

# Cloud-Robust Classification of Remote Sensing Time Series

Φ-week 2019

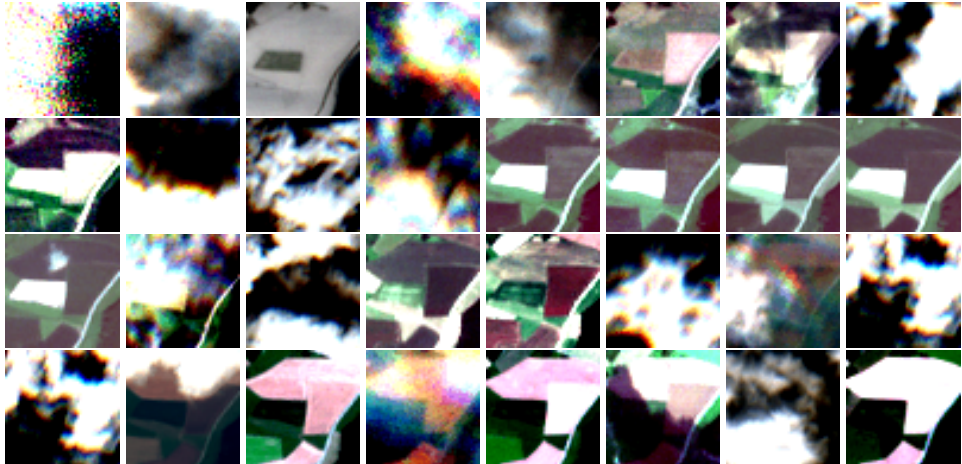
Marc Rußwurm, Marco Körner

10th September 2019, ESA ESRIN, Frascati, Italy



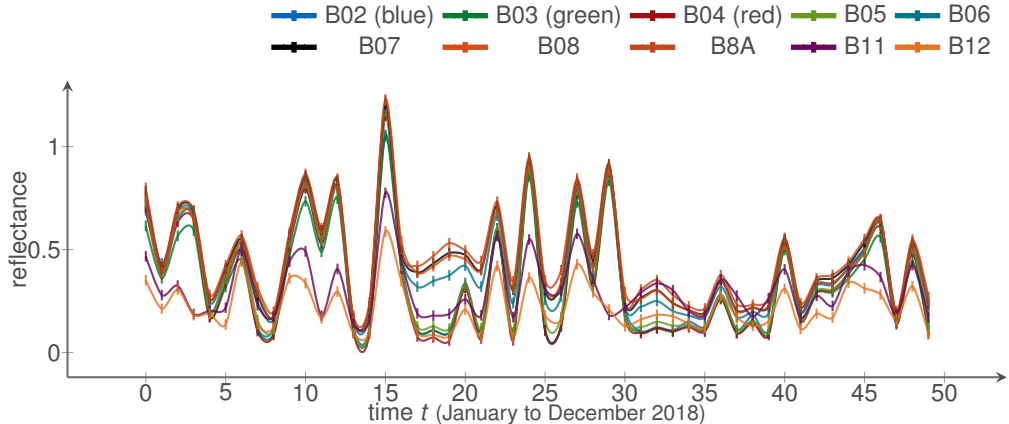


## Cloud coverage



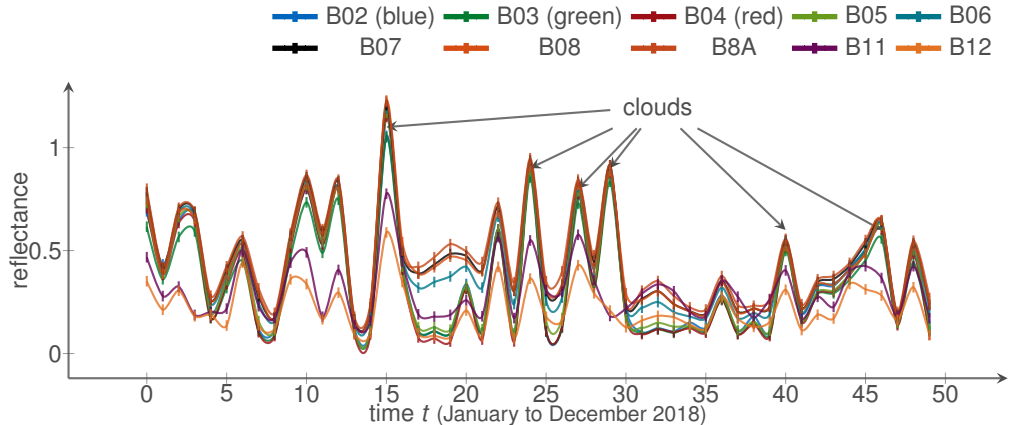
# Satellite Time Series

Sentinel 2 (raw), mean-aggregated pixels of a meadow field parcel



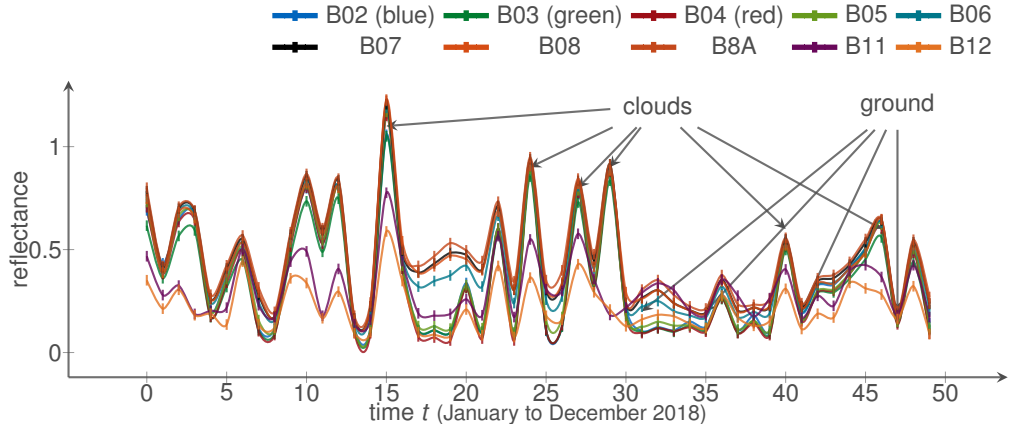
# Satellite Time Series

Sentinel 2 (raw), mean-aggregated pixels of a meadow field parcel

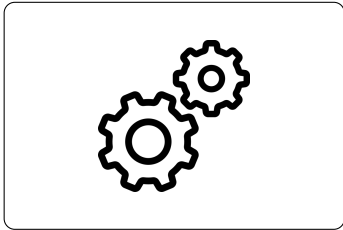


