
Using Convolutional LSTMs for Cloud-Robust Segmentation of Remote Sensing Imagery

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Dynamic spatiotemporal processes on the Earth can be observed by an increasing
2 number of optical Earth observation satellites that measure spectral reflectance at
3 multiple spectral bands in regular intervals. Clouds partially covering the surface
4 is an omnipresent challenge for the majority of remote sensing approaches that
5 are not robust regarding cloud coverage. In these approaches, clouds are typically
6 handled by cherry-picking cloud-free observations or by pre-classification of cloudy
7 pixels and subsequent masking. In this work, we demonstrate the robustness of
8 a straightforward *convolutional long short-term memory* network for vegetation
9 classification using all available cloudy and non-cloudy satellite observations. We
10 visualize the internal gate activations within the recurrent cells and find that, in
11 some cells, modulation and input gates close on cloudy pixels. This indicates that
12 the network has internalized a cloud-filtering mechanism without being specifically
13 trained on cloud labels. The robustness regarding clouds is further demonstrated
14 by experiments on sequences with varying degrees of cloud coverage where our
15 network achieved similar accuracies on all cloudy and non-cloudy datasets. Overall,
16 our results question the necessity of sophisticated pre-processing pipelines if robust
17 classification methods are utilized.

18 Supplementary material can be accessed via <https://tinyurl.com/NIPS18ST-supplement>

19 1 Introduction

20 A wide range of dynamic spatiotemporal processes of the Earth can be observed with remote sensing
21 satellites that revisit the same position on Earth at discrete time intervals. Seasonal vegetation life-
22 cycles and other land cover dynamics are typically monitored at weekly intervals at spatial resolutions
23 of several meters that allow distinguishing large single objects. Imagery acquired by these optical
24 satellites is, however, regularly covered by clouds. These coverages are typically addressed by either
25 selecting exclusively cloud-free observations or masking and removing clouds by computationally
26 sophisticated pre-processing pipelines. We investigate the robustness of convolutional long short-term
27 memory networks [8] with regard to temporal noise induced by cloud coverage for remote sensing
28 imagery.

29 2 Related Work

30 Clouds distinguish themselves from ground pixels by their the high reflectance compared to ground
31 pixels. Decision-tree based models [4, 10, 2] applied on expert-designed features are used for many
32 remote sensing applications. The `fmask` algorithm [10] and improved versions [9, 1] additionally
33 implement a projection of the detected cloud on the surface as initialization to additionally predict
34 the shadow casted by the cloud. Other approaches extract features from a time series and utilize the

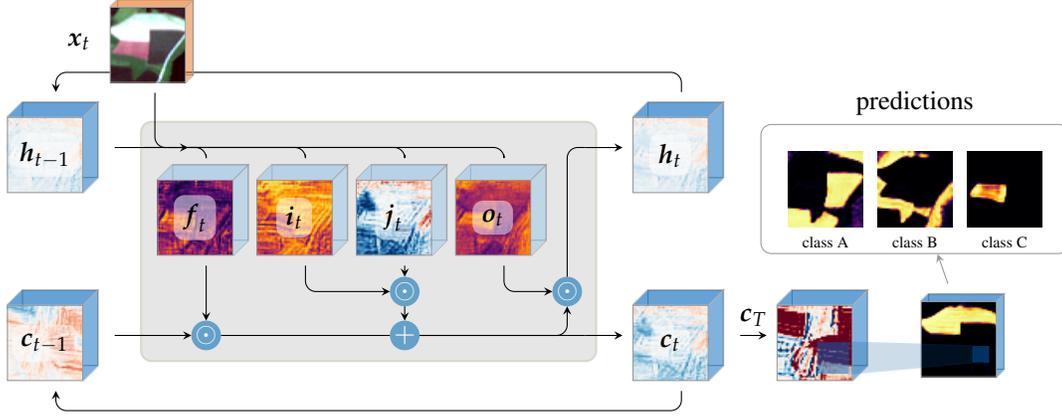


Figure 1: Illustration of the two-layer convolutional long short-term memory network (LSTM) topology. Each input image x_t of T images is passes sequentially to the LSTM encoder that extracts classification relevant features to the internal cell state tensor c_T . A second convolutional layer compresses the dimensionality to the number of classes yielding activations per class.

35 sudden increase in reflectance to identify cloudy pixels [2]. *Convolutional neural networks (CNNs)*
 36 have also shown to compare well [11] indicating that these features can be learned by deep methods.
 37 These methods have proven beneficial in the last years and are widely implemented in remote sensing
 38 approaches. However, masking single pixels by a pre-classification introduces an additional layer of
 39 procedural complexity and raises the question of how to treat these pixels accordingly in the designed
 40 framework. Overall, cloud-filtering remains a pre-processing necessity for most remote sensing
 41 approaches that are prone to fail in the presence of data noise.

