments

### 4.1. The PASTIS Dataset

We present PASTIS (Panoptic Agricultural Satellite TIme Series), the first large-scale, publicly available SITS dataset with both semantic and panoptic annotations. This dataset, as well as more information about its composition, are publicly available at <https://github.com/VSainteuf/pastis-benchmark>.

**Description.** PASTIS is comprised of 2 433 sequences of multi-spectral images of shape  $10 \times 128 \times 128$ . Each sequence contains between 38 and 61 observations taken between September 2018 and November 2019, for a total of over 2 billion pixels. The time between acquisitions is uneven with a median of 5 days. This lack of regularity is due to the automatic filtering of acquisitions with extensive cloud cover by the satellite data provider THEIA. The 10 channels correspond to the non-atmospheric spectral bands of the Sentinel-2 satellite, after atmospheric correction and re-sampling at a spatial resolution of 10 meters per pixel. The dataset spans over  $4000 \text{ km}^2$ , with images taken from four different regions of France with diverse climates and crop distributions, covering almost 1% of the French Metropolitan territory. We estimate that close to 28% of images have at least partial cloud cover.

**Annotation.** Each pixel of PASTIS is associated with a semantic label taken from a nomenclature of 18 crop types plus a background class. As is common in remote sensing applications, the dataset is highly unbalanced, with a ratio of over 50 between the most and least common classes.



(a) Image from the sequence. (b) Panoptic annotation. (c) Panoptic segmentation. (d) Semantic segmentation.

**Figure 5: Qualitative results.** We consider an image sequence (a) with panoptic annotations (b). We represent the results of our method in terms of panoptic segmentation (c) and semantic segmentation (d). The parcels’ and pixels’ color corresponds to the crop type, according to a legend given in the appendix. The predominantly correct class predictions highlight the fact that the difficulty of panoptic segmentation lies in the precise delineation of each individual parcel. We observe cases where the temporal structure of the SITS was successfully leveraged to resolve boundary ambiguities that could not be seen from a single image (cyan circle  $\circlearrowleft$ ). Conversely, some visually fragmented parcels are annotated as a single instance (red circle  $\circlearrowright$ ).

Each non-background pixel also has a unique instance label corresponding to its parcel index. In total, 124 422 parcels are individualized, each with their bounding box, pixel-precise mask, and crop type. All annotations are taken from the publicly available French Land Parcel Identification System. The French Payment Agency estimates the accuracy of the crop annotations via in situ control over 98% and the relative error in terms of surfaces under 0.3%. To allow for cross-validation, the dataset is split into 5 folds, chosen with a 1km buffer between images to avoid cross-fold contamination.

## 4.2. Semantic Segmentation

Our U-TAE has  $L = 4$  resolution levels and a LTAE with  $G = 16$  heads, see appendix for an exact configuration. For the semantic segmentation task, the feature map  $d_1$  with highest resolution is set to have  $K$  channels, with  $K$  the number of classes. We can then interpret  $d_1$  as pixel-wise predictions to be supervised with the cross-entropy loss. In this setting, we do not use the PaPs module.

**Competing Methods.** We reimplemented six of the top-performing SITS encoders proposed in the literature:

- *ConvLSTM* [35, 40] and *ConvGRU*[4]. These approaches are recurrent neural networks in which all linear layers are replaced by spatial convolutions.
- *U-ConvLSTM* [37] and *U-BiConvLSTM* [23]. To reproduce these UNet-Based architectures, we replaced the L-TAE in our architecture by either a convLSTM [41] or a bidirectional convLSTM. Skip connections are temporally averaged. In contrast to the original methods, we replaced the batch normalization in the encoders with

**Table 1: Semantic Segmentation.** We report for our method and six competing methods the model size in trainable parameters, Overall Accuracy (OA), mean Intersection over Union (mIoU), and Inference Time for one fold of  $\sim 490$  sequences (IT). The second part of the table report results from our ablation study.

Model	# param $\times 1000$	OA	mIoU	IT (s)
U-TAE (ours)	1 087	<b>83.2</b>	<b>63.1</b>	<b>25.7</b>
3D-Unet [37]	1 554	81.3	58.4	29.5
U-ConvLSTM [37]	1 508	82.1	57.8	28.3
FPN-ConvLSTM [23]	1 261	81.6	57.1	103.6
U-BiConvLSTM [23]	1 434	81.8	55.9	32.7
ConvGRU [4]	1 040	79.8	54.2	49.0
ConvLSTM [35, 40]	1 010	77.9	49.1	49.1
Mean Attention	1 087	82.8	60.1	24.8
Skip Mean + Conv	1 087	82.4	58.9	24.5
Skip Mean	1 074	82.0	58.3	24.5
BatchNorm	1 087	71.9	36.0	22.3
Single Date (August)	1 004	65.6	28.3	1.3
Single Date (May)	1 004	58.1	20.6	1.3

group normalization which significantly improved the results across-the-board.

- *3D-Unet* [37]. A U-Net in which the convolutions of the encoding branch are three-dimensional to handle simultaneously the spatial and temporal dimensions.
- *FPN-ConvLSTM* [23]. This model combines a feature pyramid network [22] to extract spatial features and a bidirectional ConvLSTM for the temporal dimension.

**Analysis.** In Table 1, we detail the performance obtained with 5-fold cross validation of our approach and the six reimplemented baselines. We report the Overall Accuracy (OA) as the ratio between correct and total predictions, and (mIoU) the class-averaged classification IoU. We observe that the convolutional-recurrent methods *ConvGRU* and *ConvLSTM* perform worse. Recurrent networks embedded in an U-Net or a FPN share similar performance, with a much longer inference time for FPN. Our approach significantly outperforms all other methods in terms of precision. In Figure 5, we present a qualitative illustration of the semantic segmentation results.

