



HAL
open science

Cloud Mask Intercomparison eXercise (CMIX): An evaluation of cloud masking algorithms for Landsat 8 and Sentinel-2

Sergii Skakun, Jan Swevers, Carsten Brockmann, Georgia Doxani, Matej Aleksandrov, Matej Batič, David Frantz, Ferran Gascon, Luis Gómez-Chova, Olivier Hagolle, et al.

► To cite this version:

Sergii Skakun, Jan Swevers, Carsten Brockmann, Georgia Doxani, Matej Aleksandrov, et al.. Cloud Mask Intercomparison eXercise (CMIX): An evaluation of cloud masking algorithms for Landsat 8 and Sentinel-2. *Remote Sensing of Environment*, 2022, 274, pp.112990. 10.1016/j.rse.2022.112990 . hal-04515790v2

HAL Id: hal-04515790

<https://uca.hal.science/hal-04515790v2>

Submitted on 11 Oct 2024

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

1 **Cloud Mask Intercomparison eXercise (CMIX): an evaluation of cloud masking**
2 **algorithms for Landsat 8 and Sentinel-2**

3 *Sergii Skakun*^{a,b}, *Jan Wevers*^c, *Carsten Brockmann*^c, *Georgia Doxani*^d,
4 *Matej Aleksandrov*^e, *Matej Batič*^e, *David Frantz*^{f,o}, *Ferran Gascon*^g, *Luis Gómez-Chova*^h,
5 *Olivier Hagolle*ⁱ, *Dan López-Puigdollers*^h, *Jérôme Louis*^j, *Matic Lubej*^e, *Gonzalo Mateo-*
6 *García*^h, *Julien Osman*^k, *Devis Peressutti*^e, *Bringfried Pflug*^l, *Jernej Puc*^e, *Rudolf Richter*
7 ^m, *Jean-Claude Roger*^{a,b}, *Pat Scaramuzza*ⁿ, *Eric Vermote*^b, *Nejc Vesel*^e, *Anže Zupanc*^e,
8 *Lojze Žust*^e

9
10 ^a Department of Geographical Science, University of Maryland, College Park, MD 20742,
11 USA

12 ^b NASA Goddard Space Flight Center Code 619, Greenbelt, MD 20771, USA

13 ^c Brockmann Consult GmbH, 21029 Hamburg, Germany

14 ^d SERCO SpA c/o European Space Agency ESA-ESRIN, 00044 Frascati, Italy

15 ^e Sinergise LTD, 1000 Ljubljana, Slovenia

16 ^f Geography Department, Humboldt-Universität zu Berlin, 10099 Berlin, Germany

17 ^g European Space Agency ESA-ESRIN, 00044 Frascati, Italy

18 ^h Image Processing Laboratory, University of Valencia, 46980 Valencia, Spain

19 ⁱ Centre d'études Spatiales de la Biosphère, CESBIO Unite mixte Université de Toulouse-
20 CNES-CNRS-IRD, 31401 Toulouse CEDEX 9, France

21 ^j Telespazio France, 31023 Toulouse, France

22 ^k Thales Services SAS, Labège, France

23 ^l DLR, German Aerospace Center, D-12489 Berlin, Germany

24 ^m DLR, German Aerospace Center, D-82234 Wessling, Germany, rudolf.richter@dlr.de

25 ⁿ KBR, contractor to the contractor to the U.S. Geological Survey (USGS) Earth Resources
26 Observation and Science Center (EROS), Sioux Falls, SD 57198, USA

27 ^o Earth Observation and Climate Processes, Trier University, 54286 Trier, Germany

28

29

30 Corresponding author:

31 Sergii Skakun

32 Department of Geographical Sciences, University of Maryland, College Park

33 1153 LeFrak Hall, College Park, MD, USA, 20742

34 skakun@umd.edu

35 **Highlights**

- 36 • Ten cloud masking algorithms for Landsat 8 and Sentinel-2 are evaluated
- 37 • Algorithm performance varied depending on the reference dataset
- 38 • Average overall accuracy for Sentinel-2 was $80.0\pm 5.3\%$ to $89.4\pm 2.4\%$
- 39 • Average overall accuracy for Landsat 8 was $79.8\pm 7.1\%$ to $97.6\pm 0.8\%$
- 40 • Performance of algorithms improved when thin/semi-transparent clouds not
- 41 considered

42

43 **Abstract.**

44 Cloud cover is a major limiting factor in exploiting time-series data acquired by optical
45 spaceborne remote sensing sensors. Multiple methods have been developed to address the
46 problem of cloud detection in satellite imagery and a number of cloud masking algorithms
47 have been developed for optical sensors but very few studies have carried out quantitative
48 intercomparison of state-of-the-art methods in this domain. This paper summarizes results of
49 the first Cloud Masking Intercomparison eXercise (CMIX) conducted within the Committee
50 Earth Observation Satellites (CEOS) Working Group on Calibration & Validation (WGCV).
51 CEOS is the forum for space agency coordination and cooperation on Earth observations,
52 with activities organized under working groups. CMIX, as one such activity, is an
53 international collaborative effort aimed at intercomparing cloud detection algorithms for
54 moderate-spatial resolution (10-30 m) spaceborne optical sensors. The focus of CMIX is on
55 open and free imagery acquired by the Landsat 8 (NASA/USGS) and Sentinel-2 (ESA)
56 missions. Ten algorithms developed by nine teams from fourteen different organizations
57 representing universities, research centers and industry, as well as space agencies (CNES,
58 ESA, DLR, and NASA), are evaluated within the CMIX. Those algorithms vary in their
59 approach and concepts utilized which were based on various spectral properties, spatial and
60 temporal features, as well as machine learning methods. Algorithm outputs are evaluated
61 against existing reference cloud mask datasets. Those datasets vary in sampling methods,
62 geographical distribution, sample unit (points, polygons, full image labels), and generation
63 approaches (experts, machine learning, sky images). Overall, the performance of algorithms
64 varied depending on the reference dataset, which can be attributed to differences the
65 reference datasets were produced. The algorithms were in good agreement for thick cloud
66 detection, which were opaque and had lower uncertainties in their identification, in contrast
67 to thin/semi-transparent clouds detection. Not only did CMIX allow identification of

68 strengths and weaknesses of existing algorithms and potential areas of improvements, but
69 also the problems associated with the existing reference datasets. The paper concludes with
70 recommendations on generating new reference datasets, metrics, and an analysis framework
71 to be further exploited and additional input datasets to be considered by future CMIX
72 activities.

73

74 Keywords: cloud, intercomparison, validation, Landsat 8, Sentinel-2, CMIX, CEOS

75

76 **1 Introduction**

77 Identification of clouds in satellite imagery acquired by passive remote sensing
78 sensors in the visible and infrared parts of the electromagnetic spectrum (EM) is an essential
79 pre-processing step in producing high-quality geoinformation products. Omission of clouds
80 can lead to errors that propagate to high-level products related to Earth surface monitoring,
81 whereas over detection of clouds can lead to a reduced number of valid observations and,
82 therefore, decrease the frequency of cloud-free data. Development of cloud masking
83 algorithms remains an area of active research in the remote sensing community (Foga et al.,
84 2017; Frantz et al., 2018; Hagolle et al., 2010; Hollingsworth et al., 1996; Irish et al., 2006;
85 López-Puigdollers et al., 2021; Qiu et al., 2019; Scaramuzza et al., 2012; Zhu et al., 2015;
86 Zhu and Woodcock, 2012). A range of algorithms utilize satellite image spectral and spatial
87 properties along with decision tree rules to distinguish cloud versus non-cloud regions (Qiu et
88 al., 2019). These algorithms rely mainly on physical properties of cloud reflectance.
89 Utilization of multi-temporal satellite images, where clouds are considered “anomalies” with
90 respect to a cloud-free reference, can generally improve cloud detection (Frantz et al., 2015;
91 Hagolle et al., 2010; Zhu & Woodcock, 2014). With the advancement of machine learning
92 (ML) and deep learning (DL) methods neural networks models are trained to detect clouds in
93 satellite imagery (Chai et al., 2019; Jeppesen et al., 2019; Mateo-García et al., 2020; Segal-
94 Rozenhaimer et al., 2020; Wieland et al., 2019; Xie et al., 2017).

95 Although a large number of cloud masking algorithms for optical satellite imagery is
96 currently available, there are a limited quantity of studies aiming at their intercomparison.
97 Three studies should be mentioned in this regard. Foga et al. (2017) compared 13 cloud
98 masking algorithms and their variants for cloud detection in Landsat 7 and Landsat 8 data.
99 Their primary objective was to select an algorithm for generating quality assurance (QA)
100 layers when producing operational Landsat data products. They found that CFMask, a C code

101 version of the Fmask algorithm (Qiu et al., 2019; Zhu et al., 2015), gave the best
102 performance, and this algorithm is currently used within the U.S. Geological Survey (USGS)
103 operational processing chain to generate Landsat Level-1 products (Wulder et al., 2019).
104 Baetens et al. (2018) compared three methods applied to Sentinel-2 data by analyzing 30
105 images and found large differences in quality, specifically when taking into account the
106 necessary dilation (buffer) of cloud masks. Tarrio et al. (2020) carried out a study comparing
107 five cloud masking algorithms for Sentinel-2 imagery. By analyzing 28 images over six
108 Sentinel-2 tiles using a sample-based approach and analyst-interpreted reference data they
109 found that none of the algorithms yielded the best performance in terms of identifying both
110 cloud and shadow. They also explored ensemble models to integrate outputs from multiple
111 algorithms and found that on average a +2.7% gain can be achieved over the best-performing
112 model, although at the expense of computational performance.

113 The main objective of this paper is to summarize results of the first Cloud Masking
114 Intercomparison eXercise (CMIX) conducted within the Committee of Earth Observation
115 Satellites (CEOS) Working Group on Calibration & Validation (WGCV). CMIX is an
116 international collaborative effort co-led by National Aeronautics and Space Administration
117 (NASA) and European Space Agency (ESA) aimed at intercomparing state-of-the-art cloud
118 masking algorithms for moderate-spatial resolution (10-30 m) spaceborne optical sensors.
119 CMIX was recommended following the first Atmospheric Correction Inter-comparison
120 eXercise (ACIX) (Doxani et al., 2018), and was conducted in conjunction with ACIX-II-
121 Land and ACIX-II-Aqua (Pahlevan et al., 2021). The focus of this effort is on open and free
122 imagery acquired by Landsat 8 Operational Land Imager (OLI) and Thermal Infrared Sensor
123 (TIRS), and Sentinel-2 MultiSpectral Instrument (MSI) sensors, with corresponding cloud
124 masking algorithms applied. Five existing cloud reference datasets for Landsat 8 and
125 Sentinel-2 are utilized to compare ten cloud masking algorithms. Within CMIX, a qualitative

126 definition of “cloud” is adopted, which provides an absolute (spectrally independent)
127 indication of cloudiness in the satellite image. Although rules defining clouds vary across
128 algorithms and reference data, ultimately all data are converted to “cloud” and “non-cloud”
129 classes to perform a consistent intercomparison. Algorithms are compared using the same set
130 of reference data and metrics under identical conditions. Cloud shadows are not considered in
131 this study, since it is typically a cloud-derived product, and its performance heavily depends
132 on accuracy of cloud detection. Consequently, efforts are primarily directed to cloud mask
133 evaluation.

134 The rest of the paper is organized as follows: a brief description of cloud reference
135 data, cloud masking algorithms, and performance metrics is provided in Section 2. Detailed
136 description of results and their implications are respectively presented in Section 3 and
137 Section 4. Section 5 offers recommendations on further activities regarding generation of
138 cloud reference data and intercomparison of algorithms.