42 Similar to our work, only a few approaches have tried to design robust methods that do not require
 43 this additional pre-classification step. One approach added pre-classified cloud labels as additional
 44 prediction targets that allowed the implemented network to distinguish cloud from ground classes
 45 [7]. Also, ensemble-based methods of supervised classifiers have shown robustness regarding the
 46 appearance of clouds [6].

47 3 Method

48 In this section, we outline the theoretical basis of *convolutional long short-term memory (convLSTM)*
 49 networks utilized in this work and detail the employed network topology.

50 3.1 Convolutional Long Short-term Neural Networks

51 *Long short-term memory networks (LSTM)* [3] implement internal gates to control the gradient-
 52 flow through time and an additional container for long-term memory c_t . This yields the LSTM
 53 update $(h_t, c_{t-1}) \leftarrow (x_t, h_{t-1}, c_{t-1})$ that map an input x_t and short-term context h_{t-1} to a hidden
 54 representation h_t . Additionally, a long-term cell state c_{t-1} is updated to c_t at each iteration and can
 55 store information for a theoretically unlimited number of iteration. Three gates control the update of
 56 the cell state

$$c_t \leftarrow c_{t-1} \odot f_t + i_t \odot j_t \quad (1)$$

57 by element-wise multiplication denoted by the *Hadamard* operator \odot . The forget gate $f_t =$
 58 $\sigma(x_t * \theta_{fx} + h_{t-1} * \theta_{fh} + \mathbf{1})$ evaluates the influence of the previous cell state c_{t-1} with a sigmoidal
 59 $\sigma(\cdot) \in]0, 1[$ activation function. The input and modulation gates

$$i_t = \sigma(x_t * \theta_{ix} + h_{t-1} * \theta_{ih}), \text{ and } j_t = \tanh(x_t * \theta_{jx} + h_{t-1} * \theta_{jh}) \quad (2)$$

60 are element-wise multiplied for the cell state update. The output gate $o_t =$
 61 $\tanh(x_t * \theta_{ox} + h_{t-1} * \theta_{oh})$ determines with the cell state the current cell output $h_t \leftarrow o_t \odot c_t$.
 62 Convolutional recurrent networks implement a convolution, denoted by $*$, instead of a matrix mul-
 63 tiplication of the formulation of recurrent networks. Each respective gate activation, referred by
 64 subscripts f, i, j, o, is controlled by trainable weights for input $\theta_{fx}, \theta_{ix}, \theta_{jx}, \theta_{ox} \in \mathbb{R}^{k \times k \times d \times r}$ and

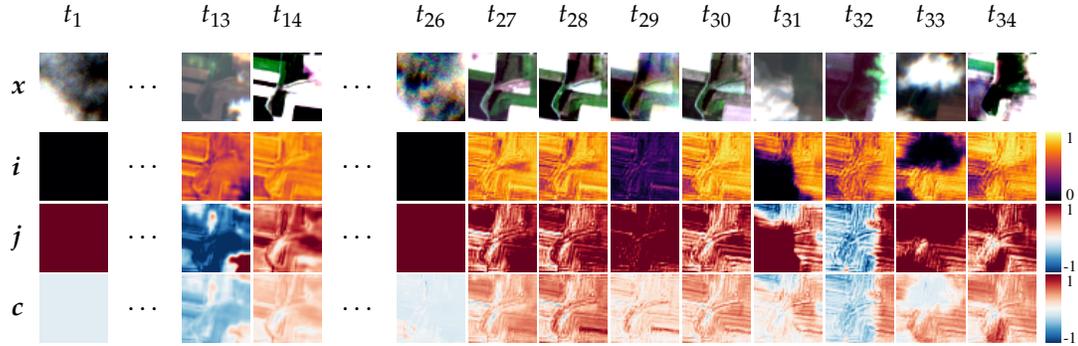


Figure 2: Activations of the cell state and selected gates of one convolutional LSTM cell that indicate that the cell has internalized a cloud-filtering scheme. The input gate i in this specific cell seems to be assigned values of zero on cloudy pixels as seen at steps $t = 13, 26, 31, 33$.

65 hidden representation $\theta_{fh}, \theta_{ih}, \theta_{jh}, \theta_{oh} \in \mathbb{R}^{k \times k \times r \times r}$ where d represents the dimensional depth of
 66 the input image, k the convolutional kernel size, and r a hyper-parameter determining the number
 67 of recurrent cells by setting dimensionality of the hidden states. With this change, image data of
 68 certain width, height and depth can be processed where convolutions partially connect the local pixel
 69 neighborhoods between layers.