**Ablation Study.** We first study the impact of using spatially interpolated attention masks to collapse the temporal dimension of the spatio-temporal feature maps at different levels of the encoder simultaneously. Simply computing the temporal average of skip connections for levels without temporal encoding as proposed by [42, 37], we observe a drop of 4.8 mIoU points (Skip Mean). This puts our method performance on par with its competing approaches. Adding a  $1 \times 1$  convolutional layer after the temporal average reduces this drop to 4.2 points (Skip Mean + Conv). Lastly, using interpolated masks but foregoing the channel grouping strategy by averaging the masks group-wise into a single attention mask per level results in a drop of 3.1 points (Mean Attention). This implies that our network is able to use the grouping scheme at different resolutions simultaneously. In conclusion, the main advantage of our proposed attention scheme is that the temporal collapse is controlled at all resolutions, in contrast to recurrent methods.

Using batch normalization in the encoder leads to a severe degradation of the performance of 27.1 points (Batch-Norm). We conclude that the temporal diversity of the acquisitions requires special considerations. This was observed for all U-Net models alike. We also train our model on a single acquisition date (with a classic U-Net and no temporal encoding) for two different cloudless dates in August and May (Single Date). We observe a drop of 24.8 and 42.5 points respectively, highlighting the crucial importance of the temporal dimension for crop classification. We also observed that images with at least partial cloud cover received on average 58% less attention than their cloud-free counterparts. This suggests that our model is able to use the attention module to automatically filter out corrupted data.

### 4.3. Panoptic Segmentation

We use the same U-TAE configuration for panoptic segmentation, and select a PaPs module with 190k parameters and a shape patch size of  $16 \times 16$ . In Table 2, we report the class-averaged Segmentation Quality (SQ), Recognition Quality (RQ), and Panoptic Quality (PQ) [19]. We observe that while the network is able to correctly detect and clas-

**Table 2: Panoptic Segmentation Experiment.** We report class-averaged panoptic metrics: SQ, RQ, PQ (see Metric Correction paragraph before references).

	SQ	RQ	PQ
<b>U-TAE + PaPs</b>	81.5	53.2	43.8
U-ConvLSTM + Paps	80.2	43.9	35.6
$S = 24$	80.7	50.6	41.3
$S = 8$	80.9	52.3	42.7
Multiplicative Saliency	74.6	49.9	37.5
Single-image	72.3	18.7	14.1

sify most parcels, the task remains difficult. In particular, the combination of ambiguous borders and hard-to-classify parcel content makes for a challenging panoptic segmentation problem. We illustrate these difficulties in Figure 5, along with qualitative results.

Replacing the temporal encoder by a U-BiConvLSTM as described in Section 4.2 (U-BiConvLSTM+PaPs), we observe a noticeable performance drop of 8.2 PQ, which is consistent with the results of Table 1. As expected, our model’s performance is not too sensitive to changes in the size  $S$  of the shape patch. Indeed, the shape patches only determine the rough outline of parcels while the pixel-precise instance masks are derived from the saliency map. Performing shape prediction with a simple element-wise multiplication as in [49] (Multiplicative Saliency) instead of our residual CNN results in a drop of over  $-6.9$  SQ. Using a single image (August) leads to a low panoptic quality. Indeed, identifying crop types and parcel borders from a single image at the resolution of Sentinel-2 is particularly difficult.

Inference on 490 sequences takes 129s: 26s to generate U-TAE embeddings, 1s for the heatmap and saliency, 90s for instance proposals, and 12s to merge them into a panoptic segmentation. Note that the training time is also doubled compared to simple semantic segmentation.

## 5. Conclusion

We introduced U-TAE, a novel spatio-temporal encoder using a combination of spatial convolution and temporal attention. This model can be easily combined with *PaPs*, the first panoptic segmentation framework operating on SITS. Lastly, we presented PASTIS, the first large-scale panoptic-ready SITS dataset. Evaluated on this dataset, our approach significantly outperformed all other approaches for semantic segmentation, and set up the first state-of-the-art for panoptic segmentation of satellite image sequences.

We hope that the combination of our open-access dataset and promising results will encourage both remote sensing and computer vision communities to consider the challenging problem of panoptic SITS segmentation, whose economic and environmental stakes can not be understated.

## Metric Correction

The values reported in this version of the article for the Panoptic Segmentation experiment differ from the version published in the ICCV 2021 proceedings. Indeed, a bug in the computation of the Recognition Quality (RQ) metric was present in the original implementation resulting in the *void* target instances not being properly ignored. Instead, all predictions matched to *void* target instances were counted as false positives, thus artificially reducing the RQ score. Since the panoptic metrics are not involved in the training loss, this bug did not impact the overall training procedure. All models of Table 2 were re-evaluated with the corrected implementation. Across methods this resulted in a  $\sim 3$  PQ increase, driven by a similar increase in RQ. Refer to [github.com/VSainteuf/utae-paps/issues/11](https://github.com/VSainteuf/utae-paps/issues/11) for more details.