139

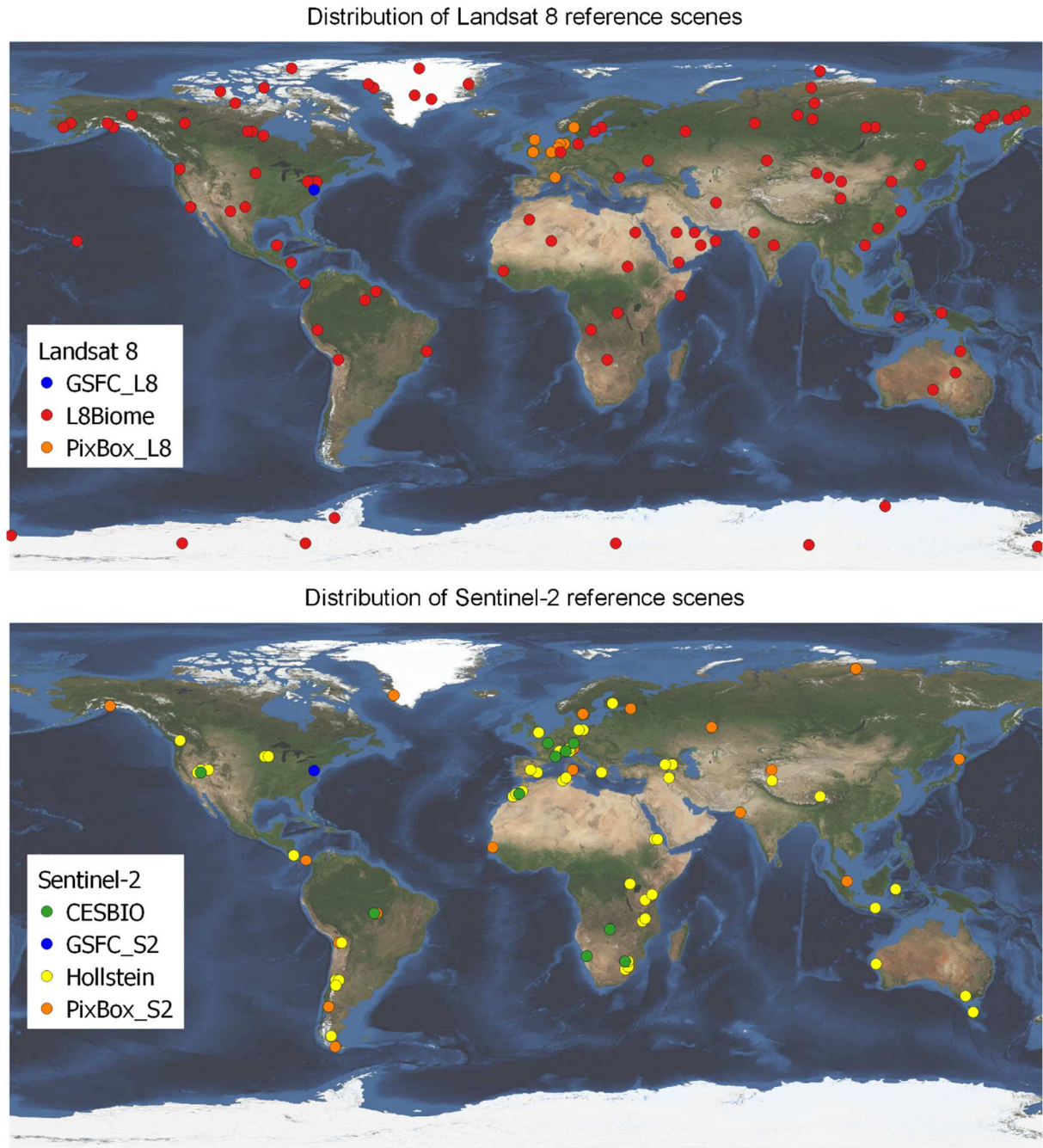
140 **2 Methods**

141 *2.1 Cloud reference datasets*

142 Intercomparison of algorithms within CMIX is performed using existing Sentinel-2
143 and Landsat 8 cloud reference datasets (Table 1), which include Hollstein (Hollstein et al.,
144 2016), PixBox (Paperin et al., 2021a, 2021b), L8Biome (Foga et al., 2017), CESBIO
145 (Baetens et al., 2019) and GSFC (Skakun et al., 2021). These datasets were
146 collected/generated for different purposes using different methodologies and cloud class
147 nomenclatures. Some of the datasets are single-pixel collections (where a minimum mapping
148 unit is a pixel), while others are the collections of connected pixel areas (polygons) or
149 correspond to whole images. For the majority of datasets, pixels were classified manually
150 through photointerpretation by an expert or a group of experts; in others, the labelling process
151 was semi-automatic with extensive manual checking during classification and post-

152 processing. Geographical distribution of Landsat 8 and Sentinel-2 scenes in the reference
153 datasets is shown in Figure 1.

154



155

156 Figure 1. Geographical distribution of the Landsat 8 and Sentinel-2 scenes in the reference
157 datasets used in CMIX.

158

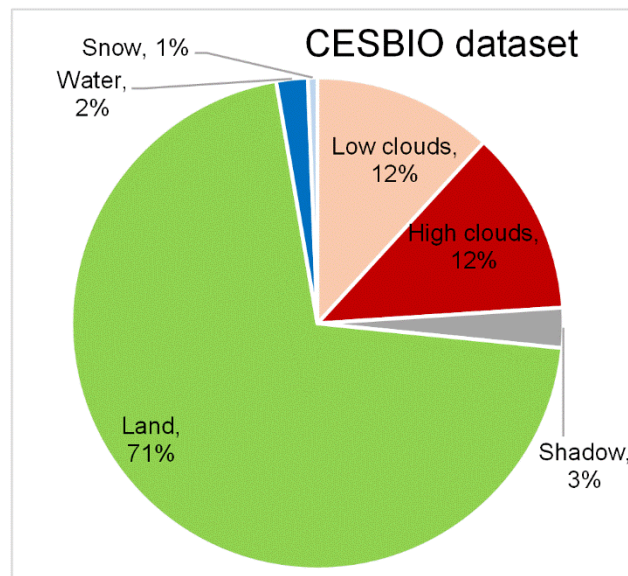
159 Table 1. Summary of cloud reference data (L8: Landsat 8, S2: Sentinel-2).

Dataset	Spatial domain	Level of automatization	Purpose	Thematic depth	Satellites	Spatial resolution	# scenes	Data Availability
CESBIO	Fully classified Sentinel-2 scenes	Classification using an iterative and supervised active learning method	Validation	6 classes	S2	60 m	S2: 30	https://zenodo.org/record/1460961
GSFC	Sample polygons	Manually selected and classified by an expert assisted by ground-based images of the sky	Validation	4 classes	L8, S2	Polygons (in vector format)	L8: 6 S2: 28	https://doi.org/10.17632/r7tnvx7d9g.1
Hollstein	Sample polygons	Manually selected and classified by an expert	Training and validation	6 classes	S2	Polygons (at 20 m)	S2: 59	https://git.gfz-potsdam.de/EnMAP/sentinel2_manual_classification_clouds
L8Biome	Fully classified Landsat 8 scenes	Manually classified by an expert	Training and validation	4 classes	L8	30 m	L8: 96	http://doi.org/10.5066/F7251GDH
PixBox	Sample pixels	Manually selected and classified by an expert	Validation	10 classes	S2, L8	S2: 10 m L8: 30 m	S2: 29 L8: 11	https://zenodo.org/record/5036991 https://zenodo.org/record/5040271

160

161 2.1.1 CESBIO dataset (Sentinel-2)

162 The CESBIO dataset was generated using an active learning method (Baetens et al.,
163 2019) using the Hollstein dataset (see section 2.1.3) as training samples. The classification
164 method was iterative, the operator constituted a first set of training samples, and iteratively
165 added other samples, where the classification results were wrong or uncertain. It provides
166 fully classified Sentinel-2 scenes into one of the following classes (Figure 2): low-altitude
167 clouds, high-altitude clouds, cloud shadows, land, water, and snow. In addition to the
168 classification map, a QA layer is provided showing the confidence of classification. Overall,
169 30 Sentinel-2 scenes were utilized in CMIX with the total number of labelled pixels
170 85,782,723 (at 60 m spatial resolution). The scenes were acquired from ten sites around the
171 world, five mainly vegetated and five arid sites. The detailed description of the CESBIO
172 dataset is given in Baetens et al. (2019).



173

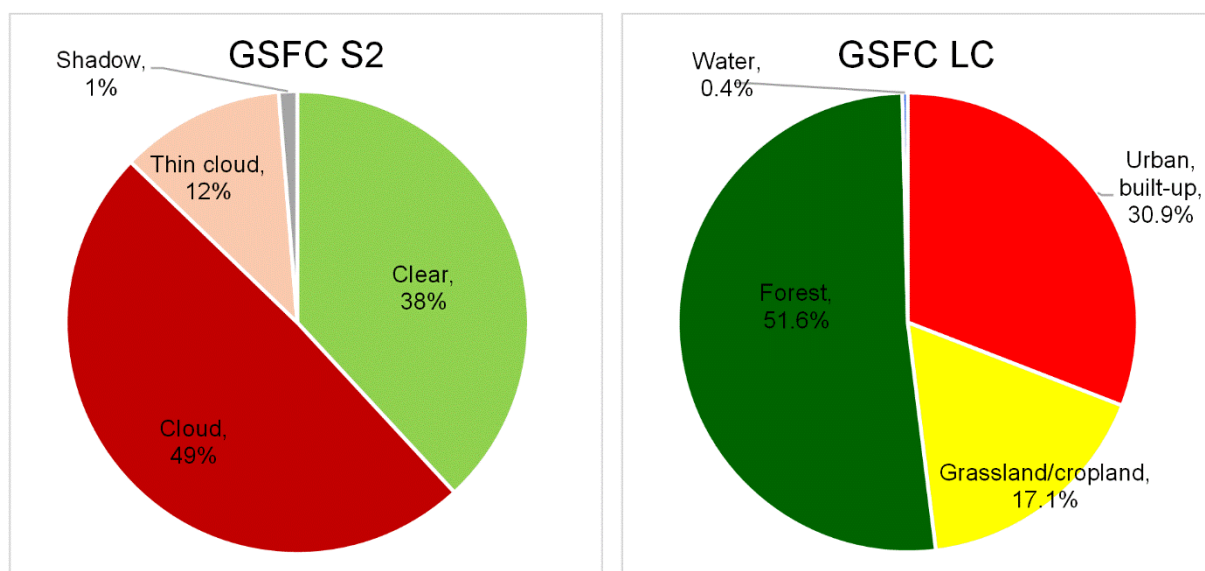
174 Figure 2. Distribution of labeled pixels in the CESBIO dataset.

175

176 2.1.2 GSFC dataset (Landsat 8, Sentinel-2)

177 GSFC cloud reference data were collected over the NASA Goddard Space Flight
178 Center (GSFC) (Skakun et al., 2021). The area is quite heterogeneous with major land cover

179 classes being forest (~52%) and impervious surfaces (31%) with patches of natural vegetation
 180 and cultivated areas (totaling 17%) (Figure 3). NASA GSFC also has an AERONET station
 181 (Holben et al., 1998), which provides aerosol optical thickness (AOT) and water vapor.
 182 Ground-based images of the sky were collected from 2017 through 2019 using a smartphone
 183 camera with a fisheye lens. These data were collected manually during the Landsat 8 and
 184 Sentinel-2 overpasses. Reference data were collected for 6 Landsat 8 and 28 Sentinel-2
 185 scenes. The objective was to capture various cloud conditions and seasonal variability.
 186 Labeling of satellite imagery was performed into cloud, thin cloud (semi-transparent),
 187 shadows, and clear classes (Figure 3). Regions within cloud boundaries were excluded from
 188 the reference data due to large uncertainties regarding the exact boundaries of clouds,
 189 especially on Sentinel-2 imagery (Skakun et al., 2021). In order to facilitate the labelling
 190 process, Sentinel-2 and Landsat 8 images were presented in various spectral combinations
 191 including true color (red-green-blue) and false color (NIR-red-green, SWIR1-NIR-red), and
 192 using a cirrus band (at 1.38 μm). The detailed description of the GSFC dataset is given in
 193 Skakun et al. (2021).

