



PROGRAMME OF THE  
EUROPEAN UNION



co-funded with



# CMIX-II: Cloud Mask Intercomparison eXercise – second edition

*Jan Wevers<sup>a</sup>, Sergii Skakun<sup>b,c</sup>, Carsten Brockmann<sup>a</sup>, Eric Vermote<sup>c</sup>*

<sup>a</sup> Brockmann Consult GmbH, 21029 Hamburg, Germany

<sup>b</sup> Department of Geographical Science, University of Maryland, College Park, MD 20742, USA

<sup>c</sup> NASA Goddard Space Flight Center Code 619, Greenbelt, MD 20771, USA

7<sup>th</sup> Sentinel-2 Validation Team Meeting

13 – 15 October 2025 | ESA – ESRIN | Frascati (RM), Italy

# Objective

- The **Cloud Masking Inter-comparison eXercise** (CMIX) is an international collaborative effort aimed at intercomparing cloud detection algorithms for moderate-spatial resolution (10-30 m) spaceborne optical sensors.
- CMIX is a joint activity by **ESA** and **NASA** in the frame of the **CEOS working group on calibration and validation** (WGCV)
- The first CMIX was started in 2018, the final meeting was held in Dec. 2019. Internal report was circulated between 2020 and 2021 leading to a final publication of the results in 2022 (Skakun et al. 2022).
- Report and paper provided recommendations and ways forward to advance inter-comparison and validation of cloud detection algorithms. Those were related to providing a quantitative definition of clouds, generating new cloud reference/validation datasets, and expanding the analysis framework.
  - These necessitated carrying out the second CMIX (**CMIX-II**).
- Like CMIX, **CMIX-II** utilizes open and free **multi-spectral** data
  - **Sentinel-2** and
  - **Landsat 8/9**
- The Focus of CMIX-II is not only on clouds but also on **cloud shadow** detection.

# Design – course of events

- Definition of the inter-comparison protocol and reference datasets
  - June 2022
- **Phase 1 – Test Dataset (TDS) exercise**
  - Provision of reference datasets samples
  - Application of the CM processors
  - Initial analysis of the results and provision to participants
  - Feedback from participants
- **Phase 2 – Main exercise**
  - Preparation of complete datasets
  - Provision of datasets to the participants
  - Application of the CM processors
  - Analysis of the results and preparation of report
- Provision of Analysis Report (Internally): Early 2026
- 2nd Workshop of CEOS-WGCV CMIX II
  - Short after release of report @ ESA or NASA
- Preparation of a publication: Mid 2026

CMIX-II was split into two phases to make sure all participating parties agree on the design of the datasets.

# Participants

- 25 participating algorithms from 20 participating parties  
→ **largest intercomparison of cloud masking algorithms**

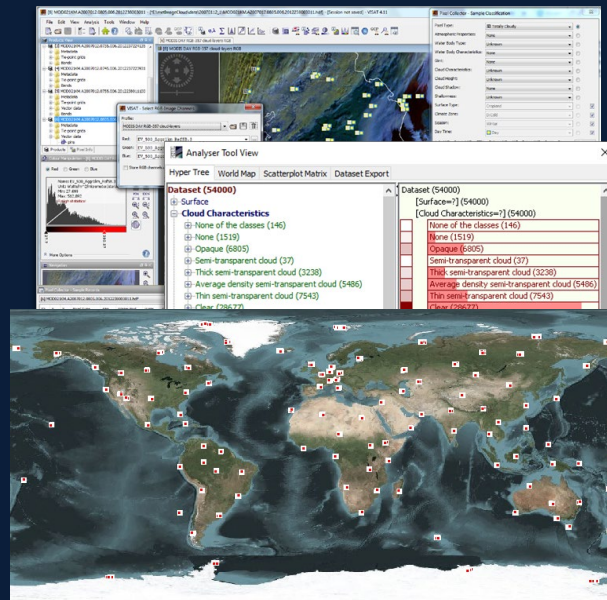
Participant	Processor/Model	Affiliation
Sergii Skakun Eric Vermote Jean-Claude Roger	LaSRC	NASA Goddard Space Flight Center / University of Maryland
Jan Wevers Jorrit Scholze	IdePix	Brockmann Consult GmbH
Alistair Francis	SEnSel v2	ESA/ESRIN
Feng Yin	SIAC	University College London
David Frantz	FORCE	Trier University
Béatrice Berthelot	MAGELLIUM	MAGCMA
Luis Gómez-Chova	UVDeepCloud	University of Valencia
Kaupo Voormansik Tetiana Shtym	KappaMask	KappaZeta Ltd.
Jerome Louis	Senzcor & prototype	Telespazio France
Christopher Brown Valerie Pasquarella	CloudScore+ & CloudScore CDF	Google, LLC
Hervé Poilvé	Overland & OneCloudDetector	Airbus Geo-Intelligence Toulouse (Airbus DS)

Participant	Processor/Model	Affiliation
Bringfried Pflug, Avi Avi Pertiwi Raquel de los Reyes	PACO	DLR
Hankui Zhang	LANA	South Dakota State University
Zhe Zhu Shi Qiu	Fmask 4.7 & Fmask 5	University of Connecticut
Ute Gangkofner	noClouds	eoConsultancy
Christian Köppl Andreas Brunn Hannah Kofler	ConstellR	ConstellR
Matthieu Molinier Jean-Eudes Gbodjo	MoCo & DeepCluster	VTT
Ruben van de Kerchove Carolien Toté Yannis Kalfas	LOS	Vito
Jakub Nalepa Bartosz Grabowski	KPLabs_nnUNet	KP Labs
Chris Rampersad Rick Chern		EarthDaily