# Satellite Time Series

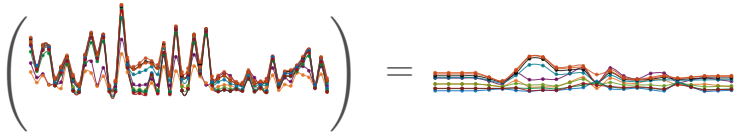
Sentinel 2 (raw), mean-aggregated pixels of a meadow field parcel



## Data Preprocessing



preprocessing



$f_{\theta_{\text{sel}}}(X)$  temporal selection (not considering winter period) where  $\theta_{\text{sel}} = \{t_{\text{start}}, t_{\text{end}}\}$

$f_{\theta_{\text{atm}}}(X)$  atmospheric correction ( $X_{\text{top-of-atmosphere}} \rightarrow X_{\text{bottom-of-atmosphere}}$ )

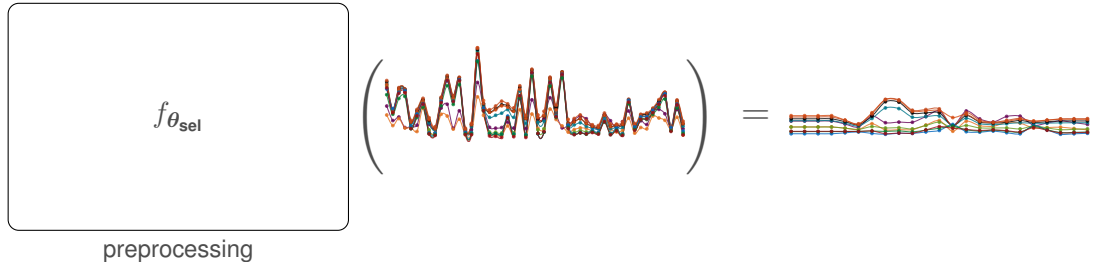
$f_{\theta_{\text{cl}}}(X)$  cloud/cloud-shadow classification (F-Mask, MAJA, CNNs, Cloud Clustering (go FDL!))