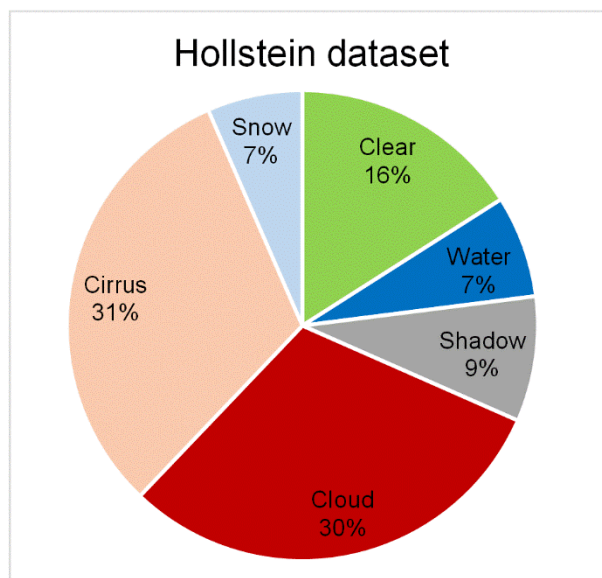


194
 195 Figure 3. Distribution of labeled pixels in the GSFC S2 dataset (left) and land cover classes
 196 (right).

197

198 2.1.3 Hollstein dataset (Sentinel-2)

199 The “S2 Hollstein dataset” is a database of manually labeled Sentinel-2A spectra of
200 clouds (Hollstein et al., 2016). By means of different spectral tools, pixels were selected and
201 classified into one of the following six classes (Figure 4): cloud (opaque clouds), cirrus
202 (cirrus, semi-transparent clouds and vapor trails), snow (snow and ice), shadow (shadows
203 from clouds, cirrus, mountains, buildings, etc.), water (lakes, rivers, seas), and clear-sky
204 (other remaining areas). Spectral tools include false-color composites of Sentinel-2 images,
205 image enhancements and graphical visualization of spectra. The aim was to create highly
206 heterogeneous classes with a balanced number of pixels. There were 59 total Sentinel-2
207 scenes and 1,593,911 reference (labelled) pixels.



208

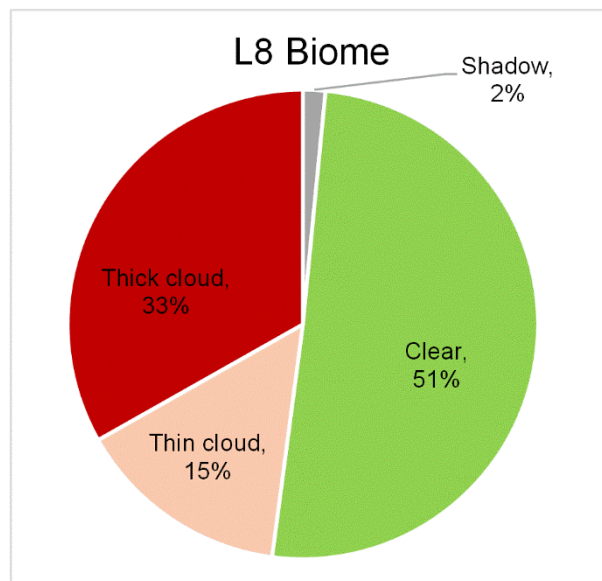
209 Figure 4. Distribution of labeled pixels in the Hollstein dataset.

210

211 2.1.4 L8Biome dataset (Landsat 8)

212 The “L8 Biome” cloud validation dataset consisted of 96 Landsat 8 scenes, which
213 were selected using a semi-random sampling by biome (Foga et al., 2017). These biomes
214 included barren, forest, grass/crops, shrubland, snow/ice, urban, water, and wetlands. For

215 each biome 12 Landsat 8 scenes were selected, and each scene was manually classified by an
216 expert into the following classes (Figure 5): clear, thin cloud, cloud, and cloud shadow. It
217 should be noted that no specific threshold was used to detect thin (semi-transparent) clouds,
218 which were primarily determined by the analyst. Also, the cloud shadow class in the
219 validation dataset was not provided for all the Landsat 8 scenes. The detailed description of
220 the L8Biome dataset is provided in Foga et al. (2017).



221

222 Figure 5. Distribution of labeled pixels in the L8Biome dataset.

223

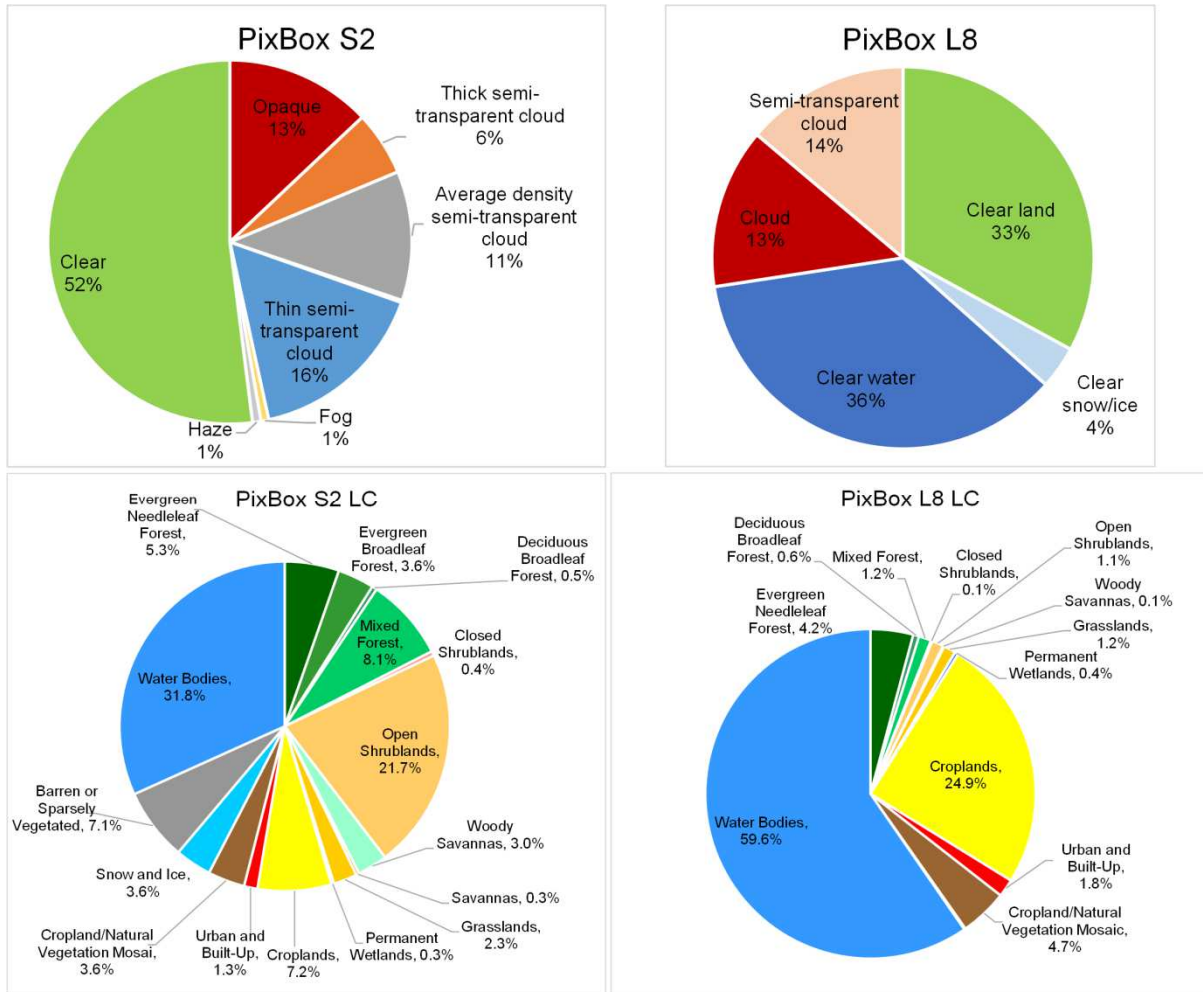
224 2.1.5 PixBox dataset (Landsat 8, Sentinel-2)

225 The overarching goal of the so called "PixBox" is to enable a quantitative assessment
226 of the quality of a pixel classification produced by an automated algorithm/procedure. Pixel
227 classification is defined as assigning a certain number of attributes to an image pixel, such as
228 cloud, clear sky, water, land, inland water, flooded, snow etc. These pixel classification
229 attributes are typically used to further guide higher level processing. PixBox is not only a
230 dataset but also includes a method comprising a procedure to define the best thematic, spatial
231 and temporal distribution for each collection purpose, a dedicated software for collecting

232 pixels, the analysis, comparing the collected reference against an automatic classification, as
233 well as the generation of a report.

234 For the PixBox Reference Dataset, a trained expert(s) manually labels pixels of an
235 image sensor into a detailed set of pre-defined classes. These are typically different cloud
236 transparencies, cloud shadow, and condition of the underlying surface (“semi-transparent
237 clouds over snow”, “clouds over bright scattering water”). The collected dataset includes 10’s
238 of thousands of pixels because it necessitates representation for all classes, and for various
239 observation and environmental conditions such as climate zones, solar illumination, viewing
240 angles, etc. Prior to the collection process the expert is provided with a detailed list of
241 distribution of categories and classes that needs to be fulfilled. During the collection process
242 the growing database is constantly checked against this reference. Quality control of the
243 collected pixels is important in order to detect misclassifications and systematic errors.

244 PixBox is a commercially sold product/service of Brockmann Consult GmbH. The
245 following two PixBox datasets have been made freely available to be used for CMIX
246 (Paperin 2021a, Paperin 2021b). The Sentinel-2 PixBox dataset contained 17,351 pixels (at
247 10 m) manually collected from 29 Sentinel-2A/B Level 1C products (top-of-atmosphere
248 reflectance—TOA reflectance). The Landsat 8 PixBox dataset contained 20,500 pixels (at
249 30 m) manually collected from 11 Landsat-8 Level 1 products (TOA reflectance). The
250 Sentinel-2 PixBox dataset is spatially, temporally, and thematically evenly distributed, while
251 the Landsat 8 dataset has a strong spatial focus on the Northern European coastal areas.
252 Distribution of labelled pixels and corresponding land cover classes for the PixBox datasets
253 are shown in Figure 6.



254

255 Figure 6. Distribution of labeled pixels and land cover classes in the PixBox dataset.

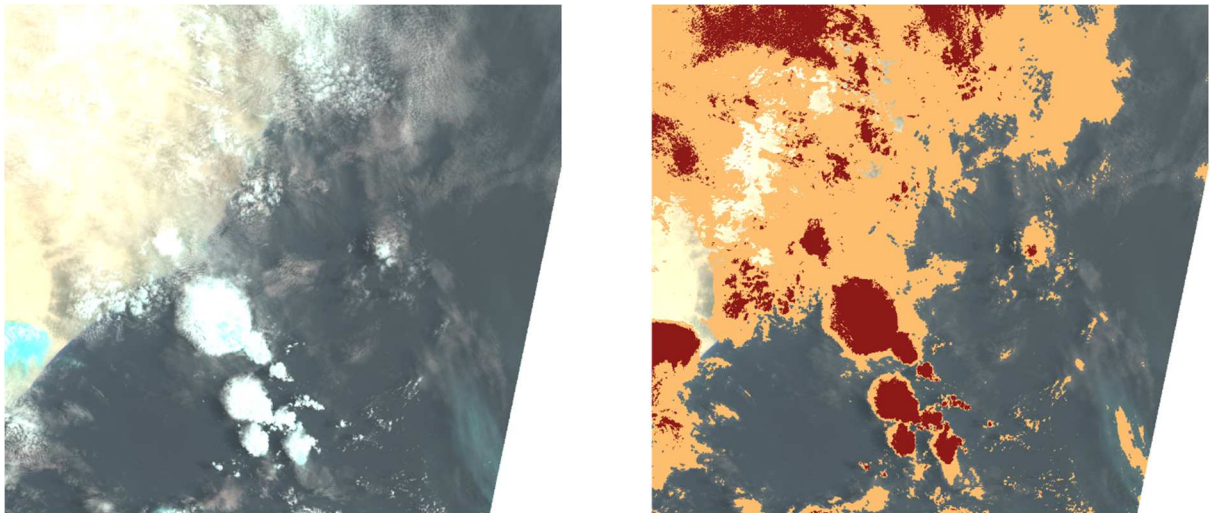
256

257 *2.1.6 Summary of strengths and limitations of cloud reference datasets*

258 Table 2 summarizes the strengths and limitations of cloud reference datasets used in
 259 this study. Reference data incorporating global coverage and a wide range of image
 260 conditions (L8Biome, PixBox, Hollstein) are based on the photointerpretation of images by
 261 an expert or a group of experts. This can introduce some subjectivity in labelling clouds,
 262 especially for thin/semi-transparent clouds that can be wavelength-dependent and fog
 263 (Scaramuzza et al., 2011) (Figure 7), and it is usually difficult to draw the exact boundary
 264 between this type of clouds and clear pixels. Another approach is to use high-quality pixels
 265 (with no uncertainties in cloud detection) and subsequently apply machine learning

266 algorithms to extrapolate classification for the whole scene through an iterative process until
267 the classification results assessed by an expert are deemed to be satisfactory (CESBIO)
268 (Figure 8). The quality of the resulting map, however, can still depend on the training data
269 and classification method used. A third approach (GSFC dataset) is to utilize ground-based
270 imagery of the sky to produce a training/validation cloud dataset, either through manual or
271 automatic labelling (Figure 8). While such an approach would potentially decrease
272 subjectivity in identifying clouds, a network of such sites with sky cameras would be required
273 (similar to the Aeronet network) in order to capture various geographical conditions.