## References

- [1] Nicolas Audebert, Bertrand Le Saux, and Sébastien Lefèvre. Semantic segmentation of earth observation data using multimodal and multi-scale deep networks. In *ACCV*, 2016.
- [2] Adeline Bailly, Simon Malinowski, Romain Tavenard, Laetitia Chapel, and Thomas Guyet. Dense bag-of-temporal-siftwords for time series classification. In *International Workshop on Advanced Analysis and Learning on Temporal Data*. Springer, 2015.
- [3] Simon Bailly, Sébastien Giordano, Loic Landrieu, and Nesrine Chehata. Crop-rotation structured classification using multi-source Sentinel images and LPIS for crop type mapping. In *IGARSS*, 2018.
- [4] Nicolas Ballas, Li Yao, Chris Pal, and Aaron Courville. Delving deeper into convolutional networks for learning video representations. *ICLR*, 2016.
- [5] Christopher Boshuizen, James Mason, Pete Klupar, and Shannon Spanhake. Results from the planet labs flock constellation. *AIAA/USU Conference on Small Satellites*, 2014.
- [6] Alessandro Michele Censi, Dino Ienco, Yawogan Jean Eudes Gbodjo, Ruggero Gaetano Pensa, Roberto Interdonato, and Raffaele Gaetano. Spatial-temporal GraphCNN for land cover mapping. *IEEE Access*, 2021.
- [7] Dawa Derksen, Jordi Ingla, and Julien Michel. Spatially precise contextual features based on superpixel neighborhoods for land cover mapping with high resolution satellite image time series. In *IGARSS*, 2018.
- [8] Matthias Drusch, Umberto Del Bello, Sébastien Carlier, Olivier Colin, Veronica Fernandez, Ferran Gascon, Bianca Hoersch, Claudia Isola, Paolo Laberinti, Philippe Martimort, et al. Sentinel-2: Esa's optical high-resolution mission for gmes operational services. *Remote sensing of Environment*, 2012.
- [9] Angel Garcia-Pedrero, Consuelo Gonzalo-Martin, and M Lillo-Saavedra. A machine learning approach for agricultural parcel delineation through agglomerative segmentation. *International journal of remote sensing*, 2017.
- [10] Vivien Sainte Fare Garnot and Loic Landrieu. Lightweight temporal self-attention for classifying satellite images time series. In *International Workshop on Advanced Analytics and Learning on Temporal Data*. Springer, 2020.
- [11] Kaiming He, Georgia Gkioxari, Piotr Dollár, and Ross Girshick. Mask R-CNN. In *ICCV*, 2017.
- [12] Dino Ienco, Raffaele Gaetano, Claire Dupaquier, and Pierre Maurel. Land cover classification via multitemporal spatial data by deep recurrent neural networks. *Geoscience and Remote Sensing Letters*, 2017.
- [13] Jordi Ingla, Marcela Arias, Benjamin Tardy, Olivier Hagolle, Silvia Valero, David Morin, Gérard Dedieu, Guadalupe Sepulcre, Sophie Bontemps, Pierre Defourny, et al. Assessment of an operational system for crop type map production using high temporal and spatial resolution satellite optical imagery. *Remote Sensing*, 2015.
- [14] Sergey Ioffe and Christian Szegedy. Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *ICML*, 2015.
- [15] Shunping Ji, Chi Zhang, Anjian Xu, Yun Shi, and Yulin Duan. 3d convolutional neural networks for crop classification with multi-temporal remote sensing images. *Remote Sensing*, 2018.
- [16] Andreas Kamlaris and Francesc X Prenafeta-Boldú. Deep learning in agriculture: A survey. *Computers and electronics in agriculture*, 2018.
- [17] Dahun Kim, Sanghyun Woo, Joon-Young Lee, and In So Kweon. Video panoptic segmentation. In *CVPR*, 2020.
- [18] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *ICLR*, 2015.
- [19] Alexander Kirillov, Ross Girshick, Kaiming He, and Piotr Dollár. Panoptic feature pyramid networks. In *CVPR*, 2019.
- [20] Nataliia Kussul, Mykola Lavreniuk, Sergii Skakun, and Andrii Shelestov. Deep learning classification of land cover and crop types using remote sensing data. *Geoscience and Remote Sensing Letters*, 2017.
- [21] Hei Law and Jia Deng. Cornernet: Detecting objects as paired keypoints. In *ECCV*, 2018.
- [22] Tsung-Yi Lin, Piotr Dollár, Ross Girshick, Kaiming He, Bharath Hariharan, and Serge Belongie. Feature pyramid networks for object detection. In *CVPR*, 2017.
- [23] 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 multiday crop recognition from multitemporal image sequences. *ISPRS Journal of Photogrammetry and Remote Sensing*, 2021.
- [24] Khairiya Mudrik Masoud, Claudio Persello, and Valentyn A Tolpekin. Delineation of agricultural field boundaries from Sentinel-2 images using a novel super-resolution contour detector based on fully convolutional networks. *Remote sensing*, 2020.
- [25] Nando Metzger, Mehmet Ozgur Turkoglu, Stefano D'Aronco, Jan Dirk Wegner, and Konrad Schindler. Crop classification under varying cloud cover with neural ordinary differential equations. *arXiv preprint arXiv:2012.02542*, 2020.