70 3.2 Network architecture

71 We utilize this single-layer convolutional LSTM neural network to encode a sequence of T satellite
 72 images to the fixed length representation c_T , as illustrated in Fig. 1. To balance the influence of
 73 the sequence order, we also encode the reversed sequence and append the final cell states. In initial
 74 published experiments, we found 256 recurrent cells to be optimal and used this hyper-parameter of
 75 dimensionality for the hidden tensors within the LSTM network.

76 After sequential encoding, the combined cell state is passed to a second convolution layer that
 77 compresses the dimensionality from 2×256 hidden dimensions to the number of classes. Applying
 78 softmax normalization produces activations that can be interpreted as network-confidences per class
 79 and are illustrated in the figure. We used convolutional kernels of $3 \times 3px$ in size throughout the
 80 network. To train, we evaluate the cross-entropy between the last layer and a one-hot representation
 81 of the ground truth labels. The influence of each weight on the evaluated loss is determined by
 82 back-propagated gradients and iterative adjustments are determined by the Adam optimizer[5].

83 4 Results

84 The primary objective for this network was to identify the type of cultivated crops in an area of
 85 interest of $100 \text{ km} \times 40 \text{ km}$. Hence, we trained our network end-to-end on label data describing the
 86 crop-type on distinct field parcels. No additional label information about cloud coverages was used.
 87 We used a sequence of 46 SENTINEL 2 satellite observations from the year 2016 for this objective.
 88 This satellite measures the reflectances of 13 spectral bands at 10 m, 20 m, and 60 m resolution. To
 89 harmonize the data sources, we bi-linearly interpolated these to 10 m resolution and rasterized the
 90 crop labels accordingly. In this section, we evaluate the robustness of the proposed network regarding
 91 cloud coverage.

92 4.1 Long-short term memory cell activations

93 We trained the network on field crop labels for thirty epochs using raw sequences of cloudy and
 94 non-cloudy observations. The top row of Fig. 2 shows an particular example input sequence of
 95 $T = 34$ images of $48 \times 48px$ in size. The following rows illustrate activations of the internal
 96 convolutional LSTM gates i, j and cell state c given each input element. While all of the 256 recurrent
 97 cells likely contribute to the classification decision, only few were visually interpretable similar to
 98 the shown example. Following Eq. (1), the cell state is updated with new information based on

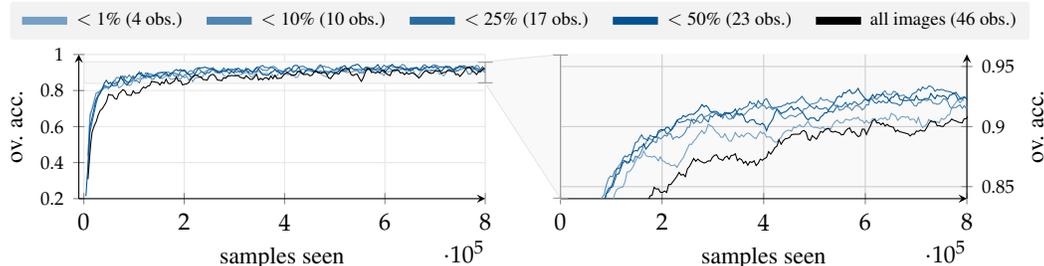


Figure 3: Overall accuracy over the training progress of the same convolutional LSTM network topology trained on datasets with different degrees of cloud coverage.

99 the input and modulation gates i, j . The activations in Fig. 2 of these gates in the second and third
 100 row show that the input gate i approaches zero at pixels that are covered by clouds. This effect can
 101 be observed at time steps $t = \{13, 26, 31, 33\}$. At time $t = 32$ the input gate seems unchanged,
 102 however, the modulation gate j changes sign. Overall, these results indicate that the convolutional
 103 recurrent network has internalized a mechanism for cloud-filtering. More activation examples can be
 104 obtained from the supplementary material.

105 4.2 Experiments with varying degrees of cloud coverage

106 In this experiment, we trained the network on datasets with different degrees of cloud coverage. To
 107 determine the cloud coverage per observation, all satellite observations have been pre-processed using
 108 the `fmask` algorithm implemented in the `Sen2Cor` software, as being common practice in remote
 109 sensing. This yields a per-pixel cloud classification label. With this, a cloud coverage pixel ratio per
 110 observation can be calculated. Based on this, several sub-datasets have been created with either all 46
 111 observations, the 26 images covered with less than 50%, 17 images with less than 25%, 10 with less
 112 than 10%, and 4 completely cloud-free images.