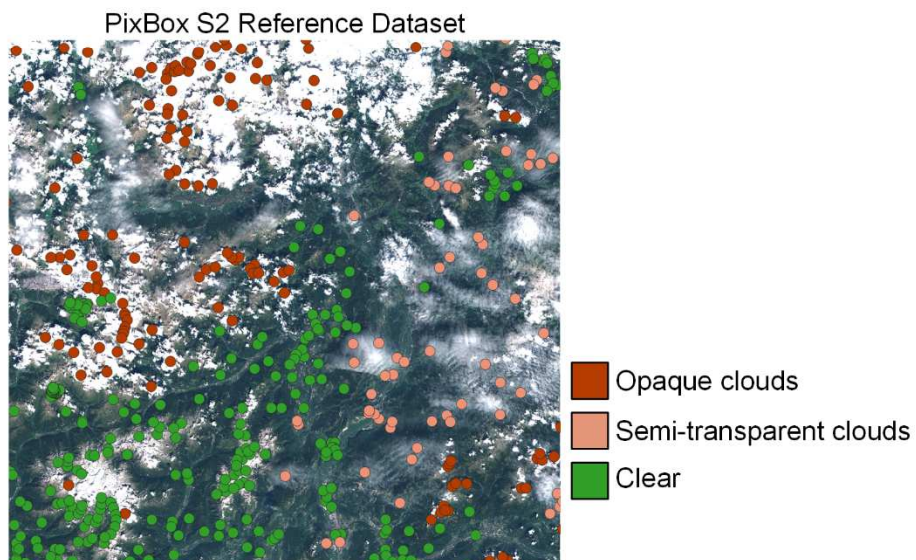
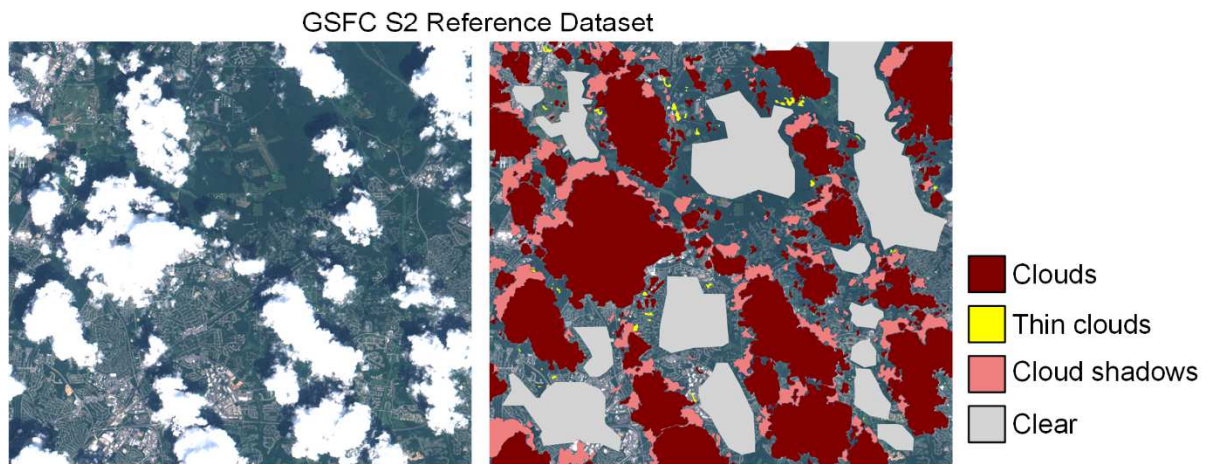
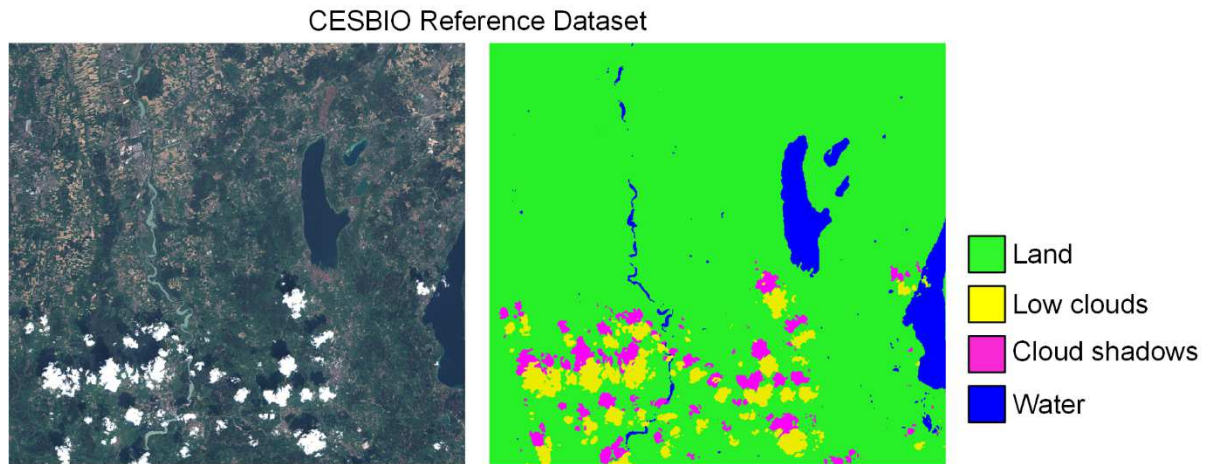
274



275

276 Figure 7. Part of the L8Biome scene (LC81570452014213LGN00) with some thin clouds not
277 labelled. Thin clouds are shown in orange, and thick clouds in maroon.

278



279

280 Figure 8. Examples of labeled data in the three datasets: CESBIO (fully labeled images);

281 GSFC (polygons avoiding uncertain areas, such as cloud boundaries); PixBox (sample-based

282 approach).

283

284 Table A1 (Appendix A) provides a list of classes from the reference datasets that were
 285 used to define cloud and non-cloud pixels in the CMIX. Most of the datasets were balanced
 286 in terms of cloud and non-cloud pixels, except of CESBIO, which had 24% of cloud pixels
 287 (Figure 2). CESBIO, GSFC and Hollstein datasets were primarily over the land surface, while
 288 the majority of PixBox datasets was over the water surface: 32% for S2 and 60% for L8.

289

290 Table 2. Strengths and limitations of cloud reference datasets.

Dataset	Strengths	Limitations
CESBIO	– All pixels in the scene are classified using an iteratively supervised machine learning approach	– Based on expert knowledge (potential bias). Small number of locations (limited spatial coverage) – Cloud and non-cloud areas unbalanced
GSFC	– Assisted with ground-based imagery – Over the same territory (can be potentially used for temporal consistency analysis)	– Limited field of view and single location – Surface classes limited to the location of sky camera – Cloud boundaries excluded
Hollstein	– Manual classification of polygons using spectral features	– Lack of sample quality – Low level of detail – Based on expert knowledge (potential bias) – Cloud edges not sampled
L8Biome	– Global coverage with stratified sampling – All pixels in the scene are classified	– Based on expert knowledge (potential bias)
PixBox	– High level of detail – High level of classification precision – Global coverage with stratified sampling	– Single pixel, thus a comparably small dataset – Based on expert knowledge (potential bias)

291

292 **2.2 Cloud masking algorithms**

293 This subsection briefly describes the main concepts utilized in each of the cloud
 294 masking algorithms with a summary presented in Table 3.

295 Table 3. Summary of cloud masking algorithms (L8: Landsat 8, S2: Sentinel-2). The “Objective” column shows the intended performance of
 296 algorithm in terms of cloud omission/commission errors. “Balanced” means the algorithms aims at balancing omission/commission errors.
 297 “Cloud-free conservative” means the algorithm aimed at minimizing cloud omission errors.

Processor	Organization	Methodology	Objective	Spatial resolution, m	Temporality	Buffer for clouds	Shadow detection	References
ATCOR	DLR	Spectral tests	Balanced	L8: 30 S2: 20	Mono	100 m	Yes	Richter & Schläpfer (2019a)
CD-FCNN	University of Valencia	Machine learning	Balanced	L8: 30 S2: 10/20/60	Mono	No	No	Mateo-García et al., (2020), López-Puigdollers et al. (2021)
Fmask 4.0 CCA	USGS	Spectral tests	Balanced	L8: 30 S2: 20	Mono	L8: 90 m S2: 60 m	Yeas	Foga et al. (2017), Qiu et al. (2019), Zhu et al. (2015)
FORCE	Humboldt-Universität zu Berlin / Trier University	Spectral test + parallax (S2 only)	Cloud-free conservative	L8: 30 S2: 10	Mono	300 m	Yes	Frantz (2019), Frantz et al. (2018), Frantz et al. (2016), Zhu et al. (2015), Zhu & Woodcock (2012)
IdePix	Brockmann Consult	Spectral tests	Balanced	S2: 20	Mono	Not used (user-defined)	Yes	Wevers et al. (2021)
InterSSIM	Sinergise	Machine learning + spatio-temporal context	Cloud-free conservative	S2: 10	Multi	160 m	No	Puc & Žust (2019)
LaSRC	NASA / University of Maryland	Spectral tests	Cloud-free conservative	L8: 30 S2: 10	Mono	L8: 150 m S2: 50 m	Yes	Skakun et al. (2019), Skakun et al. (2021), Vermote et al. (2016)
MAJA	CNES / CESBIO	Multi-temporal and spectral tests	Cloud-free conservative	S2: 240	Multi	240 m	Yes	Hagolle et al. (2010), Hagolle et al. (2017)
s2cloudless	Sinergise	Machine learning	Cloud-free conservative	S2: 10	Mono	160 m	No	Zupanc (2017)
sen2cor	ESA / Telespazio France	Spectral test + auxiliary data	Balanced	S2: 20	Mono	No	Yes	Louis et al. (2016), Louis (2021)

298

299 *2.2.1 ATCOR*

300 ATCOR is a generic atmospheric correction algorithm for mono-temporal multi-
301 /hyper-spectral satellite imagery in the solar reflective region (400 – 2500 nm) and thermal
302 region (8-13 μm) (Richter & Schläpfer, 2019b). The code uses MODTRAN5 look-up tables
303 for the radiative transfer functions. Separate codes exist for the processing of flat and rugged
304 terrain imagery. A preprocessing step calculates different masks (water cloud, cirrus cloud,
305 shadow, water) based on spectral tests. The cloud masking uses a buffer of 100 m. For
306 Landsat-8 and Sentinel-2 data the TOA reflectance threshold of the cirrus band is set to 0.01
307 (reflectance units). The lower threshold for thin cirrus detection was used prevent scenes with
308 very thin cirrus being classified as (thin) cirrus because other classes (e.g., water, shadow) are
309 generally of more interest than very thin cirrus. Cloud detection in ATCOR was aimed to
310 have a balance between commission and omission errors. In CMIX, ATCOR version 9.3.0
311 (2019) was used. CMIX processing of ATCOR did not use a Digital Elevation Model (DEM)
312 or any other auxiliary data. Some scenes from reference datasets were not processed by
313 ATCOR, since they were acquired with Sun elevation angle values less than 30° .