- [26] Rohit Mohan and Abhinav Valada. Efficienttps: Efficient panoptic segmentation. *International Journal of Computer Vision*, 2021.
- [27] Alejandro Newell, Kaiyu Yang, and Jia Deng. Stacked hour-glass networks for human pose estimation. In *ECCV*, 2016.
- [28] Mehmet Ozgur Turkoglu, Stefano D’Aronco, Gregor Perich, Frank Liebisch, Constantin Streit, Konrad Schindler, and Jan Dirk Wegner. Crop mapping from image time series: deep learning with multi-scale label hierarchies. *arXiv preprint arXiv:2102.08820*, 2021.
- [29] Maria Papadomanolaki, Maria Vakalopoulou, and Konstantinos Karantzalos. A deep multi-task learning framework coupling semantic segmentation and fully convolutional LSTM networks for urban change detection. *Transactions on Geoscience and Remote Sensing*, 2021.
- [30] Charlotte Pelletier, Geoffrey I Webb, and François Petitjean. Temporal convolutional neural network for the classification of satellite image time series. *Remote Sensing*, 2019.
- [31] François Petitjean, Camille Kurtz, Nicolas Passat, and Pierre Gançarski. Spatio-temporal reasoning for the classification of satellite image time series. *Pattern Recognition Letters*, 2012.
- [32] Yuchu Qin, Antonio Ferraz, Clément Mallet, and Corina Ivan. Individual tree segmentation over large areas using airborne lidar point cloud and very high resolution optical imagery. In *IGARSS*, 2014.
- [33] Christoph Rieke. Deep learning for instance segmentation of agricultural fields. [https://github.com/chrieke/InstanceSegmentation\\_Sentinel2](https://github.com/chrieke/InstanceSegmentation_Sentinel2), 2017.
- [34] Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomedical image segmentation. In *MICCAI*, 2015.
- [35] Marc Rußwurm and Marco Körner. Convolutional LSTMs for cloud-robust segmentation of remote sensing imagery. *NeurIPS Workshops*, 2018.
- [36] Marc Rußwurm and Marco Körner. Self-attention for raw optical satellite time series classification. *ISPRS*, 2020.
- [37] Rose 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 *CVPR Workshops*, 2019.
- [38] Vivien Sainte Fare Garnot, Loic Landrieu, Sébastien Giordano, and Nesrine Chehata. Time-space tradeoff in deep learning models for crop classification on satellite multispectral image time series. In *IGARSS*, 2019.
- [39] Vivien Sainte Fare Garnot, Loic Landrieu, Sébastien Giordano, and Nesrine Chehata. Satellite image time series classification with pixel-set encoders and temporal self-attention. In *CVPR*, 2020.
- [40] 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. In *NeurIPS*, 2015.
- [41] 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. *arXiv preprint arXiv:1506.04214*, 2015.
- [42] 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*, 2019.
- [43] Romain Tavenard, Simon Malinowski, Laetitia Chapel, Adeline Bailly, Heider Sanchez, and Benjamin Bustos. Efficient temporal kernels between feature sets for time series classification. In *ECML-KDD*. Springer, 2017.
- [44] Pavel Tokmakov, Cordelia Schmid, and Karteek Alahari. Learning to segment moving objects. *International Journal of Computer Vision*, 2019.
- [45] Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N Gomez, Łukasz Kaiser, and Illia Polosukhin. Attention is all you need. *NeurIPS*, 2017.
- [46] Francesco Vuolo, Martin Neuwirth, Markus Immitzer, Clement Atzberger, and Wai-Tim Ng. How much does multi-temporal Sentinel-2 data improve crop type classification? *International journal of applied earth observation and geoinformation*, 2018.
- [47] Fabien H Wagner, Ricardo Dalagnol, Yuliya Tarabalka, Tasiana YF Segantini, Rogério Thomé, and Mayumi Hirye. Unet-id, an instance segmentation model for building extraction from satellite images—case study in the joanópolis city, brazil. *Remote Sensing*, 2020.
- [48] François Waldner and Foivos I Diakogiannis. Deep learning on edge: extracting field boundaries from satellite images with a convolutional neural network. *Remote Sensing of Environment*, 2020.
- [49] Yuqing Wang, Zhaoliang Xu, Hao Shen, Baoshan Cheng, and Lirong Yang. Centermask: single shot instance segmentation with point representation. In *CVPR*, 2020.
- [50] Yuxin Wu and Kaiming He. Group normalization. In *ECCV*, 2018.
- [51] Linjie Yang, Yuchen Fan, and Ning Xu. Video instance segmentation. In *CVPR*, 2019.
- [52] Lexiang Ye and Eamonn Keogh. Time series shapelets: a new primitive for data mining. In *ACM SIGKDD*, 2009.
- [53] Fisher Yu, Dequan Wang, Evan Shelhamer, and Trevor Darrell. Deep layer aggregation. In *CVPR*, 2018.
- [54] Yuan Yuan and Lei Lin. Self-supervised pre-training of transformers for satellite image time series classification. *Journal of Selected Topics in Applied Earth Observations and Remote Sensing*, 2020.
- [55] Tiebiao Zhao, Haoyu Niu, Erick de la Rosa, David Doll, Dong Wang, and YangQuan Chen. Tree canopy differentiation using instance-aware semantic segmentation. In *ASABE Annual International Meeting*, 2018.
- [56] Xingyi Zhou, Dequan Wang, and Philipp Krähenbühl. Objects as points. *arXiv preprint arXiv:1904.07850*, 2019.

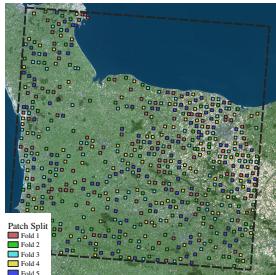
## Supplementary Material

In this appendix, we provide additional information on the PASTIS dataset and our exact model configuration. We also provide complementary qualitative experimental results.

### A.1. PASTIS Dataset



(a) Location of the four tiles.



(b) Selected patches.



(c) Single patch.

**Figure 6: Data Location.** Spatial distribution of the four Sentinel tiles used in PASTIS 6a, and of the selected patches of tile T30UXV 6b. We show an example of patch in 6c, and highlight with red circles examples of parcels that are mostly outside of the patch’s extent and thus annotated with the void label. The green circle  $\circlearrowright$  highlight a parcel partially cut off by the patch borders, but with sufficient overlap to be kept as a valid parcel.

**Overview.** The PASTIS dataset is composed of 2433 square  $128 \times 128$  patches with 10 spectral bands and at 10m resolution, obtained from the open-access Sentinel-2 platform.<sup>1</sup> For each patch, we stack all available acquisitions

<sup>1</sup><https://scihub.copernicus.eu>

between September 2018 and November 2019, forming our four dimensional multi-spectral SITS:  $T \times C \times H \times W$ . The publicly available French Land Parcel Identification System (FLPIS) allows us to retrieve the extent and content of all parcels within the tiles, as reported by the farmers. Each patch pixel is annotated with a semantic label corresponding to either the parcels’ crop type or the background class. The pixels of each unique parcel in the patch receive a corresponding instance label.