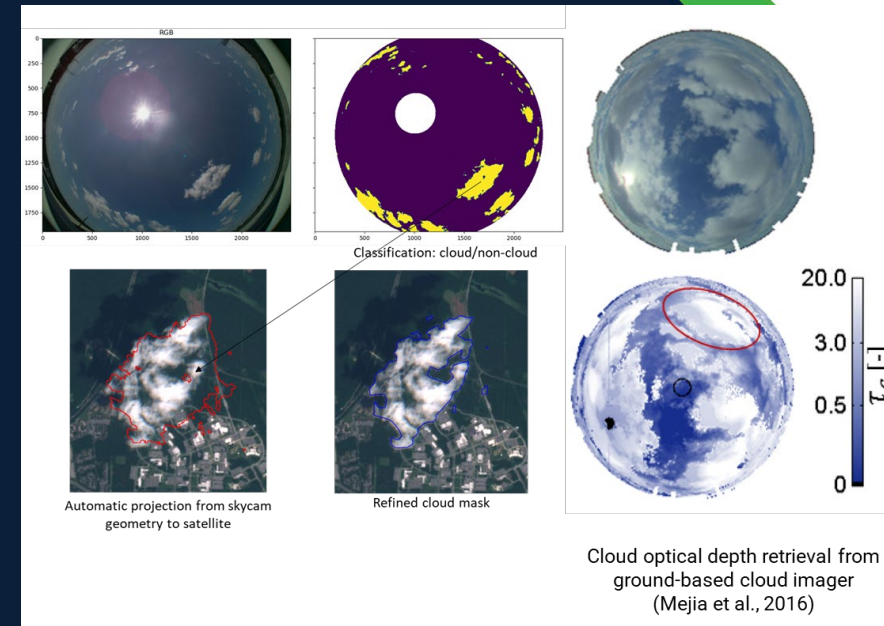
# Datasets

- Sky Camera Network
- PixBox (expert pixel collection) dataset (incl. COD)
- Multi-temporal (time series) - identifying potential systematic errors
- Collaborative dataset using IRIS (active learning): classification of subsets by the participants

## PixBox

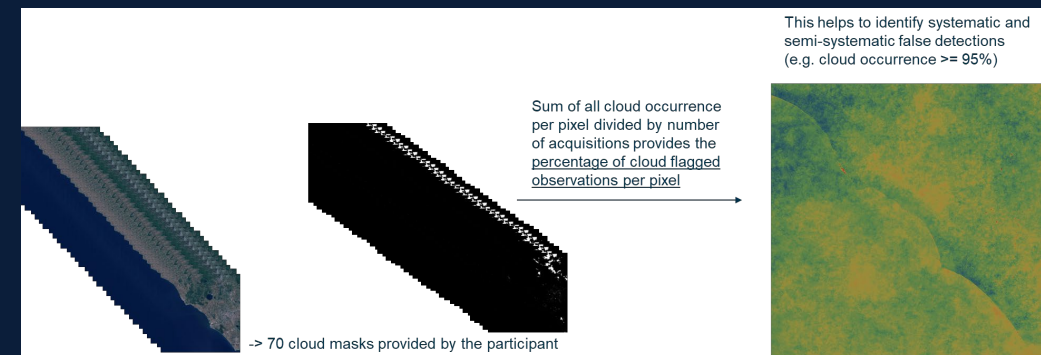
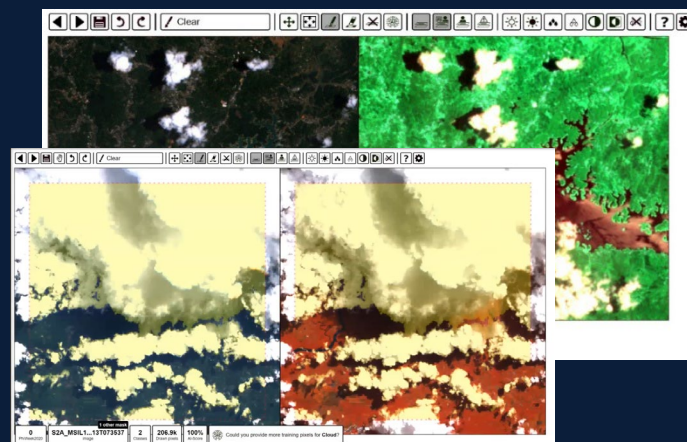
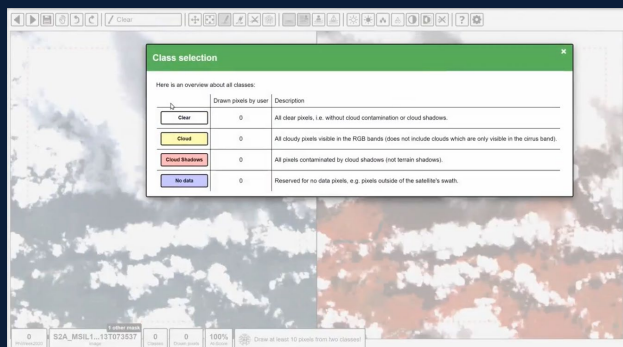