314

315 *2.2.2 CD-FCNN*

316 The cloud detection approach based on deep learning, proposed by the Image and
317 Signal Processing (ISP) group of the University of Valencia, is applicable to multispectral
318 images from moderate spatial resolution satellites, including Landsat 8 and Sentinel-2.
319 Training accurate global cloud detection models based on deep learning requires large
320 datasets of annotated images, which must reflect the high variability of clouds, surface, and
321 atmospheric conditions. This is a major difficulty since high-quality labeled datasets usually
322 do not exist or are not publicly available for most satellite sensors. For Landsat 8, the
323 L8Biome dataset matches these requirements (Jeppesen et al., 2019). However, similar global

324 datasets do not exist for Sentinel-2 yet. (Sentinel-2 Cloud Mask Catalogue (Francis et al.,
325 2020) was made available after CMIX was initiated). Therefore, Landsat 8 datasets
326 (L8Biome, 80%, and L8SPARCS, 20%) were used to train fully convolutional neural
327 networks (FCNN) that may be transferred to perform cloud detection in Sentinel-2 images.
328 L8SPARCS (Spatial Procedures for Automated Removal of Cloud and Shadow) (USGS,
329 2016) was created for the validation of the cloud detection approach proposed by Hughes &
330 Hayes (2014). It consists of 80 Landsat-8 sub-scenes manually labeled in five different
331 classes: cloud, cloud-shadow, snow/ice, water, flooded, and clear-sky. The size of each sub-
332 scene is 1000x1000 pixels.

333 After a minimum adaptation of Sentinel-2 data, in terms of band selection and spatial
334 resolution, the models trained on Landsat 8 data are directly applied to Sentinel-2 images.
335 The proposed neural network architecture is based on a modified U-Net with significantly
336 less training parameters and lower computational cost (Mateo-García et al., 2020). It seeks to
337 provide both faster inference time and accurate detection through a lightweight architecture
338 with a moderate number of parameters, i.e., approximately 96,000 parameters, which is
339 around 1% of original U-Net parameters. Moreover, this modified version of U-Net works
340 seamlessly with Landsat-8 and Sentinel-2 images thanks to a transfer learning strategy over
341 both sensors. In this way, all input bands, regardless of the sensor, are homogenized and
342 resampled to 30m overlapping patches of 32x32 pixels, which are used for training the
343 networks in a 64-batch size configuration. Models are trained to minimize a pixel-wise binary
344 cross-entropy cost function, between ground truth and predictions, using the Adam stochastic
345 gradient descent optimization algorithm. An initial learning rate of 10^{-5} , a weight decay of $5 \times$
346 10^{-4} and 120 epochs were used to train the final network. The TensorFlow framework was
347 used to implement and train the models on a GPU (average of 800 s/epoch in all
348 configurations). Training and testing details can be found in López-Puigdollers et al. (2021);

349 in addition, the pre-trained model and a Python-implementation of the proposed cloud
350 detection algorithm for Landsat-8 and Sentinel-2 is provided in a public repository
351 (<https://github.com/IPL-UV/DL-L8S2-UV>).

352 Since we propose to use the same model for Landsat-8 and Sentinel-2, we are
353 restricted to bands available in both sensors. In this context, three different bands
354 configurations were tested: "RGBI" corresponds to bands B2, B3, B4 and B5 of Landsat-8
355 and B2, B3, B4 and B8 of Sentinel-2; "RGBISWIR" to bands B2, B3, B4, B5, B6 and B7 of
356 Landsat-8 and B2, B3, B4, B8, B11 and B12 of Sentinel-2; and "ALLNT" includes all
357 "RGBISWIR" bands plus the coastal aerosols and cirrus bands (B1 and B9 in Landsat-8, B1
358 and B10 in Sentinel-2, respectively). After internal testing, the network selected for
359 benchmarking in CMIX was the "RGBISWIR" network. Further results about the different
360 band configurations can be found in López-Puigdollers et al. (2021).

361 The CD-FCNN output is given by a sigmoid activation function that provides
362 continuous values, which could be interpreted as probabilities, between 0 and 1. In order to
363 compare with the rest of the methods, these values are binarized into "non-cloud" (0) or
364 "cloud" (1) classes for each pixel. We set a default 0.5 threshold to obtain the binary cloud
365 mask assuming unbiased data. However, this threshold has a crucial importance in terms of
366 balance between commission and omission and errors. In Landsat-8 images both errors are
367 similar, but performance may decrease in complex scenarios with presence of ambiguous
368 pixels, e.g. over snow, urban areas or coastal lines. Adjusting this threshold for a specific
369 dataset may improve the tradeoff between omission and commission errors depending on the
370 requirements of the application, i.e. cloud or cloud-free conservative applications. The
371 resulting cloud mask is spatially resampled from the native Landsat 8 resolution of 30 m to
372 the corresponding Sentinel-2 resolutions of 10, 20 and 60 m. Throughout the entire process

373 the work is done at a pixel level, and no spatial dilation of the cloud mask is considered at
374 any stage.

375

376 *2.2.3 Fmask 4.0 CCA*

377 Function of Mask (Fmask) 4.0 is a cloud assessment algorithm used with Landsat and
378 Sentinel-2 imagery (Qiu et al., 2019). An earlier version, Fmask 3.3, is applied operationally
379 to create cloud masks for USGS Landsat products. The algorithm provisionally identifies
380 cloud pixels using spectral tests, then matches those pixels to provisional cloud shadow pixels
381 using sensor geometry, the Digital Elevation Model (DEM) of the terrain, and an iterative
382 search of altitudes (in Landsat imagery). Fmask was designed to provide a balance between
383 cloud commission and omission errors. Fmask 4.0 is available under an MIT license at
384 <https://github.com/GERSL/Fmask>.

385

386 *2.2.4 FORCE*

387 FORCE (Framework for Operational Radiometric Correction for Environmental
388 monitoring, <https://github.com/davidfrantz/force>) is developed as an ‘all-in-one’ open-source
389 software solution for the mass-processing and analysis of Landsat and Sentinel-2 image
390 archives (Frantz, 2019). FORCE includes a mono-temporal Level 2 processing system for
391 Analysis Ready Data (ARD) generation which includes: radiometric correction, cloud
392 masking, and data cube generation (Frantz et al., 2016). The cloud masking has branched
393 from Fmask version 1.6.3 (Zhu & Woodcock, 2012), and since then has been developed in
394 parallel (Frantz et al., 2015; Frantz et al., 2016; Frantz et al., 2018). Parts of the updates in
395 Zhu et al. (2015) were incorporated. A darkness filter was implemented to mitigate false
396 positives in bifidly structured dryland areas, where the scene-based temperature distribution
397 tests for Landsat can result in commission errors of cold image parts (Frantz et al., 2015).

398 Cirrus masking is based on an elevation-dependent equation (Baetens et al., 2019). The most
399 notable difference to the original Fmask, however, is the complete replacement of the cloud
400 probability module for Sentinel-2 with a new algorithm that makes use of the Cloud
401 Displacement Index, which is formulated to enhance parallax effects in highly correlated NIR
402 bands (Frantz et al., 2018). The FORCE cloud masking aims to aggressively detect clouds
403 and cloud shadows to increase cloud producer's accuracy at the deliberate expense of cloud
404 commission for its safe operation in time-series applications. Circular buffers are used to
405 reduce false negatives (300 m for opaque clouds). FORCE provides quality bits whereby 12
406 quality indicators with respect to atmospheric conditions are provided (Frantz, 2019).
407 Multiple indicators can be set simultaneously for each pixel, e.g., snow and cloud. This
408 quality product is generated at 30 m and 10 m resolution for Landsat and Sentinel-2,
409 respectively. FORCE v. 3.0-dev was used in CMIX.

410

411 *2.2.5 Idepix*

412 IdePix (Identification of Pixel properties) is a multi-sensor pixel identification tool
413 available as a SNAP (Sentinel Application Platform) plugin (Wevers, 2021). It provides pixel
414 identification algorithms for a wide variety of sensors such as Sentinel-2 MSI, Sentinel-3
415 OLCI, MERIS, Landsat-8, MODIS, VIIRS, Proba-V or SPOT VGT. IdePix classifies pixels
416 into a series of categories (flags) for further processing using a mono-temporal approach and
417 background information. Its uniqueness consists of a certain set of flags, which are calculated
418 for all instruments (common flags), complemented by instrument specific flags (instrument
419 flags). The technical design of all IdePix is instrument specific and can include decision trees,
420 probabilistic combination of calculated features or neural networks. The Sentinel-2 IdePix is
421 mainly based on a decision tree technique for cloud calculation as well as geometric
422 calculations for cloud and mountain shadows. In contrast to many other pixel identification

423 tools the final IdePix classification is non-exclusive and therefore allows multiple classes to
424 be set for a single pixel. This means a single pixel can have multiple properties such as land
425 and cloud (semi-transparent cloud over land), land and snow (land covered with snow), or
426 land, snow and cloud (semi-transparent cloud over snow covered land). This type of
427 implementation allows the most versatile usage of the flagging and combinations according
428 to users' needs compared to a standard integer flag allowing a single status per pixel.
429 Sentinel-2 IdePix derives water cloud flags and cirrus cloud flags ($B10 > 0.01$ & elevation $<$
430 2000 m) on multiple confidence levels, as well as cloud shadow, mountain shadow, snow/ice
431 and water flags. The pixel identification (IdePix) for Sentinel-2 is only working at single
432 resolution (i.e., 10 m, 20 m, 60 m). Cloud boundary pixels are flagged using a dilation filter.
433 In principle, cloud boundaries are regarded as neighbor pixels of a cloud as identified before
434 by the processor; therefore, a buffer is set around the cloud. The width of this boundary (in
435 number of pixels) can be set by the user. Usage of the buffering functionality was not
436 however utilized for CMIX to validate the sole performance of the cloud detection algorithm.

437

438 *2.2.6 S2cloudless*

439 The s2cloudless is an automated cloud-detection algorithm for Sentinel-2 imagery
440 (Zupanc, 2017) based on a gradient boosting algorithm. It was developed by the EO Research
441 team at Sinergise and is published under the MIT License on [https://github.com/sentinel-](https://github.com/sentinel-hub/sentinel2-cloud-detector)
442 [hub/sentinel2-cloud-detector](https://github.com/sentinel-hub/sentinel2-cloud-detector). The model was trained on a large training dataset with a global
443 coverage. The algorithm is monotemporal, does not consider any spatial context, and
444 therefore can be executed at any resolution. The s2cloudless algorithm can, unlike many
445 other algorithms, be executed also on averaged Sentinel-2 reflectance values over arbitrary
446 user-defined geometries and still provide meaningful results. The input features are Sentinel-
447 2 Level-1C TOA reflectance values of the following ten bands: B01, B02, B04, B05, B08,

448 B8A, B09, B10, B11, B12 and output of the algorithm is a cloud probability map. Users of
449 the algorithm can convert the cloud probability map to a cloud mask by thresholding the
450 cloud probability map. The recommended value for the threshold is 0.4 to minimize cloud
451 omission errors. Users can optionally apply additional morphological operations during the
452 conversion of the cloud probability map to the cloud mask. These operations are as follows:
453 convolution of the probability map and dilation of the binary cloud mask with a disk. We
454 recommend convolving cloud probability maps at 10 m (160 m) resolution with a disk with a
455 radius of 22 (2) px and dilate cloud masks with a disk with radius 11 (1) px. Sentinel Hub
456 (<https://www.sentinel-hub.com>, details in EO Research team (2020)) and Google Earth
457 Engine ([https://developers.google.com/earth-](https://developers.google.com/earth-engine/datasets/catalog/COPERNICUS_S2_CLOUD_PROBABILITY)
458 [engine/datasets/catalog/COPERNICUS_S2_CLOUD_PROBABILITY](https://developers.google.com/earth-engine/datasets/catalog/COPERNICUS_S2_CLOUD_PROBABILITY)) provide precomputed
459 s2cloudless cloud probability maps and masks to their users for the entire Sentinel-2 archive.