**Dataset Extent.** The SITS of PASTIS are taken from 4 different Sentinel-2 tiles in different regions of the French metropolitan territory as depicted in Figure 6a. These regions cover a wide variety of climates and culture distributions. Sentinel tiles span  $100 \times 100$ km and have a spatial resolution of 10 meter per pixel. Each pixel is characterized by 13 spectral bands. We select all bands except the atmospheric bands B01, B09, and B10. Each of these tiles is subdivided in square patches of size  $1.28 \times 1.28$ km ( $128 \times 128$  pixels at 10m/pixel), for a total of around 24,000 patches. We then select 2,433 patches ( 10% of all available patches, see Figure 6b), favoring patches with rare crop types in order to decrease the otherwise extreme class imbalance of the dataset.

**Nomenclature** The FLPIS uses a 73 class breakdown for crop types. We select classes with at least 400 parcels and with samples in at least 2 of the 4 Sentinel-2 tiles. This leads us to adopt a 18 classes nomenclature, presented in Figure 7. Parcels belonging to classes not in our 18-classes nomenclature are annotated with the *void* label, see below.

**Patch Boundaries.** The FLPIS allows us to retrieve the pixel-precise borders of each parcel. We also compute bounding boxes for each parcel. The parcels’ extents are cropped along the extent of their  $128 \times 128$  patch, and the bounding boxes are modified accordingly. Parcels whose surface is more than 50% outside of the patch are annotated with the *void* label, see Figure 6c.

**Void and Background Labels.** Pixels which are not within the extent of any declared parcel are annotated with the background “stuff” label, corresponding to all non-agricultural land uses. For the semantic segmentation task, this label becomes the 20-th class to predict. In the panoptic setting, this label is associated with pixels not within the extent of any predicted parcel. We do not compute the panoptic metrics for the background class, since our focus is on retrieving the parcels’ extent rather than an extensive land-cover prediction. In other words, the reported panoptic metrics are the “things” metrics, which already penalize parcels predicted for background pixels by counting them as false positives.

Label and Color	Class Name	Number of parcels
0	Background	-
1	Meadow	31292
2	Soft winter wheat	8206
3	Corn	13123
4	Winter barley	2766
5	Winter rapeseed	1769
6	Spring barley	908
7	Sunflower	1355
8	Grapevine	10640
9	Beet	871
10	Winter triticale	1208
11	Winter durum wheat	1704
12	Fruits, vegetables, flowers	2619
13	Potatoes	551
14	Leguminous fodder	3174
15	Soybeans	1212
16	Orchard	2998
17	Mixed cereal	848
18	Sorghum	707
19	Void label	35924

Figure 7: Color code of our class nomenclature, and the number of parcel per class.

The void class is reserved for *out-of-scope* parcels, either because their crop type is not in our nomenclature or because their overlap with the selected square patch is too small. We remove these parcels from all semantic or panoptic metrics and losses. Predicted parcels which overlap with an IoU superior to 0.5 with a void parcel are not counted as false positive or true positive, but are simply ignored by the metric, as recommended in [19].

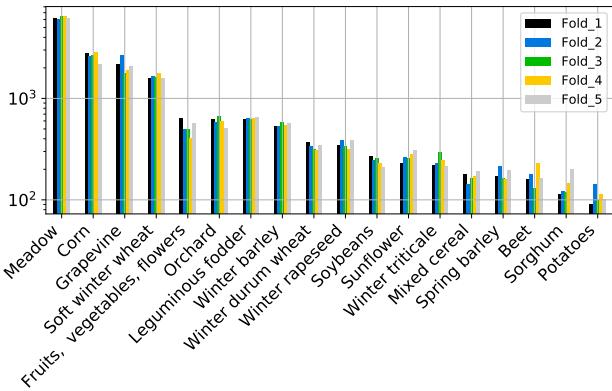


Figure 8: Class distribution for the five folds (in log-scale).

**Cross-Validation.** The 2,433 selected patches are randomly subdivided into 5 splits, allowing us to perform cross-validation. The official 5-fold cross-validation scheme used for benchmarking is given in Table 3. In order to avoid heterogeneous folds, each fold is constituted of

patches taken from all four Sentinel tiles. We also chose folds with comparable class distributions, as measured by their pairwise Kullback-Leiber divergence. We show the resulting class distribution for each fold in Figure 8. Finally, we prevent adjacent patches from being in different folds to avoid data contamination. Geo-referencing metadata of the patches and parcels is included in PASTIS, allowing for the constitution of geographically consistent folds to evaluate spatial generalization. However, this is out of the scope of this paper.

Fold	Train	Val	Test
I	1-2-3	4	5
II	2-3-4	5	1
III	3-4-5	1	2
IV	4-5-1	2	3
V	5-1-2	3	4

Table 3: Official 5-fold cross validation scheme. Each line gives the repartition of the splits into train, validation and test set for each fold.

**Temporal Sampling.** The temporal sampling of the sequences in PASTIS is irregular: depending on their location, patches are observed a different number of times and at different intervals. This is a result of both the orbit schedule of Sentinel-2 and the policy of Sentinel data providers not to process tile observations identified as covered by clouds for more than 90% of the tile’s surface. As this corresponds to the *real world* setting, we decided to leave the SITS as is, and thus to encourage methods that can favourably address this technical challenge. As a result, the proposed SITS are constituted of 33 to 61 acquisitions. In order to assess how our model handles lower sampling frequencies, we limited the number of available acquisitions at inference time<sup>2</sup>, and observed a drop of performance of -0.7, -2.0, -5.5, and -14.6 points of mIoU with 32, 24, 16, and 8 available dates, respectively.