$f_{\text{int}}(X)$  temporal interpolation to generate equal sample times

$f_{\theta_{\dots}}$  **many more problem-specific chained building blocks**



# Data Preprocessing



$f_{\theta_{sel}}(X)$  temporal selection (not considering winter period) where  $\theta_{sel} = \{t_{start}, t_{end}\}$

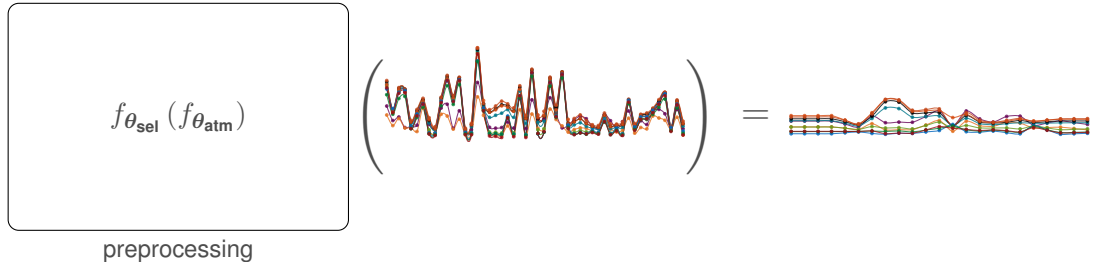
$f_{\theta_{atm}}(X)$  atmospheric correction ( $X_{top-of-atmosphere} \rightarrow X_{bottom-of-atmosphere}$ )

$f_{\theta_{cl}}(X)$  cloud/cloud- shadow classification (F-Mask, MAJA, CNNs, Cloud Clustering (go FDL!))

$f_{int}(X)$  temporal interpolation to generate equal sample times

$f_{\theta_{...}}$  many more problem-specific chained building blocks

# Data Preprocessing



$f_{\theta_{sel}}(X)$  temporal selection (not considering winter period) where  $\theta_{sel} = \{t_{start}, t_{end}\}$

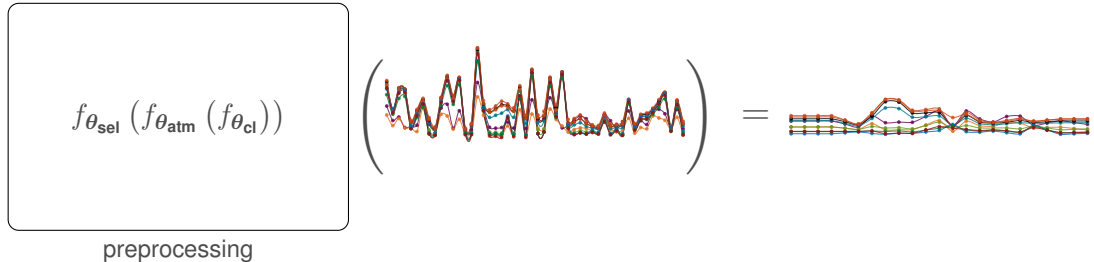
$f_{\theta_{atm}}(X)$  atmospheric correction ( $X_{top-of-atmosphere} \rightarrow X_{bottom-of-atmosphere}$ )

$f_{\theta_{cl}}(X)$  cloud/cloud- shadow classification (F-Mask, MAJA, CNNs, Cloud Clustering (go FDL!))

$f_{int}(X)$  temporal interpolation to generate equal sample times

$f_{\theta_{...}}$  many more problem-specific chained building blocks

## Data Preprocessing



$f_{\theta_{\text{sel}}}(\mathbf{X})$  temporal selection (not considering winter period) where  $\theta_{\text{sel}} = \{t_{\text{start}}, t_{\text{end}}\}$

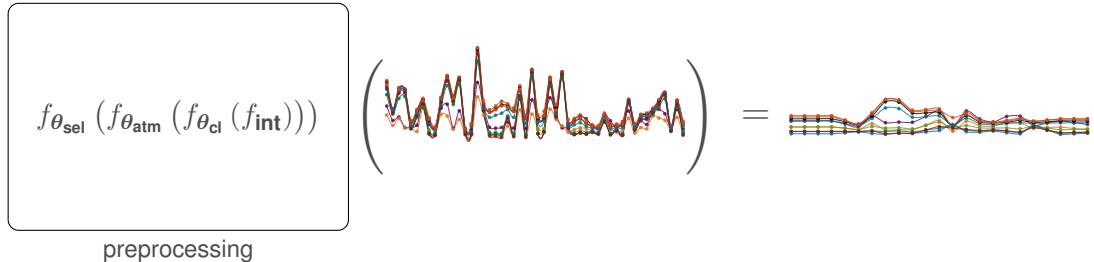
$f_{\theta_{\text{atm}}}(\mathbf{X})$  atmospheric correction ( $\mathbf{X}_{\text{top-of-atmosphere}} \rightarrow \mathbf{X}_{\text{bottom-of-atmosphere}}$ )

$f_{\theta_{\text{cl}}}(\mathbf{X})$  cloud/cloud-shadow classification (F-Mask, MAJA, CNNs, Cloud Clustering (go FDL!))