## Sky Camera



## Multi-temporal

### Collaborative dataset



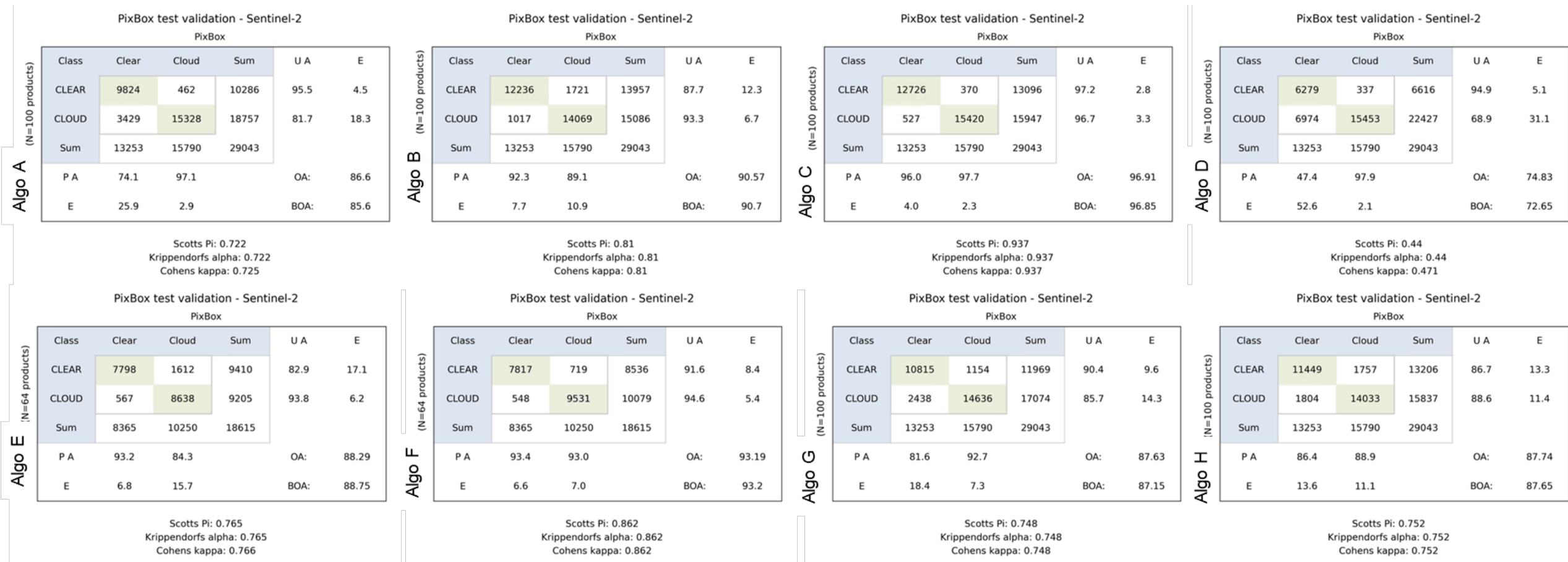
- Per pixel: confusion matrix and estimated OA, PA, UA
- Visual inspection (potentially study an impact on SR, especially with transparent/cirrus clouds)
- Comparisons to sky imagery outputs:
  - Per-pixel comparison, when a precise co-registration between sky and satellite images will be feasible.
  - Cloud fraction: cloud fractions derived from sky and satellite imagery in a specified ROI will be estimated and will be compared. A separate analysis of completely clear and fully cloudy scenes will be performed. Analysis will be performed using a varying threshold for defining clouds (fuzzy logic analysis).
- Time series analysis:
  - Temporal analysis on cloud (mask) occurrence per pixel per algorithm, to identify potential systematic and semi-systematic errors.
  - Statistic analysis between number of systematic errors and confusion matrix performance.

# Results – PixBox dataset

- CMIX-II is actively running in Phase 2. Therefore, only anonymized results will be shown.
- On the next slides the performances of the algorithms for Cloud/Clear and Cloud/Clear/Cloud shadow will be shown.
  - 4 algorithms do not provide a cloud shadow flag
- Questions to be answered:
  - Machine learning is superior to classic spectral test based approaches?
  - High accuracies of cloud detection automatically leads to good cloud shadow detections?