460 The s2cloudless cloud masks for CMIX were provided in a binary mode (1 – cloud
461 and 0 – non-cloud) using the latest (v0.1) model and default values for threshold and
462 morphological operations.

463

464 *2.2.7 InterSSIM*

465 The InterSSIM cloud detection algorithm is a multi-temporal extension of the
466 s2cloudless algorithm (section 2.2.6), but unlike s2cloudless, the InterSSIM algorithm takes
467 temporal and spatial contexts into account. The algorithm was developed by the EO Research
468 Team at Sinergise (Puc & Žust, 2019) and integrated into the eo-learn Python library
469 published under the MIT License on <https://github.com/sentinel-hub/eo-learn>. The input data
470 and parameters for the InterSSIM are same as in s2cloudless (see section 2.2.6) with the
471 addition of prior satellite observations. The algorithm works on the ten Sentinel-2 TOA
472 bands, and in addition to cloud probabilities from the s2cloudless model incorporates

473 additional features: spatially averaged reflectance values, minimum and mean reflectance
474 values over all prior observations, and maximum, mean, and standard deviation of structural
475 similarity indices computed between the observation for which cloud mask is being predicted
476 and every other prior observations. The output of the algorithm is a cloud probability map for
477 the target timeframe, which can be converted into a cloud mask with the same procedure as in
478 the case of the s2cloudless algorithm.

479 The InterSSIM cloud masks for CMIX were provided in a binary mode (1 – cloud and
480 0 – non-cloud) using the latest (v0.1) s2cloudless model with default parameter values.

481

482 2.2.8 *LaSRC*

483 The Land Surface Reflectance Code (LaSRC) is a generic atmospheric correction
484 algorithm aimed at removing atmospheric effects associated with optical satellite imagery
485 acquisitions (Doxani et al., 2018; Vermote et al., 2016). The code is based on the inversion of
486 the 6SV radiative transfer code (Kotchenova et al., 2006; Vermote et al., 1997). Within the
487 atmospheric correction process, LaSRC generates several quality assurance (QA) layers,
488 including a cloud mask. The main metric for deriving a cloud mask is a per-pixel inversion
489 residual error (Skakun et al., 2019; Skakun et al., 2021; Vermote et al., 2016), which shows
490 the goodness of aerosol optical thickness (AOT) estimation process. For both Landsat 8 and
491 Sentinel-2, we used a threshold of 0.05 for the residual to identify cloudy pixels and to
492 minimize cloud omission errors, so only high-quality pixels will be used for further
493 processing. Pixels adjacent to clouds within 5 pixels are separately masked as “adjacent to
494 clouds”. For S2, a conservative threshold of 0.003 (reflectance units) was used for the cirrus
495 band. Therefore, for LaSRC pixels identified as cloud or adjacent were used as “cloud”,
496 whereas all others were used as “non-cloud”. In CMIX, LaSRC version 3.5.5 was used.

497

498

499

500

501 2.2.9 MAJA

502 MAJA is applicable to satellites which perform repetitive observations at similar
503 viewing angles, such as Sentinel-2. It was developed by CNES with methods designed by
504 CESBIO with a few modules provided by DLR. MAJA is an open-source software.

505 MAJA's cloud and shadow detection methods include several tests, which use the
506 multi-spectral and multi-temporal properties of surfaces, clouds, and shadows to classify
507 different types of pixels. The methods are described in Hagolle et al. (2010) and Hagolle et
508 al. (2017). The main cloud test detects the pixels for which the surface reflectance in the blue
509 band increases sharply. The cloud masks obtained with MAJA are dilated by 240 m, firstly to
510 account for the parallax effects due to differences in observation angles between spectral
511 bands, and secondly for the adjacency effects of clouds and for their 'fuzzy' borders. MAJA
512 aims at a sensible reliability for surface reflectance monitoring, its tests and thresholds are
513 therefore optimized to minimize cloud or cloud shadow omission (aiming at maximizing
514 producer's accuracy for clouds, but balanced for cirrus clouds), without excessively
515 degrading the commission error. Cirrus band is used to detect high clouds using the following
516 equation: $\text{Cirrus} > 0.007 + 0.007 \times h^2$. where h is the pixel altitude in km above sea level.

517 In CMIX, the cloud masks for Sentinel-2 were computed at 240 m resolution to
518 optimize the computation time, but this can prevent MAJA from detecting very small clouds.
519 In the more recent MAJA versions the clouds and shadows masks are computed at 120 m,
520 which should further improve MAJA's performance. MAJA has been intensively validated
521 and some of its validation data sets (Baetens et al., 2019) were used in the CMIX experiment.
522 Due to the necessity to process times series of data with a processed data volume 10 times

523 greater than the other algorithms, the MAJA team was not able to process all the data sets
524 submitted to CMIX, and it was decided to only produce the datasets acquired when both
525 Sentinel-2A and -2B satellites were operational.

526

527 *2.2.10 Sen2Cor*

528 Sen2Cor is a processor for Sentinel-2 Level 2A product generation; it performs the
529 atmospheric correction of the Top-Of-Atmosphere (TOA) Level 1C input data. It is
530 composed of two main modules: an atmospheric correction module and a scene classification
531 module that provides a “Scene Classification Map” (SCL), which is used internally in the
532 atmospheric correction module to distinguish between cloudy, clear and water pixels. The
533 Sen2Cor processor is used by the European Space Agency to generate Sentinel-2 Level-2A
534 products within the Sentinel-2 ground segment. Sen2Cor software is available for download
535 at <https://step.esa.int/main/third-party-plugins-2/sen2cor/>. The code is open source and
536 written in Python.

537 The Sen2Cor version 2.8 cloud screening algorithm (Louis et al., 2016; Louis, 2021)
538 uses the reflective properties of scene features (TOA reflectance). Potential cloudy pixels
539 undergo a sequence of filtering based on spectral bands thresholds, ratios, and indexes
540 computations (Normalized Difference Snow Index – NDSI, Normalized Difference
541 Vegetation Index –NDVI). Sen2Cor was designed to provide a balance between cloud
542 omission and commission errors. In addition, it includes a cirrus and cloud shadow detection
543 algorithm. A series of additional steps to improve the quality of the classification are
544 automatically triggered using a priori information: digital elevation model (DEM)
545 information, ESA CCI Water Bodies Map v4.0 (Lamarche et al., 2017), ESA CCI Land
546 Cover Map v.2.0.7 (2015) and a snow climatology.

547 In CMIX, Sen2cor version 2.8 was used. SCL classes 8, 9 and 10 were used for cloud
 548 and the remaining SCL classes for non-cloud.

549

550

551 **2.3 Performance metrics**

552 A standard set of classification metrics derived from confusion matrices (Table 4) was
 553 used to compare cloud masking algorithms and included (Table 5) overall accuracy (OA) and
 554 balanced OA (BOA), producer’s (PA) and user’s accuracies (UA). BOA (Brodersen et al.,
 555 2010) was used in addition to OA since some of the reference datasets were imbalanced in
 556 terms of cloud/clear pixels and therefore BOA would be a better indicator of algorithms
 557 performance.

558

559 Table 4. Confusion matrix for cloud validation.

		Reference	
		Cloud	Non-cloud
Map	Cloud	$n_{\text{cloud_as_cloud}}$	$n_{\text{ncloud_as_cloud}}$
	Non-cloud	$n_{\text{cloud_as_ncloud}}$	$n_{\text{ncloud_as_ncloud}}$

560

561 Table 5. Main performance metrics.

Metric	Equation
Overall accuracy (OA)	$\frac{n_{\text{cloud_as_cloud}} + n_{\text{ncloud_as_ncloud}}}{n_{\text{cloud_as_cloud}} + n_{\text{ncloud_as_ncloud}} + n_{\text{ncloud_as_cloud}} + n_{\text{cloud_as_ncloud}}} \quad (1)$
Balanced OA (BOA)	$0.5 \left(\frac{n_{\text{cloud_as_cloud}}}{n_{\text{cloud_as_cloud}} + n_{\text{cloud_as_ncloud}}} + \frac{n_{\text{ncloud_as_ncloud}}}{n_{\text{ncloud_as_cloud}} + n_{\text{ncloud_as_ncloud}}} \right) \quad (2)$
PA (for clouds)	$\frac{n_{\text{cloud_as_cloud}}}{n_{\text{cloud_as_cloud}} + n_{\text{cloud_as_ncloud}}} \quad (3)$
UA (for clouds)	$\frac{n_{\text{cloud_as_cloud}}}{n_{\text{cloud_as_cloud}} + n_{\text{ncloud_as_cloud}}} \quad (4)$

562

563 Performance metrics were estimated from confusion matrices that incorporated all
564 valid pixels over all scenes available in the dataset. PA is complementary to the omission
565 error, which shows a fraction of missed clouds; UA is complementary to the commission
566 error, which shows a fraction of over detected clouds. High PA (cloud-free, non-cloud or
567 clear conservative) means that after elimination of clouds, the users results will be minimally
568 affected by remaining clouds, while high UA (cloud conservative) means that the cloud
569 masks will not discard supernumerary valid pixels.

570

571 **3 Results**

572 *3.1 Performance of cloud masking algorithms for Sentinel-2*

573 *3.1.1 CESBIO reference dataset*

574 Table 6 and Figure 9 show performance metrics when applying cloud masking
575 algorithms on the Sentinel-2 CESBIO dataset. Several observations can be made when
576 analyzing these results. The number of reference pixels varied, since the CESBIO dataset was
577 generated at 60 m spatial resolution, and processors produced masks at various spatial
578 resolution: 10 m (FORCE, InterSSIM, LaSRC and S2cloudless), 20 m (ATCOR, Idepix,
579 Fmask 4.0 CCA, Sen2Cor), 60 m (CD-FCNN, interpolated from 30 m), and 240 m (MAJA).
580 Cloud and non-cloud classes were imbalanced in the reference dataset (of all labelled pixels
581 24.3% were clouds), therefore it results in the OA to be biased towards the non-cloud
582 (dominant) class. Therefore, the balanced OA (BOA) is a more appropriate metric. Overall,
583 BOA varied from 79.5% to 90.5%, an average of $85.9 \pm 3.7\%$. When not considering MAJA
584 (whose developers generated the CESBIO dataset), the highest cloud PA was 85.6%, with the
585 average being $75.9 \pm 8.7\%$, meaning that most algorithms missed almost 24% of clouds
586 identified in the CESBIO dataset. Average cloud UA without MAJA was $85.1 \pm 10.6\%$,

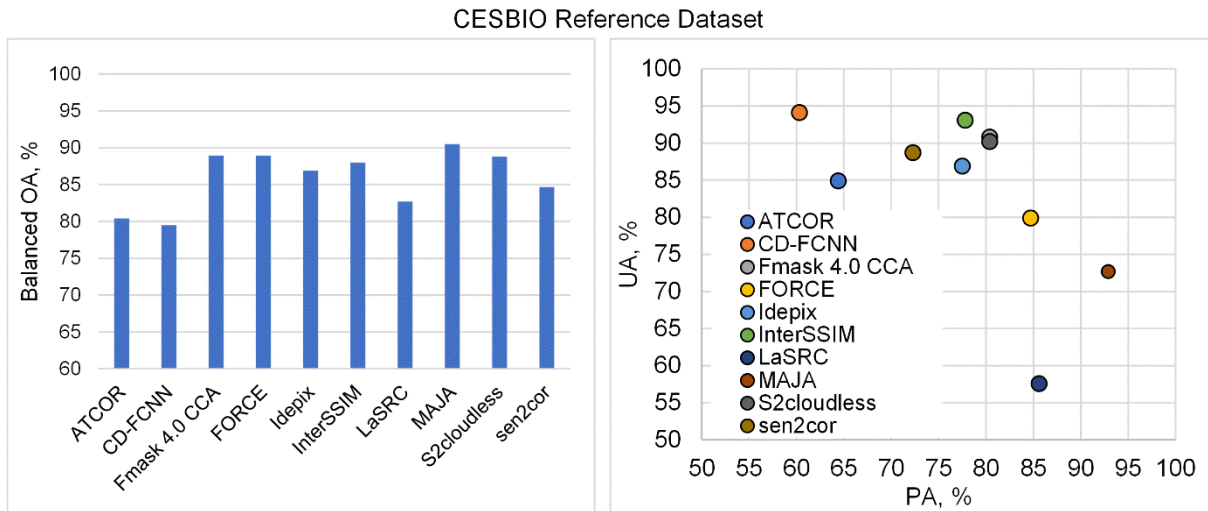
587 meaning an average of 15% over detection of clouds, which may lie in the dilated parts of the
 588 cloud masks (FORCE, MAJA), or be associated with a stricter detection of cirrus clouds
 589 (LaSRC). Overall, the performance of cloud masking algorithms varied for this dataset by an
 590 average 11-12% of PA and UA, as measured by the coefficient of variation (CV), which is a
 591 ratio between standard deviation and average.