$f_{\text{int}}(\mathbf{X})$  temporal interpolation to generate equal sample times

$f_{\theta_{\dots}}$  many more problem-specific chained building blocks

## Data Preprocessing



$f_{\theta_{sel}}(X)$  temporal selection (not considering winter period) where  $\theta_{sel} = \{t_{start}, t_{end}\}$

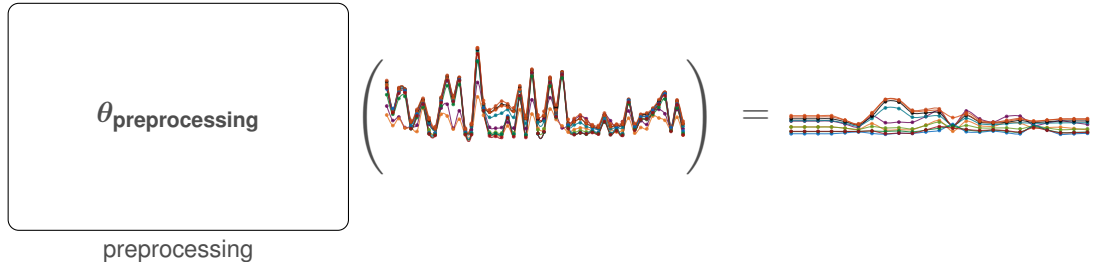
$f_{\theta_{atm}}(X)$  atmospheric correction ( $X_{top-of-atmosphere} \rightarrow X_{bottom-of-atmosphere}$ )

$f_{\theta_{cl}}(X)$  cloud/cloud-shadow classification (F-Mask, MAJA, CNNs, Cloud Clustering (go FDL!))

$f_{int}(X)$  temporal interpolation to generate equal sample times

$f_{\theta_{...}}$  many more problem-specific chained building blocks

# Data Preprocessing



$f_{\theta_{\text{sel}}}(\mathbf{X})$  temporal selection (not considering winter period) where  $\theta_{\text{sel}} = \{t_{\text{start}}, t_{\text{end}}\}$

$f_{\theta_{\text{atm}}}(\mathbf{X})$  atmospheric correction ( $\mathbf{X}_{\text{top-of-atmosphere}} \rightarrow \mathbf{X}_{\text{bottom-of-atmosphere}}$ )

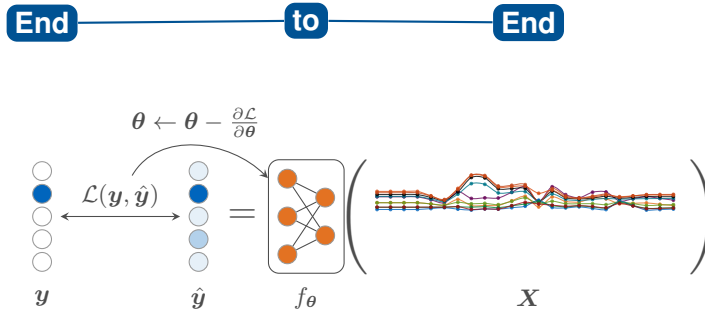
$f_{\theta_{\text{cl}}}(\mathbf{X})$  cloud/cloud-shadow classification (F-Mask, MAJA, CNNs, Cloud Clustering (go FDL!))

$f_{\text{int}}(\mathbf{X})$  temporal interpolation to generate equal sample times

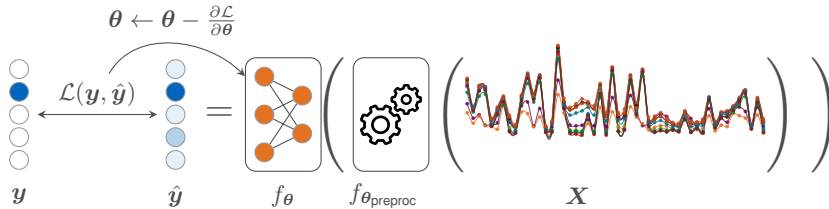
$f_{\theta_{\dots}}$  **many more problem-specific chained building blocks**

Deep Learning Models are trained differently...