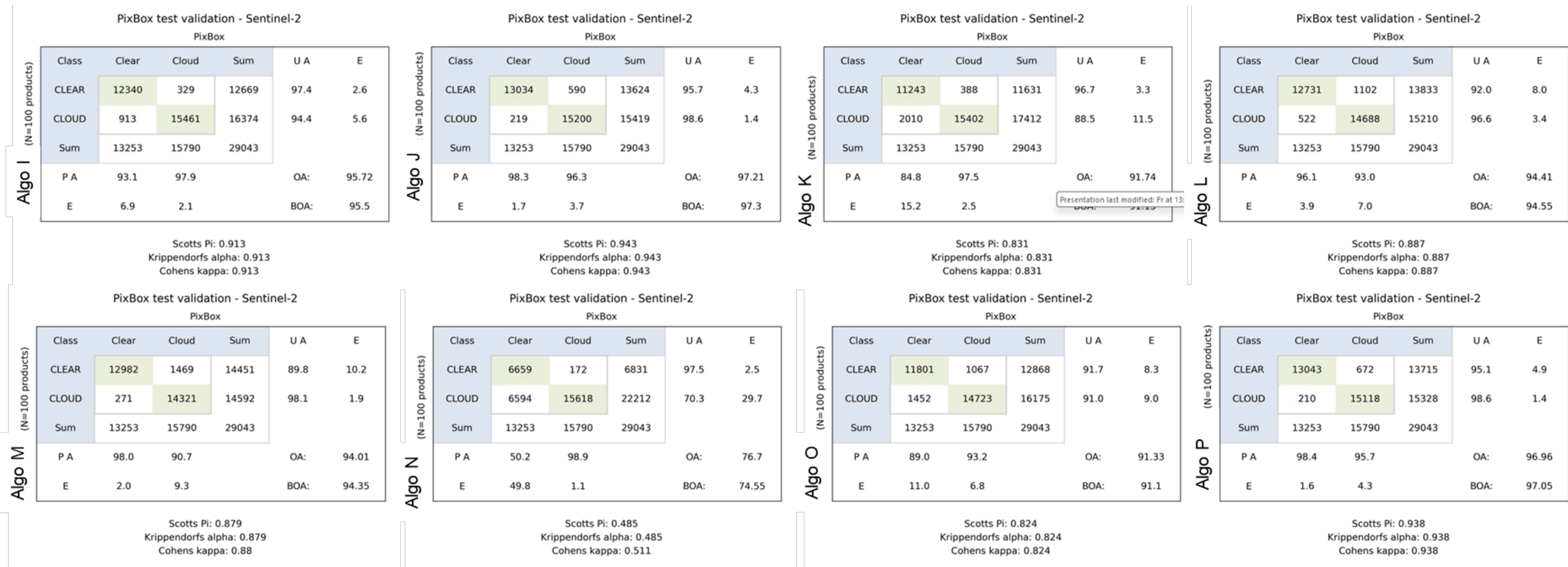
# Results – Cloud vs. Clear

CMIX-II is actively running in Phase 2. Therefore, only anonymized results will be shown.



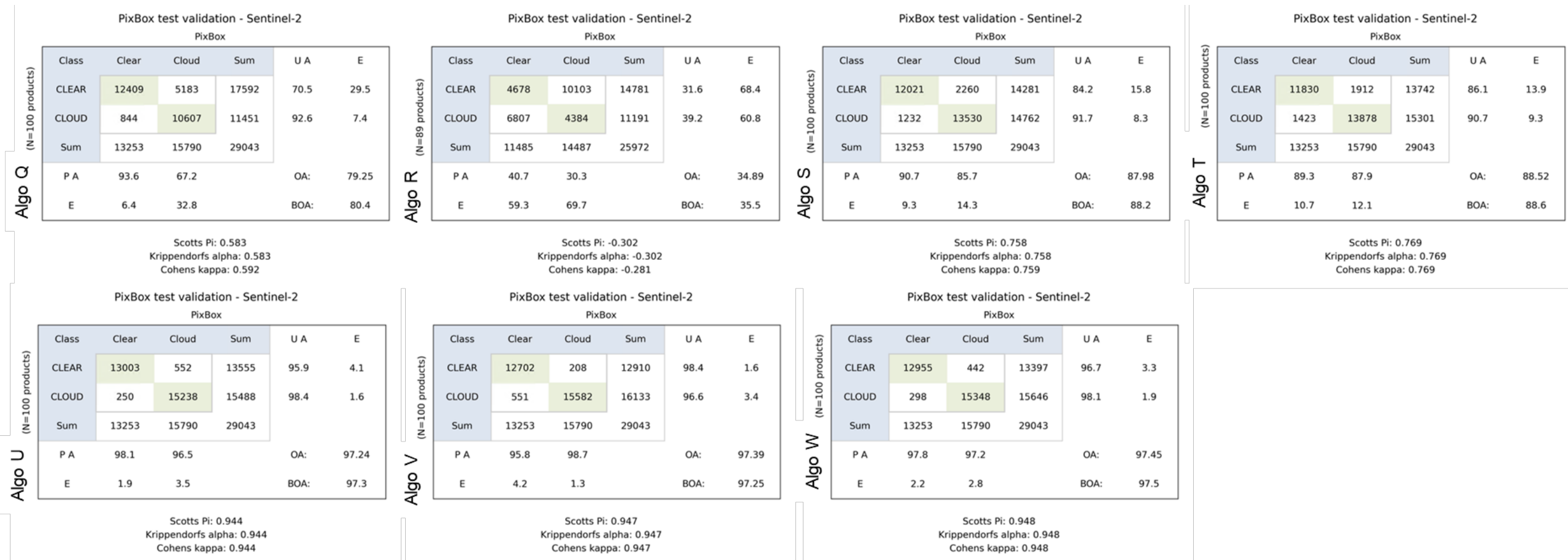
# Results – Cloud vs. Clear

CMIX-II is actively running in Phase 2. Therefore, only anonymized results will be shown.



# Results – Cloud vs. Clear

CMIX-II is actively running in Phase 2. Therefore, only anonymized results will be shown.