592

593 Table 6. Performance metrics of Sentinel-2 cloud masking algorithms for the CESBIO
 594 dataset. All algorithms, except MAJA, processed all 30 reference scenes (with 24.3% of
 595 clouds in the reference dataset), while MAJA processed 28 references scenes (25.6%). Here,
 596 and in Table 7 through Table 14: in bold are the numbers with the highest value for the
 597 particular metric (column-wise); * denotes algorithms which did not process the whole
 598 dataset; algorithms that are underscored were produced by the same team as the reference
 599 dataset.

Processor	Cloud			
	OA	BOA	PA	UA
ATCOR	88.6	80.4	64.4	84.9
CD-FCNN	89.5	79.5	60.3	94.1
Fmask 4.0 CCA	93.3	88.9	80.4	90.8
FORCE	91.1	88.9	84.7	79.9
Idepix	91.7	86.9	77.5	86.9
InterSSIM	93.2	88.0	77.8	93.1
LaSRC	81.2	82.7	85.6	57.6
<u>MAJA*</u> (28/30)	89.2	90.5	92.9	72.7
S2cloudless	93.1	88.8	80.4	90.2
sen2cor	91	84.7	72.3	88.7
<i>Average</i>	<i>90.2</i>	<i>85.9</i>	<i>77.6</i>	<i>83.9</i>
<i>Standard deviation</i>	<i>3.4</i>	<i>3.7</i>	<i>9.3</i>	<i>10.7</i>

600



601

602 Figure 9. Comparison of BOA values and distribution of PA/UA for the CESBIO reference
 603 dataset.

604

605 *3.1.2 GSFC S2 reference dataset*

606 Table 7 and Figure 10 show the results of comparing algorithm outcomes against the
 607 Sentinel-2 GSFC dataset. MAJA provided only 10 images out of 28 images. In the S2 GSFC
 608 dataset, cloud and non-cloud are almost balanced (approx. 61% of reference pixels are
 609 identified as clouds), therefore there is minimal difference between OA and BOA. BOA
 610 varied from 80.7% to 96.8% with LaSRC being the outlier (developers of LaSRC produced
 611 the GSFC data), with average being $85.7 \pm 2.8\%$ (not considering LaSRC). Average values of
 612 cloud PA and UA not considering LaSRC were $73.7 \pm 5.6\%$ and $98.2 \pm 2.7\%$, respectively,
 613 meaning large omission errors. It is worth noting that FORCE and MAJA, whose PA was
 614 better than the UA for the other reference datasets, have the opposite result for the GSFC
 615 reference, due to the strict classification of very thin clouds as clouds in the GSFC
 616 data set. The reason for all algorithms producing lower accuracies compared to LaSRC is that
 617 they did not identify thin (semi-transparent and cirrus) clouds, which, in turn, LaSRC was
 618 masking out using a rather conservative threshold (0.003 in reflectance units; for
 619 LaSRCv3.5.5) applied for the cirrus band (B10). As the cirrus cloud masking method is very

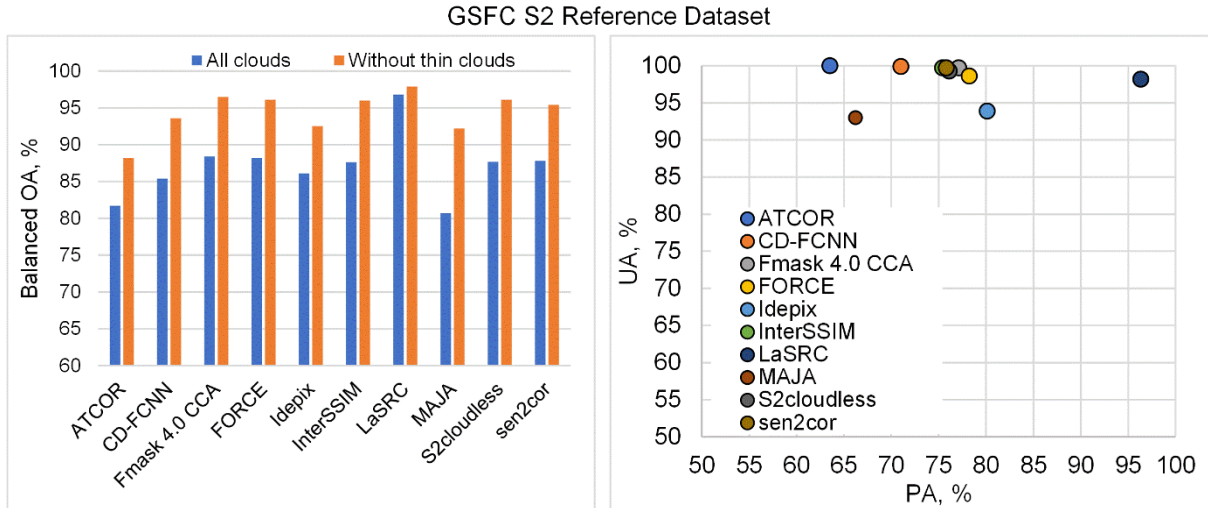
620 simple, all methods could obtain similar performances, at the expense of masking an
621 important part of usable pixels. Those clouds were labelled as thin, since they were clearly
622 visible in the ground-based images. If thin clouds are removed from the analysis (Table 7), all
623 algorithms showed much better performance: average BOA was $94.4 \pm 2.9\%$ (an average gain
624 $+7.4 \pm 2.6\%$) and cloud PA was $90.8 \pm 5.9\%$ (an average gain $+14.8 \pm 5.2\%$), while cloud-UA
625 remained essentially the same $98.1 \pm 2.7\%$. These results show the differences between
626 algorithms in defining and identifying thin (semi-transparent) cirrus clouds, at the same time
627 mostly agreeing on thick clouds. Variation in algorithms performance was 8% for cloud PA
628 (6% without thin clouds) and 3% for cloud UA.

629

630 Table 7. Performance metrics of Sentinel-2 cloud masking algorithms for the GSFC S2
631 dataset. All algorithms, with exception of MAJA, processed all 28 reference scenes (with
632 60.6% and 55.5% of clouds in reference data for all clouds and without thin clouds,
633 respectively), while MAJA processed 10 images (49.2% and 40.8%).

	All types of clouds				Without thin clouds			
	OA	BOA	Cloud		OA	BOA	Cloud	
Processor	OA	BOA	PA	UA	OA	BOA	PA	UA
ATCOR	77.9	81.7	63.5	100	86.9	88.2	76.4	100
CD-FCNN	82.4	85.4	71	99.9	92.9	93.6	87.3	99.9
Fmask 4.0 CCA	86	88.4	77.1	99.7	96.1	96.5	93.3	99.7
FORCE	86.1	88.2	78.2	98.6	95.9	96.1	94	98.5
Idepix	84.8	86.1	80.1	93.9	92.5	92.5	92.9	93.6
InterSSIM	85	87.6	75.4	99.7	95.6	96	92.4	99.7
<u>LaSRC</u>	96.7	96.8	96.3	98.2	98	97.9	98.5	97.8
MAJA* (10/28)	80.9	80.7	66.2	93	92.7	92.2	89.1	92.7
S2cloudless	85.2	87.7	76.1	99.3	95.7	96.1	93	99.3
sen2cor	85.2	87.8	75.8	99.7	95	95.4	91.2	99.7
<i>Average</i>	<i>85.0</i>	<i>87.0</i>	<i>76.0</i>	<i>98.2</i>	<i>94.1</i>	<i>94.5</i>	<i>90.8</i>	<i>98.1</i>
<i>Standard deviation</i>	<i>4.6</i>	<i>4.1</i>	<i>8.4</i>	<i>2.4</i>	<i>2.9</i>	<i>2.7</i>	<i>5.6</i>	<i>2.6</i>

634



635

636 Figure 10. Comparison of BOA values and distribution of PA/UA (for all clouds) for the
 637 GSFC S2 reference dataset.

638

639 3.1.3 Hollstein reference dataset

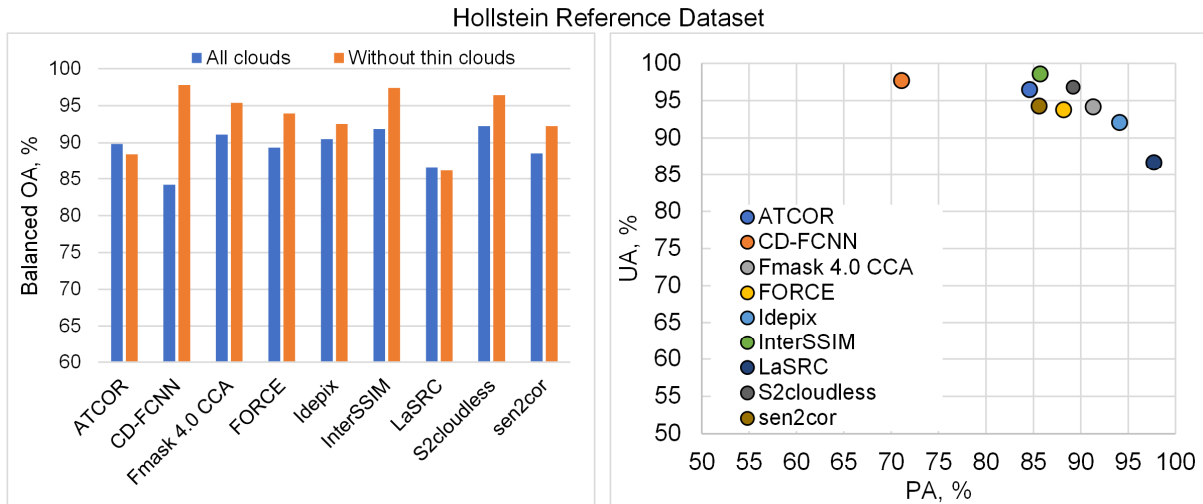
640 Table 8 and Figure 11 show algorithms performance for the Hollstein data depending
 641 on the opaque and semi-transparent/cirrus clouds. BOA varied from 84.2% to 92.3% (average
 642 $89.4 \pm 2.4\%$) for all cloud types and 86.2 to 97.8% ($93.4 \pm 3.8\%$) for opaque clouds only. Not
 643 considering semi-transparent/cirrus clouds improved algorithms performance, especially for
 644 cloud PA: an average gain $+8.0 \pm 8.1\%$. Variation of performance was comparable to the
 645 GSFC results with 8% (5% for opaque only) for PA and 4% (7%) for UA. Note that the
 646 Hollstein dataset was used to set radii of disks with which the cloud probability mask and
 647 binary cloud mask are convoluted and dilated, respectively, by the s2cloudless algorithm.
 648 MAJA was not evaluated against the Hollstein data set, as the images were acquired before
 649 Sentinel-2B launch.