# End-to-End Learning of Deep Neural Networks

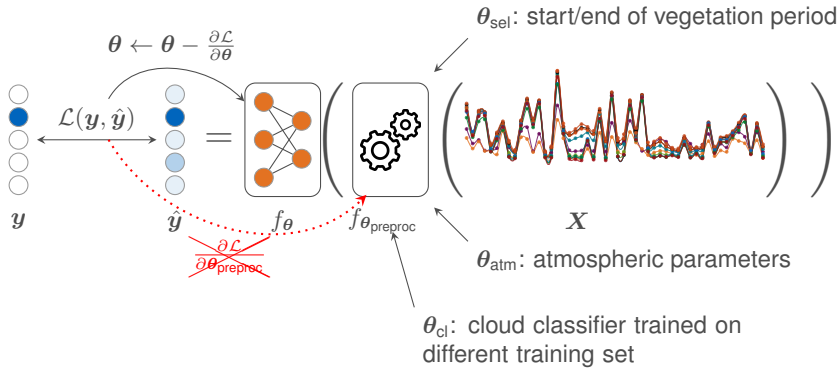


# End-to-End Learning of Deep Neural Networks

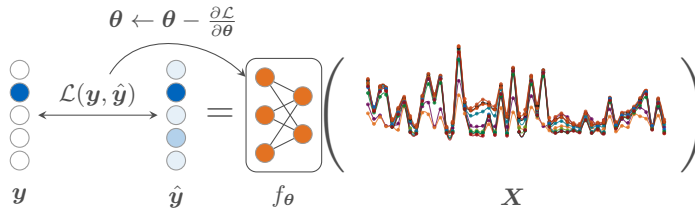




# End-to-End Learning of Deep Neural Networks



# End-to-End Learning of Deep Neural Networks

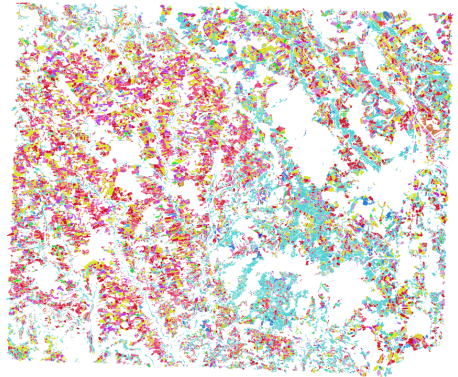


Let's look at some quantitative results...

# Crop Type Dataset northern Bavaria

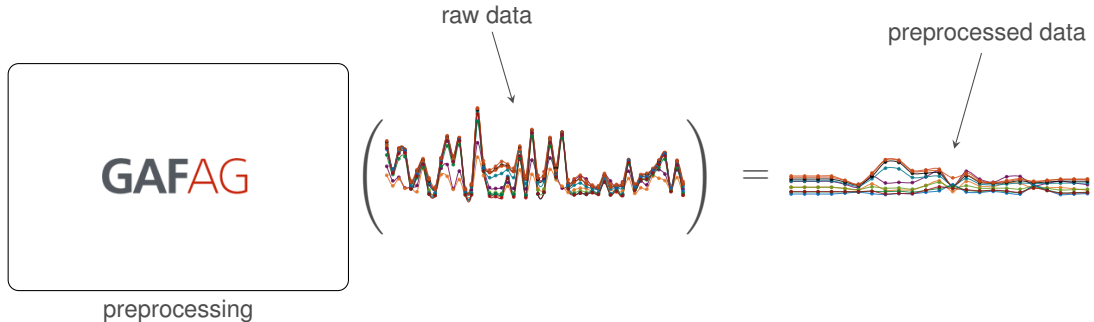
Common project with **GAFAG**

crop type labels by the  
**Bavarian Ministry of Agriculture**  
**49k field parcels** of 2018  
**34 crop categories**



Parcels colored by crop type (40 km × 40 km)

## Two Datasets: Raw and Preprocessed from the same Examples



In this case: Preprocessing Engine of **GAFAG**

## Four state-of-the-art deep Models for Time Series Classification

|            | LSTM-RNN <sup>1</sup> | Transformer <sup>1</sup> | MS-ResNet <sup>3</sup> | TempCNN <sup>4</sup> |
|------------|-----------------------|--------------------------|------------------------|----------------------|
| Mechanism  | Recurrence            | Self-Attention           | Convolution            | Convolution          |
| Parameters | 100k                  | 600k                     | 2M                     | 433k                 |

<sup>1</sup> **Hochreiter, S., & Schmidhuber, J. (1997).** Long short-term memory. Neural computation, 9(8), 1735-1780.

<sup>2</sup> **Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., & Polosukhin, I. (2017).** Attention is all you need. In Advances in neural information processing systems (pp. 5998-6008).

<sup>3</sup> **Wang, F., Han, J., Zhang, S., He, X., & Huang, D. (2018).** Csi-net: Unified human body characterization and action recognition. arXiv preprint arXiv:1810.03064.

<sup>4</sup> **Pelletier, C., Webb, G. I., & Petitjean, F. (2019).** Temporal convolutional neural network for the classification of satellite image time series. Remote Sensing, 11(5), 523.