# Results – Cloud vs. Clear

- Does Machine Learning (ML) outperform classical approaches?

- Yes
- But: ML method as such is not a guarantee for good performance -> Strongly depending on training data and know-how

Anonymys_Name	PixBox BOA	ML Algorithm	ML type	Use of time series
Algo A	85.6	Yes	Deep learning model	Yes
Algo B	90.7	Yes	Deep learning model	Yes
<b>Algo C</b>	<b>96.85</b>	<b>Yes</b>	<b>Deep convolution (U-Net)</b>	<b>No</b>
Algo D	72.65	Yes	Self-supervised learning algorithms (DL)	No
Algo E	88.75	No	-	No
Algo F	93.2	Yes	Physics-Informed Machine Learning (PIML)	
Algo G	87.15	No	-	No
Algo H	87.65	No	-	No
<b>Algo I</b>	<b>95.5</b>	<b>Yes</b>	<b>Deep learning segmentation methods</b>	<b>No</b>
<b>Algo J</b>	<b>97.3</b>	<b>Yes</b>	<b>U-Net</b>	<b>No</b>
Algo K	91.15	Yes	U-Net	No
Algo L	94.55	Yes	EvoNet	?
Algo M	94.35	Yes	U-Net	No
Algo N	74.55	Yes	Self-supervised learning algorithms (DL)	No
Algo O	91.1	No	-	No
<b>Algo P</b>	<b>97.05</b>	<b>Yes</b>	<b>U-Net CNN with a regnety006 encoder</b>	<b>No</b>
Algo Q	80.4	?	?	No
Algo R	35.5	No	-	No
Algo S	88.2	No	-	No
Algo T	88.6	?	-	No
<b>Algo U</b>	<b>97.3</b>	<b>Yes</b>	<b>SegFormer</b>	<b>No</b>
<b>Algo V</b>	<b>97.25</b>	<b>Yes</b>	<b>Deep learning for semantic segmentation</b>	<b>No</b>
<b>Algo W</b>	<b>97.5</b>	<b>Yes</b>	<b>U-Net CNN with a regnety006 encoder</b>	<b>No</b>

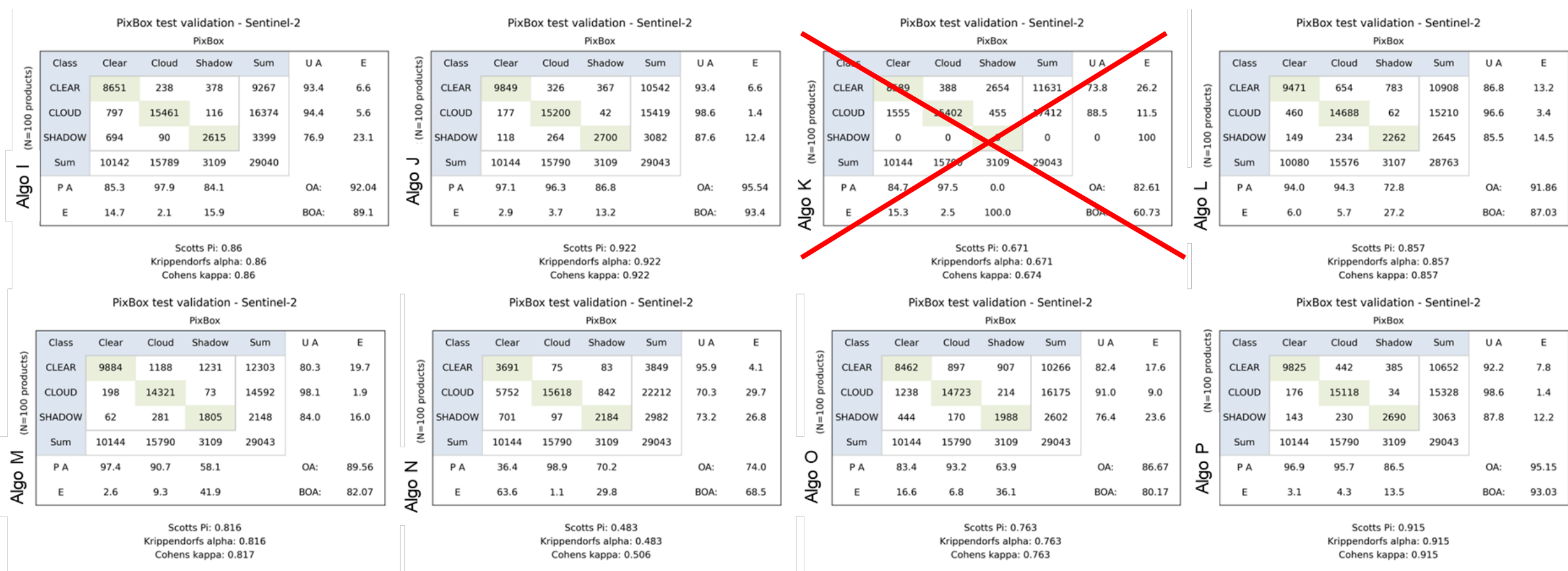
# Results – Cloud vs. Clear vs. Cloud Shadow

CMIX-II is actively running in Phase 2. Therefore, only anonymized results will be shown.