**Clouds Cover.** Even after the automatic filtering of predominantly cloudy acquisitions, some patches are still partially or completely obstructed by cloud cover. We opt to not apply further pre-processing or cloud detection, and produce the raw data in PASTIS. Our reasoning is that an adequate algorithm should be able to learn to deal with such acquisitions. Indeed, robustness to cloud-cover has been experimentally demonstrated for deep learning methods by Rußwurm and Körner [35, 36].

<sup>2</sup>This can be interpreted as the test set having an increased cloud cover.

Table 4: Width of the feature maps outputted at each level of the encoding and decoding branches of the spatial module.

Encoder		Decoder	
$e_1$	64	$d_1$	32
$e_2$	64	$d_2$	32
$e_3$	64	$d_3$	64
$e_4$	128	$d_4$	128

## A.2. Implementation Details

In this section, we detail the exact configuration of our method as well as the competing algorithms evaluated.

**Training Details.** Across our experiments, we use Adam [18] optimizer with default parameters and a batch size of 4 sequences. The semantic segmentation experiments use a fixed learning rate of 0.001 for 100 epochs. For the panoptic segmentation experiments, we start with a higher learning rate of 0.01 for 50 epochs, and decrease it to 0.001 for the last 50 epochs.

**U-TAE.** In Table 4, we report the width of the feature maps outputted by each level of the U-TAE’s encoder and decoder. In both networks, we use the same convolutional block shown in Figure 9 and constituted of one  $3 \times 3$  convolution from the input to the output’s width, and one residual  $3 \times 3$  convolution. In the encoding branch, we use Group Normalisation with 4 groups and Batch Normalisation in the decoding branch.

For the temporal encoding, we chose a L-TAE with 16 heads, and a key-query space of dimension  $d_k = 4$ . We use Group Normalisation with 16 groups at the input and output of the L-TAE, meaning that the inputs of each head are layer-normalized.

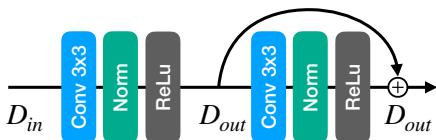


Figure 9: Structure of the convolutional block used in the spatial encoder-decoder network. This block maps a feature map with  $D_{in}$  channels to a feature map with  $D_{out}$  channels.

**Recurrent Models.** We use the same U-Net architecture for our models and *U-BiConvLSTM* and *U-ConvLSTM*, but simply replace the L-TAE by a ConvLSTM or BiConvLSTM respectively. The hidden state’s size of the biConvLSTM is chosen as 32 in both directions, and 64 for the

convLSTM. For the recurrent-convolutional methods *ConvLSTM* and *ConvGRU* not using a U-Net, we set hidden sizes of 160 and 188 respectively.

**3D-Unet.** For this network, we use the official PyTorch implementation of Rustowicz *et al.* [37]. This network is constituted of five successive 3D-convolution blocks with spatial down-sampling after the 2nd and 4th blocks. Each convolutional block doubles the number of channels of the processed feature maps, and the innermost feature maps have a channel dimension set to 128. Leaky ReLu and 3D Batch Normalisation are used across the convolutional blocks of this architecture. The sequence of feature maps is averaged along the temporal dimension to produce the final embedding of the input image sequence. In their implementation, the authors used a linear layer to collapse the temporal dimension, yet this was not a valid option for PASTIS as the sequences have highly variable lengths and the sequence indices do not correspond to the same acquisition date from one sequence to another.

**FPN-ConvLSTM.** For this architecture, the input sequence of images is first mapped to feature maps of channel dimension 64 with two consecutive  $3 \times 3$  convolution layers, followed by Group Normalization and ReLu. A 5-level feature pyramid is then constructed for each date of the sequence by applying to the feature maps 4 different  $3 \times 3$  convolution of respective dilation rates 1, 2, 4 and 8, and computing the spatial average of the feature map. These 5 maps are concatenated along the channel dimension, and processed by a ConvLSTM with a hidden state size of 88. We found it beneficial to use a supplementary convolution before the ConvLSTM to reduce the number of channels of the feature pyramid by a factor 2.

**PaPs module.** In the PaPs module, the saliency and heatmap predictions are obtained with two separate convolutional blocks operating on the high resolution feature map  $d_1$  with 32 channels. These blocks are composed of two convolutional layers of width 32 and 1 respectively. We use Batch Normalisation and ReLu after the first convolution, and a sigmoid after the second.

The 256-dimensional multi-scale feature vector ( $128 + 64 + 32 + 32$ ) is mapped to the shape, class and size predictions by three different MLPs described in Table 5. The inner layers use Batch Normalisation and ReLu activation.

The residual CNN used for shape refinement is composed of three convolutional layers :  $1 \mapsto 16 \mapsto 16 \mapsto 1$ , with ReLu activation and instance normalisation on the first layer only.

**Handling Sequences of Variable Lengths.** All models are trained on batches of sequences of variable length. To

Table 5: Configuration of the four MLPs of PaPs

MLP	Layers	Final Layer
Shape	$256 \mapsto 128 \mapsto S^2$	-
Size	$256 \mapsto 128 \mapsto 2$	Softplus
Class	$256 \mapsto 128 \mapsto 64 \mapsto K$	Softmax

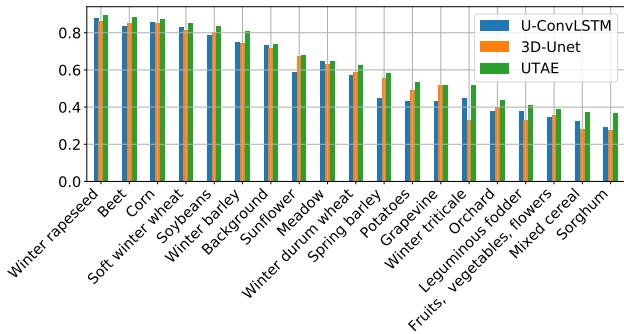


Figure 10: Per class IoU of the three best performing semantic segmentation models. Our U-TAE outperforms the other two approaches on every classes, and brings noticeable improvement on hard classes such as Mixed cereal and Sorghum.