113 We trained the network on these pre-filtered datasets. In Fig. 3 one can observe that the overall
 114 accuracy over the training process remains remarkably similar for all of the sub-sampled datasets.
 115 The right graph shows a zoomed view and reveals some differences between the dataset performances.
 116 Datasets containing observations and the four completely cloud-free observations have been slightly
 117 worse classified than the intermediate ones of 10%, 25%, and 40% coverage. It seems that the
 118 rejection of completely covered observations was beneficial as indicated by the slightly worse
 119 accuracy on the dataset of all observations. Similarly, the four cloud-free observations may have
 120 missed some characteristic vegetation-related events. Intuitively, these results show a trade-off
 121 between restrictions on cloud coverage and sequence length and demonstrate that cherry-picking
 122 single cloud-free observations may lead to inferior classification accuracy. Overall, these results
 123 demonstrate the robustness of the convolutional long short-term memory network to handle data
 124 containing temporal noise induced by cloud coverage.

125 5 Conclusion

126 Noise in temporal data is a common challenge for a variety of disciplines. In this work, we focused on
 127 noise induced by cloud coverage in multi-temporal remote sensing imagery. Most Earth observation
 128 approaches either select few completely cloud-free observations or use a pre-classification to mask
 129 cloudy pixels. The experiments of this work showed that this cloud-induced temporal noise can be
 130 learned purely from the data in an end-to-end fashion with an appropriate model design. We utilized
 131 long short-term memory cells that are popularly used in natural language processing tasks, such as
 132 translation or text generation in a straightforward two-layer network. Our results demonstrate this
 133 model design is able to consistently extract the classification-relevant features from observations
 134 between cloudy observations.

135 Our work questions the necessity of sophisticated, partly hand-crafted pre-processing pipelines for
 136 remote sensing imagery. These results show that methods perform well in the seemingly unrelated
 137 field of remote sensing and Earth observation. To encourage further research with spatiotemporal
 138 data in remote sensing and related fields, we will publish source code and data upon acceptance.

139 **References**

- 140 [1] David Frantz, Achim Röder, Thomas Udelhoven, and Michael Schmidt. Enhancing the de-
141 tectability of clouds and their shadows in multitemporal dryland landsat imagery: extending
142 fmask. *IEEE Geoscience and Remote Sensing Letters*, 12(6):1242–1246, 2015.
- 143 [2] Olivier Hagolle, Mireille Huc, D Villa Pascual, and Gérard Dedieu. A multi-temporal method
144 for cloud detection, applied to formosat-2, venus, landsat and sentinel-2 images. *Remote
145 Sensing of Environment*, 114(8):1747–1755, 2010.
- 146 [3] Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*,
147 9(8):1735–1780, 1997.
- 148 [4] André Hollstein, Karl Segl, Luis Guanter, Maximilian Brell, and Marta Enesco. Ready-to-
149 use methods for the detection of clouds, cirrus, snow, shadow, water and clear sky pixels in
150 sentinel-2 msi images. *Remote Sensing*, 8(8):666, 2016.
- 151 [5] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint
152 arXiv:1412.6980*, 2014.
- 153 [6] Chuc Duc Man, Thuy Thanh Nguyen, Hung Quang Bui, Kristofer Lasko, and Thanh Nhat Thi
154 Nguyen. Improvement of land-cover classification over frequently cloud-covered areas using
155 landsat 8 time-series composites and an ensemble of supervised classifiers. *International
156 Journal of Remote Sensing*, 39(4):1243–1255, 2018.
- 157 [7] Marc Rußwurm and Marco Körner. Temporal vegetation modelling using long short-term
158 memory networks for crop identification from medium-resolution multi-spectral satellite images.
159 In *CVPR Workshops*, pages 1496–1504, 2017.
- 160 [8] Shi Xingjian, Zhourong Chen, Hao Wang, Dit-Yan Yeung, Wai-Kin Wong, and Wang-chun
161 Woo. Convolutional lstm network: A machine learning approach for precipitation nowcasting.
162 In *Advances in neural information processing systems*, pages 802–810, 2015.
- 163 [9] Zhe Zhu, Shixiong Wang, and Curtis E Woodcock. Improvement and expansion of the fmask
164 algorithm: Cloud, cloud shadow, and snow detection for landsats 4–7, 8, and sentinel 2 images.
165 *Remote Sensing of Environment*, 159:269–277, 2015.
- 166 [10] Zhe Zhu and Curtis E Woodcock. Object-based cloud and cloud shadow detection in landsat
167 imagery. *Remote sensing of environment*, 118:83–94, 2012.
- 168 [11] Anze Zupanc. Improving cloud detection with ma-
169 chine learning. [https://medium.com/sentinel-hub/
170 improving-cloud-detection-with-machine-learning-c09dc5d7cf13](https://medium.com/sentinel-hub/improving-cloud-detection-with-machine-learning-c09dc5d7cf13), 2017.