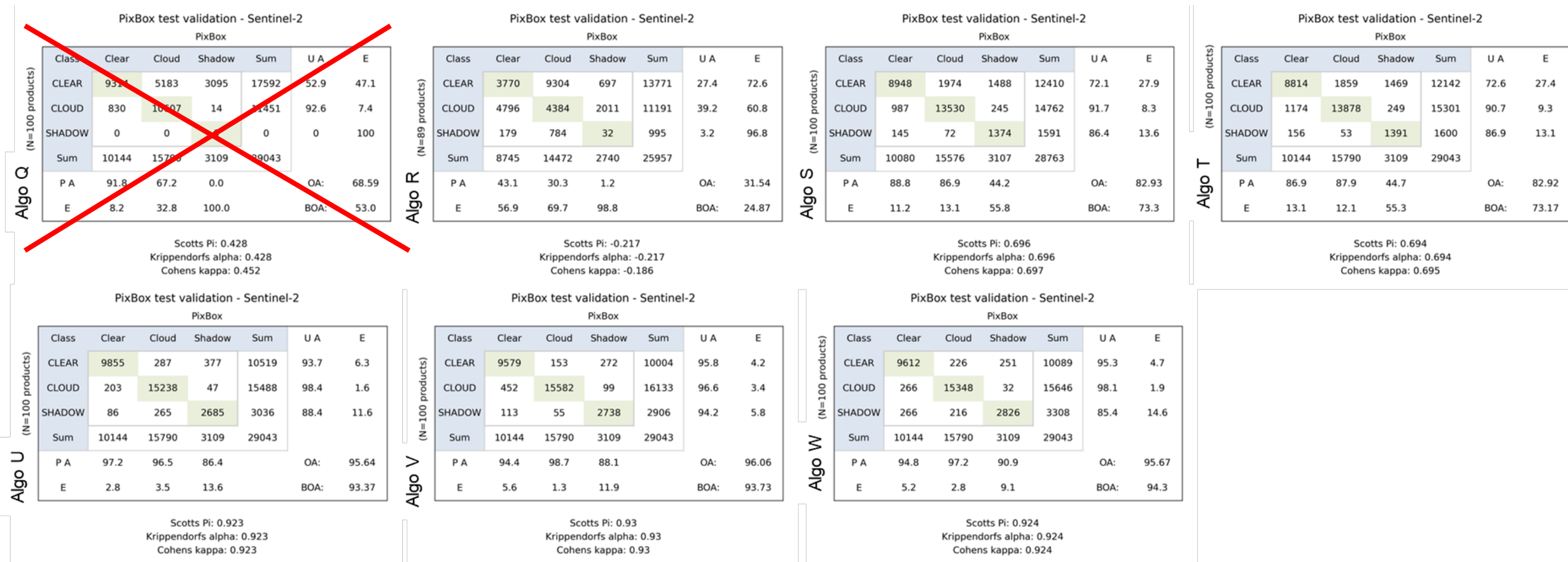
# Results – Cloud vs. Clear vs. Cloud Shadow

CMIX-II is actively running in Phase 2. Therefore, only anonymized results will be shown.



# Results – Cloud vs. Clear vs. Cloud Shadow

CMIX-II is actively running in Phase 2. Therefore, only anonymized results will be shown.



# Results – Cloud vs. Clear vs. Cloud Shadow

- Does an overall good cloud detection (BOA), or a high Producer Accuracy (PA) of cloud, guarantee a good cloud shadow detection?
  - No!
  - But it is the essential first step, since most cloud shadow detection algorithms rely on geometric approaches, using the cloud mask as a reference.

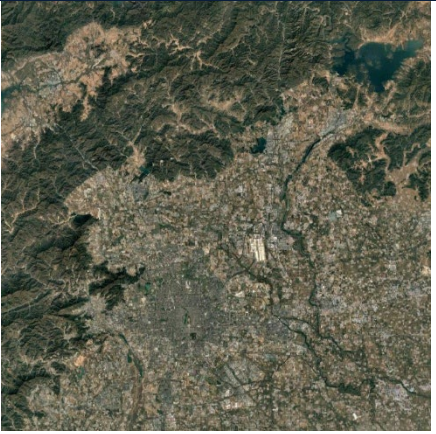
BOA for Cloud/Clear ↓ ↓ PA Cloud for Cloud/Clear/Shadow

Anonymys_Name	PixBox BOA	PA Cloud	PA Cloud Shadow
Algo A	85.6	0	0
Algo B	90.7	0	0
Algo C	96.85	97.7	71.4
Algo D	72.65	97.9	63.9
Algo E	88.75	84.3	49.4
Algo F	93.2	93	66.9
Algo G	87.15	93.2	45.8
Algo H	87.65	88.9	66.1
<b>Algo I</b>	<b>95.5</b>	<b>97.9</b>	<b>84.1</b>
<b>Algo J</b>	<b>97.3</b>	<b>96.3</b>	<b>86.8</b>
Algo K	91.15	0	0
Algo L	94.55	94.3	72.8
Algo M	94.35	90.7	58.1
Algo N	74.55	98.9	70.2
Algo O	91.1	93.2	63.9
<b>Algo P</b>	<b>97.05</b>	<b>95.7</b>	<b>86.5</b>
Algo Q	80.4	0	0
Algo R	35.5	30.3	1.2
Algo S	88.2	86.9	44.2
Algo T	88.6	87.9	44.7
<b>Algo U</b>	<b>97.3</b>	<b>96.5</b>	<b>86.4</b>
<b>Algo V</b>	<b>97.25</b>	<b>98.7</b>	<b>88.1</b>
<b>Algo W</b>	<b>97.5</b>	<b>97.2</b>	<b>90.9</b>