## Preprocessed versus Raw Data

| <b>accuracy</b> | RNN (LSTM) <sup>2</sup> | Transformer <sup>3</sup> | MS-ResNet <sup>1</sup>  | TempCNN <sup>4</sup>    |
|-----------------|-------------------------|--------------------------|-------------------------|-------------------------|
| preprocessed    | <b>.804</b> $\pm .0031$ | .804 $\pm .0011$         | <b>.849</b> $\pm .0041$ | <b>.836</b> $\pm .0012$ |
| raw             | .801 $\pm .0026$        | <b>.842</b> $\pm .0043$  | .836 $\pm .0033$        | .799 $\pm .0027$        |
| $\Delta$        | .003 $\pm .0041$        | -.038 $\pm .0045$        | .013 $\pm .0055$        | .038 $\pm .0029$        |

| <b>kappa</b> | RNN (LSTM) <sup>2</sup> | Transformer <sup>3</sup> | MS-ResNet <sup>1</sup>  | TempCNN <sup>4</sup>    |
|--------------|-------------------------|--------------------------|-------------------------|-------------------------|
| preprocessed | <b>.759</b> $\pm .0037$ | .759 $\pm .0017$         | <b>.816</b> $\pm .0048$ | <b>.799</b> $\pm .0015$ |
| raw          | .756 $\pm .0037$        | <b>.808</b> $\pm .0052$  | .800 $\pm .0039$        | .750 $\pm .0036$        |
| $\Delta$     | .003 $\pm .0048$        | -.049 $\pm .0054$        | .016 $\pm .0060$        | .049 $\pm .0036$        |

### Experiments:

mean  $\pm$  standard deviation of 10 models  
trained from different random initialization

### Findings:

remarkably similar results on preprocessed  
and raw data ( $\Delta \leq 5\%$ )

# Attention



Photo by zhengtao tang on Unsplash



# Self-Attention in Deep Learning



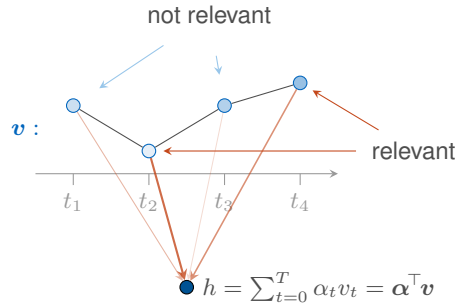
Photo by zhengtao tang on Unsplash

# Attention

Given **values**  $v$  as a sequence of observations.

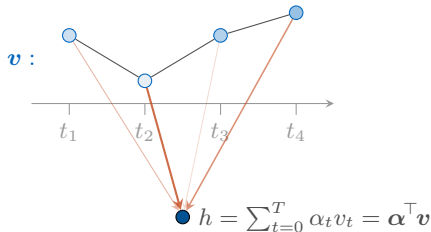
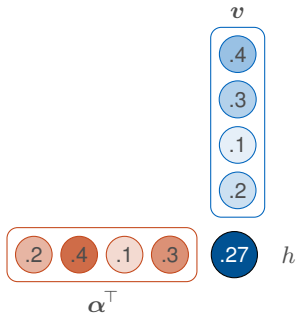
We want to **calculate a result**  $h$  based only on **classification-relevant** observations.

This is realized by an weighted sum over **attention scores**  $\alpha$



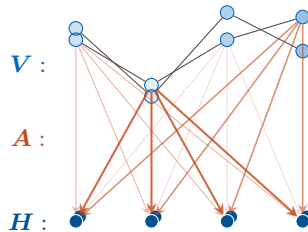
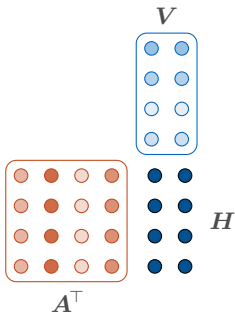
$$h = \text{Attention}(\alpha, v) = \alpha^\top v = \sum_{t=0}^T \alpha_t v_t, \quad \alpha \in [0, 1]^{T=4}, v \in \mathbb{R}^T$$

# Attention



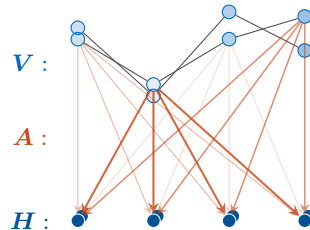
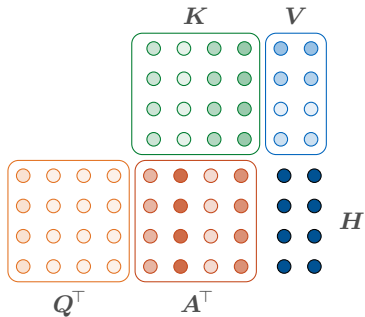
$$h = \text{Attention}(\alpha, v) = \alpha^\top v = \sum_{t=0}^T \alpha_t v_t, \quad \alpha \in [0, 1]^{T=4}, v \in \mathbb{R}^T$$