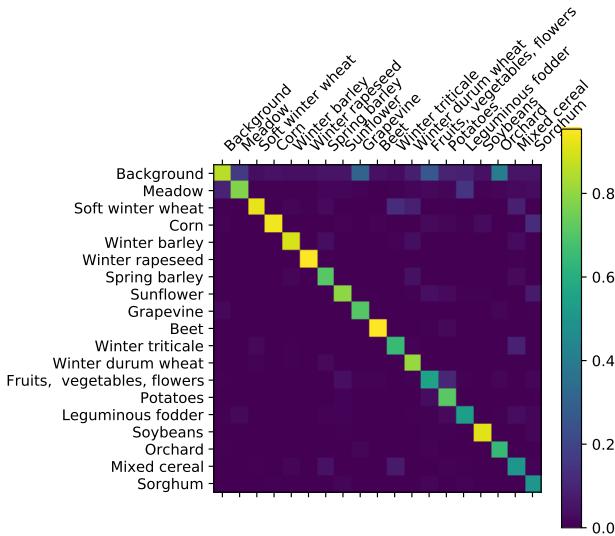


Figure 11: Confusion matrix of U-TAE for semantic segmentation on PASTIS. The color of each pixel at line  $i$  and column  $j$  corresponds to the proportion of samples of the class  $i$  that were attributed to the class  $j$ .

facilitate the handling of batches by the GPU, we append all-zeroes images at the end of shorter sequences to match the length of the longer sequence in the batch. We retain a padding mask to prevent the spatial and temporal encoding of padded values, and to exclude these padded values from

temporal averages.

### A.3. Additional Results

In Figure 10, we show the class-wise performance of the three best performing semantic segmentation models, displaying an improvement of U-TAE compared to the other methods across all crop types. We also show on Figure 11 the confusion matrix of U-TAE. Unsurprisingly, confusions seem to occur between semantically close classes such as different cereal types, or *Sunflower* and *Fruits, Vegetable, Flower*.

In Figure 12, we present qualitative results illustrating the predicted panoptic and semantic segmentations compared to the ground truth. In particular, we show some failure cases in which thin or visually fragmented parcels are not recovered correctly.

In Figure 13, we illustrate the results of the semantic segmentation for our method and three other competing approaches: *3D-Unet*, *U-BiConvLSTM*, and *convGRU*. We show how our multi-scale temporal attention masks allow our predictions to be both pixel-precise and consistent for large parcels.

Finally, we present in Figure 14 an example of inference using a single image from the sequence. As expected for mono-temporal segmentation, the parcel classification is poor. Furthermore, we show a case of a border that is essentially invisible on a single image, but that our full model is able to detect using the entire sequence of satellite images.