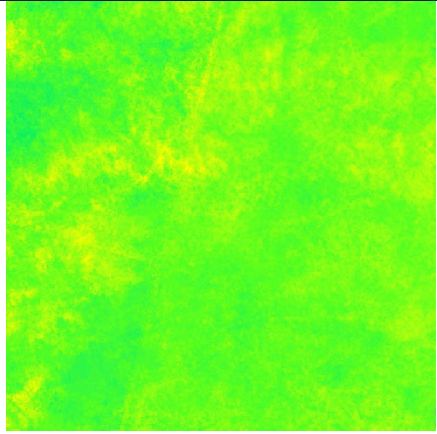
# Results – Multi-temporal dataset

## Multi-temporal analysis:

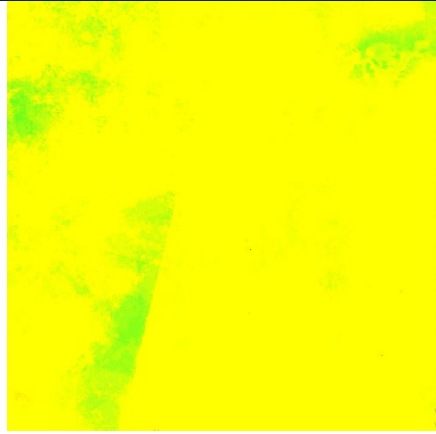
- A pixel-based analysis (like PixBox) does not give any insight into the spatial behavior of the algorithms, and it is not capable of identifying systematic and semi-systematic false detections.
- By processing an annual time series over a given location and calculating the percentage of the cloud flagged observations in relation to all observations, those systematic issues can be identified.
- 13 participants provided the results for the multi-temporal dataset
  - Again, only anonymized results will be shown.
  - The BOA for Cloud/Clear from the PixBox dataset is noted next to the Algorithm ID
- The main question: Does a good performance for the PixBox dataset mean the absence of systematic errors?