# Attention



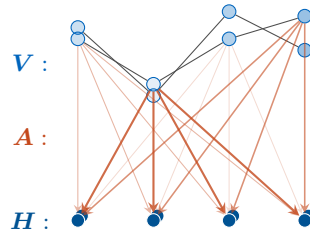
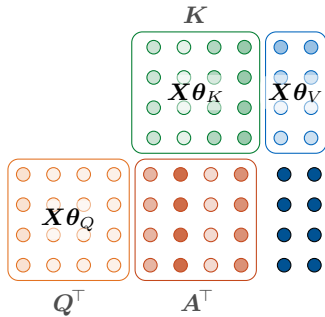
$$H = \text{Attention}(A, V) = A^T V, \quad A \in [0, 1]^{T_{in} \times T_{out}}, \quad V \in \mathbb{R}^{T \times D_v}$$

# Attention



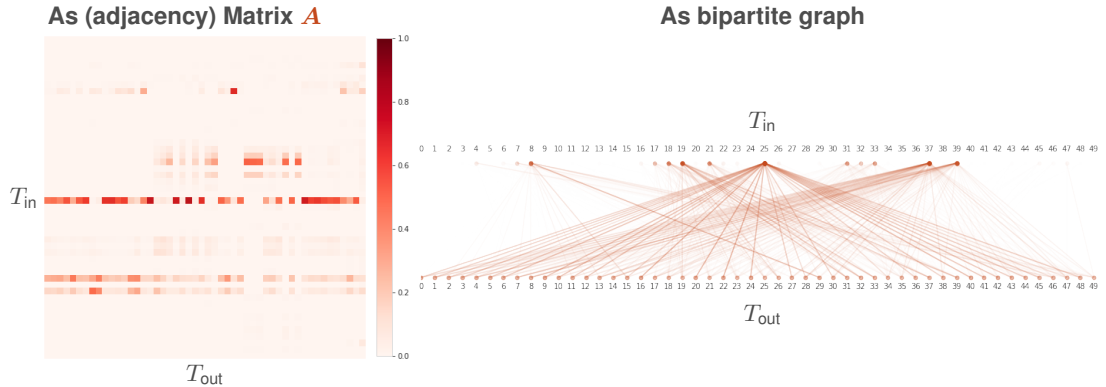
$$\text{Attention}(\mathbf{K}, \mathbf{Q}, \mathbf{V}) = \overbrace{\text{softmax}(\mathbf{Q}^T \mathbf{K})}^{\mathbf{A}^T} \mathbf{V}, \quad \mathbf{V} \in \mathbb{R}^{T \times D_v}, \mathbf{Q}, \mathbf{K} \in \mathbb{R}^{D_k \times T}, \mathbf{A} \in \mathbb{R}^{T_{in} \times T_{out}}$$

# Self-Attention



$$\text{Self-Attention}_{\theta}(X) = \text{Attention}(\underbrace{X\theta_K}_K, \underbrace{X\theta_Q}_Q, \underbrace{X\theta_V}_V) = \text{softmax}((\underbrace{X\theta_Q}_Q)(\underbrace{X\theta_K}_K))(\underbrace{X\theta_V}_V)$$