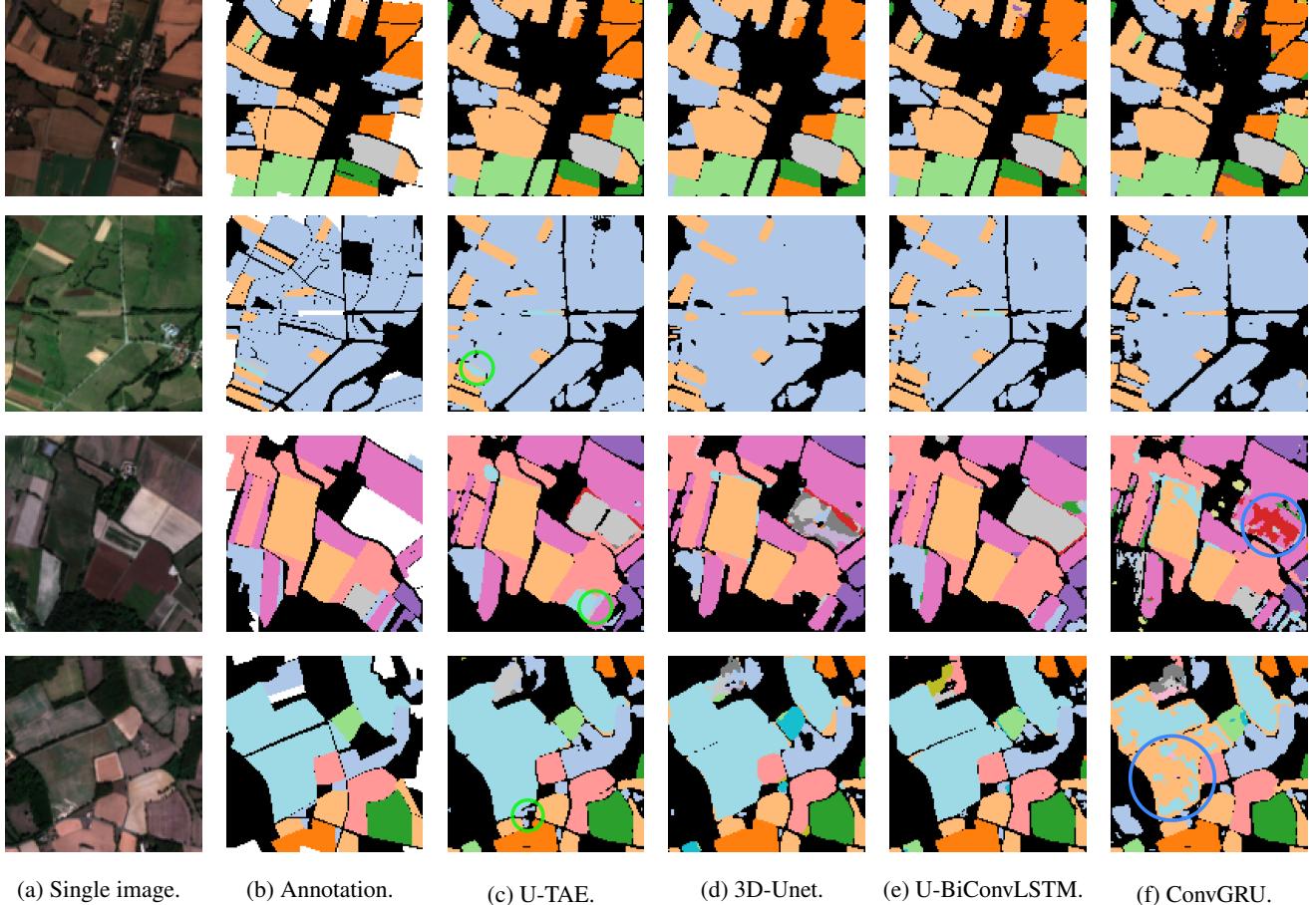
### Acknowledgments

The satellite images used in PASTIS were gathered from THEIA: “Value-added data processed by the CNES for the Theia data cluster using Copernicus data. The treatments use algorithms developed by Theia’s Scientific Expertise Centres.” The annotations of PASTIS were taken from the French LPIS produced by IGN, the French mapping agency. This work was partly supported by ASP, the French Payment Agency.

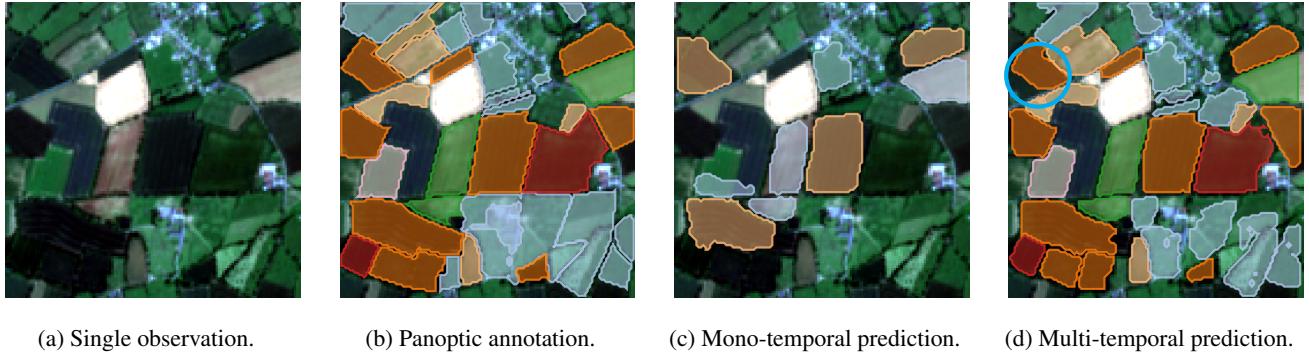


(a) Image from the sequence. (b) Panoptic annotation. (c) Panoptic segmentation. (d) Semantic segmentation.

**Figure 12: Qualitative Panoptic Segmentation Results.** We represent a single image from the sequence using the RGB channels (a), and whose ground truth parcel’s limit and types are known (b). We then represent the parcels predicted by our panoptic segmentation module (c), and the pixelwise prediction of our semantic segmentation module (d). See Figure 7 for the color to crop type correspondence. We highlight with a green circle  $\circlearrowleft$  a large, fragmented parcel declared as one single field. This leads to predictions with low confidence and a low panoptic quality. Conversely, the cyan circle  $\circlearrowright$  highlights such fragmented parcel which is correctly predicted as a single instance. This suggests that our network is able to use the temporal dynamics to recover ambiguous borders. We highlight a failure case with the red circle  $\circlearrowright$ , for which many thin parcels are not properly detected, resulting in a low panoptic quality. We observe that the semantic segmentation model struggles as well for such thin parcels. Finally, we highlight with a blue circle  $\circlearrowleft$  an example in which the panoptic prediction is superior to the semantic segmentation, indicating that detecting parcels’ boundaries and extent can be informative for their classification.



**Figure 13: Qualitative Semantic Segmentation Results.** We represent a single image from the sequence using the RGB channels (a), and whose ground truth parcel's limit and crop type are known (b). We then represent the pixelwise prediction from our approach (c), and for three other competing algorithms (d-f). The different predictions shown on this figure illustrate the importance of the resolution at which temporal encoding is performed. ConvGRU applies a recurrent-convolutional network at the highest resolution, which results in predictions with high spatial variability. As a consequence, the prediction over large parcels are inconsistent (blue circles O). Conversely, U-BiConvLSTM applies temporal encoding to feature maps with a larger receptive field, resulting in more spatially consistent predictions. Yet, this architecture often fails to retrieve small or thin parcels. In contrast, our U-TAE produces spatially consistent predictions on large parcels, while being able to retrieve such small parcels (green circles O). 3D-Unet also uses temporal encoding at different resolution levels, yet fails to recover these small parcels.



(a) Single observation. (b) Panoptic annotation. (c) Mono-temporal prediction. (d) Multi-temporal prediction.

**Figure 14: Mono-temporal Panoptic Segmentation.** We train our mono-temporal model on a single image (a), with panoptic annotation (b). We then compare the results of the mono-temporal model in (c) with the results our full model when performing inference on the full length sequence (d) from which the single patch (a) is drawn. First, we observe that many parcels are not detected by the mono-temporal model, indicating an overall low predicted quality. Second, we can see that most detected parcels are misclassified by the mono-temporal model. This is in accordance with the low semantic segmentation score of the mono-temporal model: crop types are hard to distinguish from a single observation. Last, adjacent parcels with no clear borders are predicted as a single parcel, when the multi-temporal model is able to differentiate between the two parcels (cyan circle O). This illustrates how using SITS instead of single images can help resolve ambiguous parcels delineation.