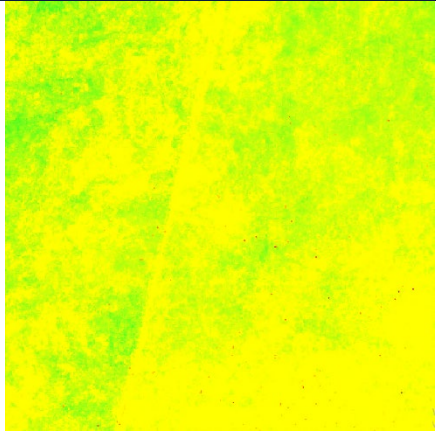
Algo A – BOA 85.6



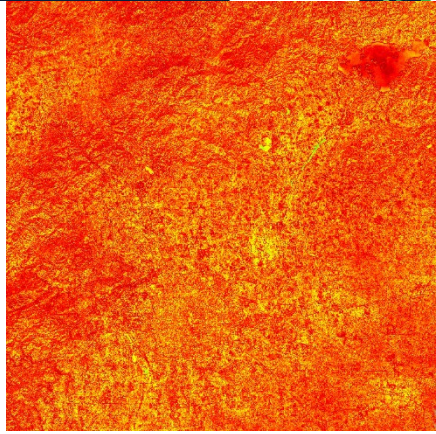
Algo B – BOA 90.7



Algo C – 96.85



Algo D – 72.65

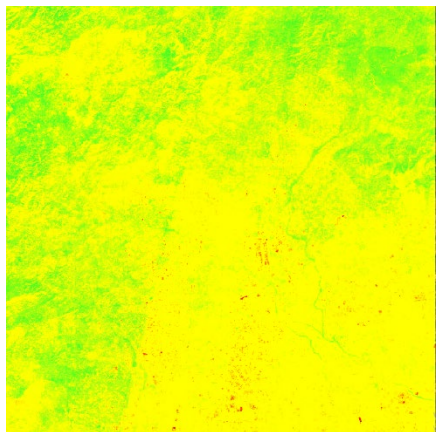


# China, Beijing

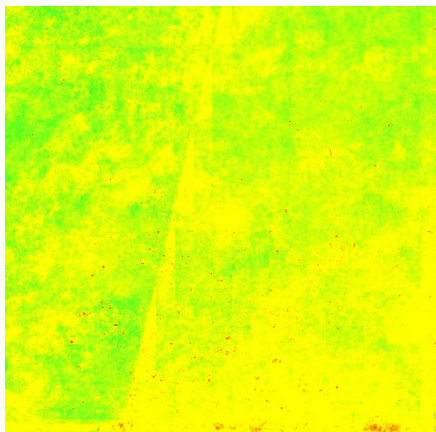
Percentage of cloud masked pixels



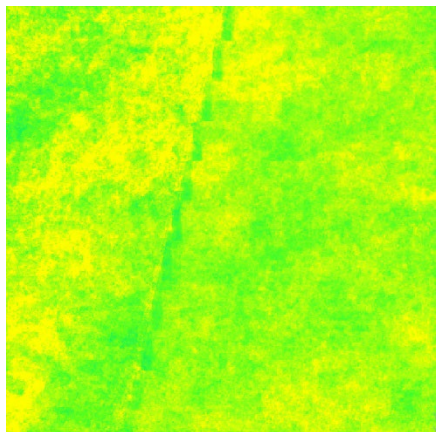
Algo H – BOA 87.65



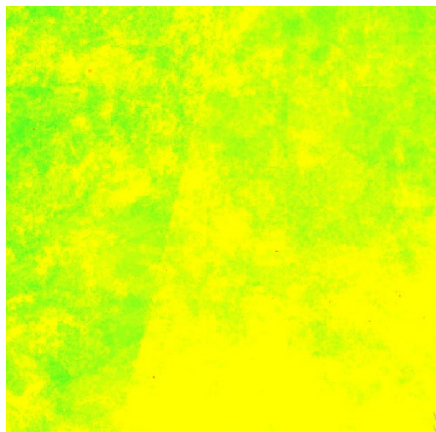
Algo I – BOA 95.5



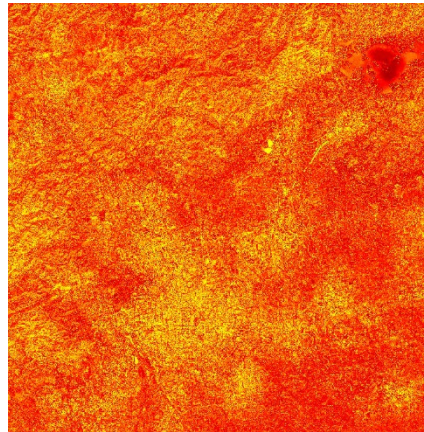
Algo J – BOA 97.3



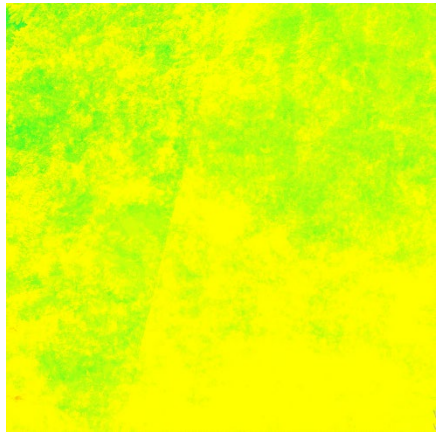
Algo K – BOA 91.15



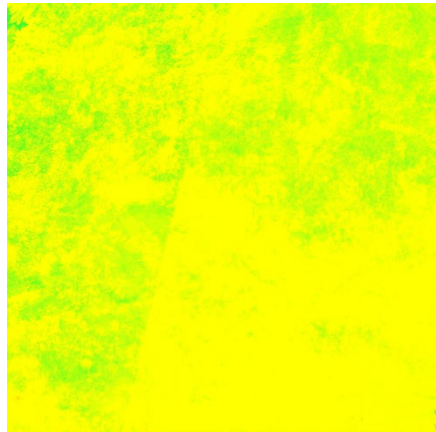
Algo N – BOA 74.55



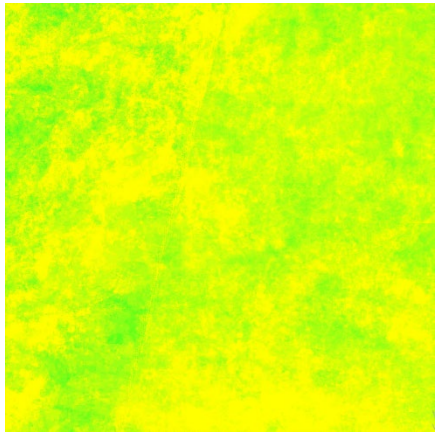
Algo S – BOA 88.2



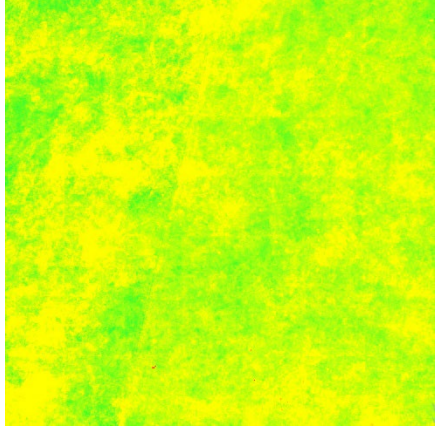
Algo T – BOA 88.6



Algo V – BOA 97.25



Algo W – BOA 97.5





# Conclusion

- The exercise is actively running with strong engagement from the participant community.
- Its outcomes are of great interest, both to the participants and to the wider cloud detection community, as highlighted by feedback at the Living Planet Symposium.
- The activity delivers valuable insights into the performance of diverse algorithms and methodological approaches.
- It also represents the largest intercomparison of cloud detection algorithms carried out to date.



BROCKMANN  
CONSULT

**Thank you for the  
attention!**



BROCKMANN  
CONSULT

## Turning Environmental Data Into Knowledge

[brockmann-consult.de](http://brockmann-consult.de)

[info@brockmann-consult.de](mailto:info@brockmann-consult.de)

+49 40 696 389 300

in 