# Visualizing Attention Scores of a Pretrained Transformer Model



# Attention Scores in Context of Input Time Series

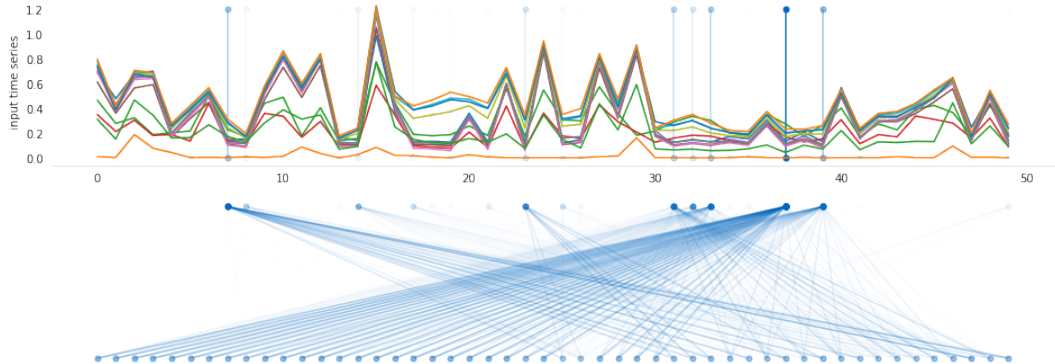
Head 1





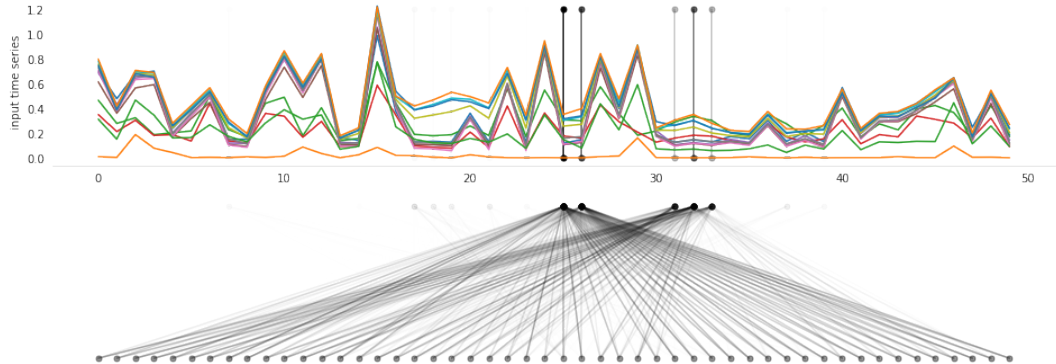
# Attention Scores in Context of Input Time Series

Head 2



# Attention Scores in Context of Input Time Series

Head 3



Let's summarize...

## Summary

What did we look at?

## Summary

What did we look at?

**end-to-end learning** in combination with **preprocessing**

## Summary

What did we look at?

**end-to-end learning** in combination with **preprocessing**

quantitative **results** on models with **raw** and **preprocessed** data

## Summary

What did we look at?

**end-to-end learning** in combination with **preprocessing**

quantitative **results** on models with **raw** and **preprocessed** data

a qualitative example on the **self-attention** mechanism

## Conclusion

What were the outcomes of these experiments?



## Conclusion

What were the outcomes of these experiments?

**classification results** on raw and preprocessed data  
were **remarkably similar**

## Conclusion

What were the outcomes of these experiments?

**classification results** on raw and preprocessed data were **remarkably similar**

do we **need** extensive **preprocessing** for deep learning models on time series data?

## Conclusion

How did deep learning models get robust to noise (e.g., clouds)?

## Conclusion

How did deep learning models get robust to noise (e.g., clouds)?

we saw how **self-attention** is used to **focus** on cloud-free observations

## Conclusion

How did deep learning models get robust to noise (e.g., clouds)?

we saw how **self-attention** is used to **focus** on cloud-free observations


gates in recurrent networks work similar (see previous work)<sup>1,2</sup>


<sup>1</sup>Rußwurm, M., & Körner, M. (2018). Multi-temporal land cover classification with sequential recurrent encoders. **ISPRS International Journal of Geo-Information**, 7(4), 129.

<sup>2</sup>Rußwurm, M., & Körner, M. (2018). Convolutional LSTMs for Cloud-Robust Segmentation of Remote Sensing Imagery. **NeurIPS2018 Workshop on Spatiotemporal Modeling**

Thank you

**Marc Rußwurm & Marco Körner**  
TUM Chair of Remote Sensing Technology  
Computer Vision Research Group

clone the  repo to the attention experiments!  
**github.com/marccoru/phiweek19**  [Open in Colab](#)

TUM Chair:  
**[www.bgu.tum.de/en/lmf/vision/](http://www.bgu.tum.de/en/lmf/vision/)**  
our official  account:  
**[github.com/tum-lmf](https://github.com/tum-lmf)**

in cooperation with **GAFAG**


 follow **@marccoru** for updates.

 **[github.com/marccoru](https://github.com/marccoru)** for code.

and **[marccoru.github.io](https://marccoru.github.io)**

# Input feature importance analysis through gradient backpropagation

if we change the input  $\mathbf{X}$ ...


$$\frac{\partial \max(\hat{\mathbf{y}})}{\partial \mathbf{X}}$$

# Input feature importance analysis through gradient backpropagation

if we change the input  $X$ ...

$$\frac{\partial \max(\hat{y})}{\partial X}$$

... how would the highest predicted score  $\max(\hat{y})$  change?



Can be implemented in four lines

```
In [18]: x_ = torch.autograd.Variable(x[None,:,:], requires_grad=True)
logprobabilities = model.forward(x_)
logprobabilities.exp().max().backward()
dydx = x_.grad
```

# Gradients from $\max(\hat{y})$ to $X$

—●— B02 (blue)    —●— B03 (green)    —●— B04 (red)    —●— B05  
—●— B06    —●— B07    —●— B08    —●— B8A  
—●— B11    —●— B12