650

651 Table 8. Performance metrics of cloud masking algorithms for the Hollstein dataset. All
 652 algorithms processed all 59 reference scenes (with 61.8% and 44.4% of clouds in reference
 653 data for all clouds and without thin clouds, respectively).

Processor	Opaque clouds and semi-transparent clouds/cirrus				Opaque clouds only			
	Cloud		Cloud		Cloud		Cloud	
	OA	BOA	PA	UA	OA	BOA	PA	UA
ATCOR	88.6	89.9	84.6	96.5	89.1	88.5	81.8	93.2
CD-FCNN	81	84.2	71.1	97.7	97.8	97.8	98.3	96.7
Fmask 4.0 CCA	91.2	91.1	91.3	94.2	94.9	95.4	99.9	89.8
FORCE	89.1	89.4	88.2	93.8	93.6	94	97.4	89.1
Idepix	91.3	90.5	94.1	92.1	91.9	92.6	98.2	85.7
InterSSIM	90.4	91.9	85.7	98.6	97.5	97.4	96.8	97.5
LaSRC	89.3	86.7	97.7	86.7	85	86.2	96.7	76
S2cloudless	91.5	92.3	89.2	96.8	96.3	96.5	97.6	94.3
sen2cor	87.9	88.6	85.6	94.3	92.2	92.3	93	89.8
<i>Average</i>	<i>88.9</i>	<i>89.4</i>	<i>87.5</i>	<i>94.5</i>	<i>93.1</i>	<i>93.4</i>	<i>95.5</i>	<i>90.2</i>
<i>Standard deviation</i>	<i>3.1</i>	<i>2.4</i>	<i>7.1</i>	<i>3.4</i>	<i>3.9</i>	<i>3.8</i>	<i>5.2</i>	<i>6.2</i>

654



655

656 Figure 11. Comparison of BOA values and distribution of PA/UA (for all clouds) for the
657 Hollstein reference dataset.

658

659 3.1.4 PixBox S2 reference dataset

660 Not all algorithms processed all 29 products of the PixBox S2 dataset; the reasons for
661 this were limitations of allowed geometries (ATCOR, 27 processed) or too sparse time-series
662 around the acquisition (MAJA, 14 processed). In order to account for the difference of
663 available products for validation, two different comparisons were made: one using all

664 available products for each algorithm and a second using only the products that all algorithms
 665 have been applied to (14 out of 29 reference scenes). We call the second dataset the least
 666 common denominator (LCD) subset, while the first is referred to as the “complete dataset”.
 667 The whole comparison could have been made only on the LCD subset, but this reduces the
 668 complete dataset by half, which reduces its utility. Therefore, the complete dataset also was
 669 used for comparison. In this comparison using the complete dataset, results for MAJA must
 670 be assessed with caution, as they are only based on 14 out of 29 products.

671 Algorithm performance for the complete PixBox dataset is provided in Table 9 and
 672 Figure 12. For all types of clouds, BOA varied from 67.5% to 85.9% (average $80.0 \pm 5.3\%$).
 673 The top two algorithms (S2cloudless and MAJA) showed a similar performance in terms of
 674 BOA; however, the tradeoff between PA and UA varied substantially for those algorithms:
 675 S2cloudless yielded PA=80.2% and UA=89.5% (more cloud omissions than commissions)
 676 and MAJA yielded PA=88.6% and UA=80.2% (less cloud omissions and more commissions,
 677 in part due to the dilation). When thin/semi-transparent clouds were not considered, all
 678 algorithms showed a better performance with an average gain in BOA of $+5.1 \pm 1.6\%$. Some
 679 algorithms (FORCE, Idepix and LaSRC) showed high commission errors (low UA), which
 680 were related to identifying snow as clouds.

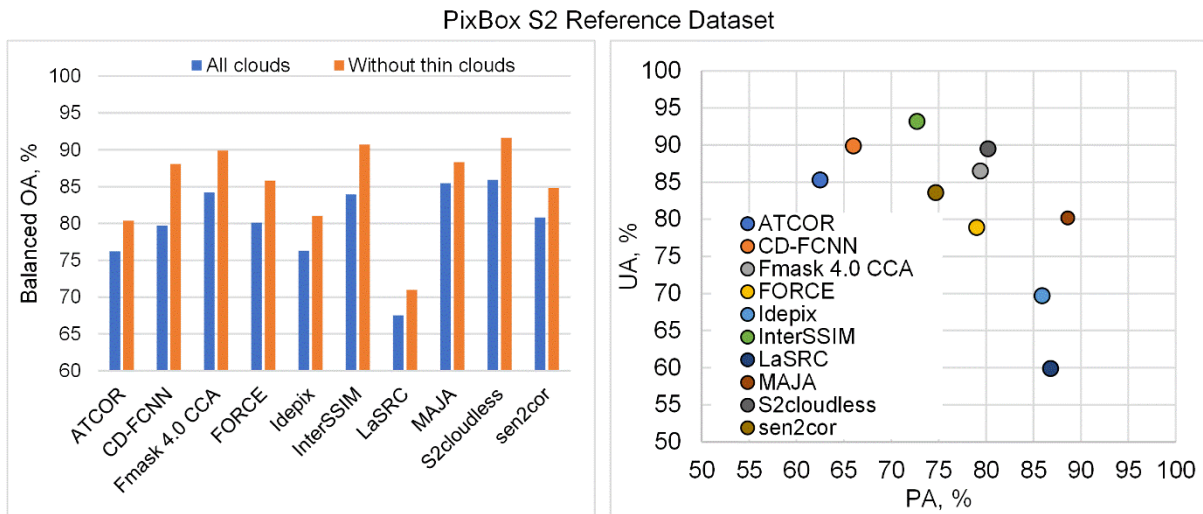
681

682 Table 9. Performance metrics of cloud masking algorithms for the complete PixBox S2
 683 dataset. ATCOR and MAJA processed 27 and 14 reference scenes, respectively, while other
 684 algorithms processed all 29 reference scenes. Fraction of cloud pixels was 47.2% and 36.8%
 685 for all cloud types and without thin clouds, respectively.

Processor	All types of clouds				Without thin clouds			
	OA	BOA	Cloud		OA	BOA	Cloud	
			PA	UA			PA	UA
ATCOR* (27/29)	76.6	76.2	62.5	85.3	82.5	80.4	70.8	81.4
CD-FCNN	80.5	79.7	66	89.9	89.5	88.1	82.7	87.9
Fmask 4.0 CCA	84.5	84.2	79.4	86.5	89.6	89.9	90.8	82.7

FORCE	80.2	80.1	79	78.9	84.6	85.8	90.4	73.6
Idepix	75.7	76.3	85.9	69.7	77.2	81	95.3	62.4
InterSSIM	84.6	84	72.7	93.2	91.9	90.7	86.2	91.3
LaSRC	66.4	67.5	86.8	59.9	65	71	93.8	51.3
MAJA* (14/29)	85.1	85.5	88.6	80.2	86.5	88.3	94.3	74.3
S2cloudless	86.3	85.9	80.2	89.5	91.6	91.6	91.6	86.4
sen2cor	81.2	80.8	74.7	83.6	85.4	84.8	82.7	78.6
Average	80.1	80.0	77.6	81.7	84.4	85.2	87.9	77.0
Standard deviation	5.7	5.3	8.3	9.6	7.7	6.0	7.1	11.7

686



687

688 Figure 12. Comparison of BOA values and distribution of PA/UA (for all clouds) for the
689 PiBox S2 reference dataset.

690

691 Table 10 shows BOA values when comparing complete and LCD PiBox dataset.
692 When restricting to the LCD, s2cloudless yielded the highest BOA in all cases Overall, the
693 differences in BOA between complete and LCD sets were below 2%. Also, algorithms
694 performance improved when thin clouds and snow were excluded from the analysis.

695

696 Table 10. Performance metrics of cloud masking algorithms for the complete and LCD
697 PiBox dataset for various scenarios.

Processor	All types of clouds		All types of clouds (excluding snow)		Without thin clouds	
	BOA	BOA	BOA	BOA	BOA	BOA

	complete	LCD	complete	LCD	complete	LCD
ATCOR	76.2	78.3	77.2	79.3	80.4	81.6
CD-FCNN	79.7	78.6	80.4	79.5	88.1	86.0
Fmask 4.0 CCA	84.2	85.1	86.3	86.9	89.9	89.7
FORCE	80.1	83.0	82.1	85.2	85.8	88.2
<u>Idepix</u>	76.3	73.8	84.0	83.0	81.0	78.8
InterSSIM	84.0	84.2	84.9	85.2	90.7	91.1
LaSRC	67.5	70.7	74.2	78.0	71.0	73.4
MAJA	85.5	85.5	86.1	86.1	88.3	88.3
S2cloudless	85.9	87.3	86.7	87.8	91.6	93.1
sen2cor	80.8	82.3	82.1	85.4	84.8	85.3
<i>Average</i>	<i>80.0</i>	<i>80.9</i>	<i>82.4</i>	<i>83.6</i>	<i>85.2</i>	<i>85.5</i>
<i>Standard deviation</i>	<i>5.3</i>	<i>5.1</i>	<i>3.9</i>	<i>3.3</i>	<i>6.0</i>	<i>5.7</i>

698

699

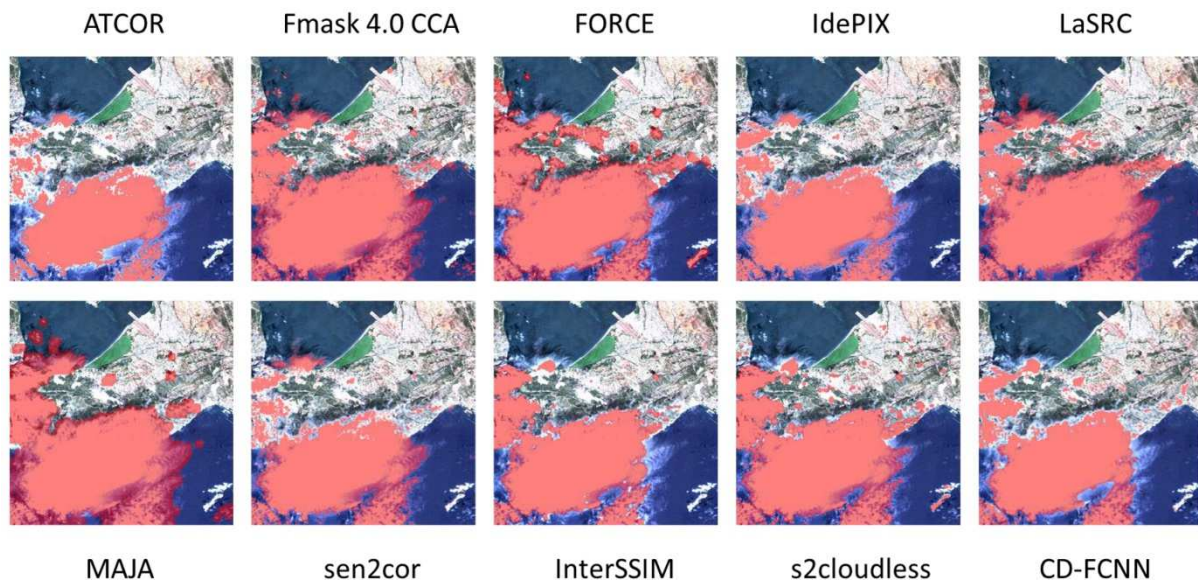
700

701

702

703

Figure 13 shows an example of cloud detection over the Sentinel-2 scene from the PixBox dataset. The scene features opaque clouds as well as semi-transparent clouds over the water. All algorithms were successful in identifying opaque clouds, while majority struggled to identify semi-transparent over the water.



704

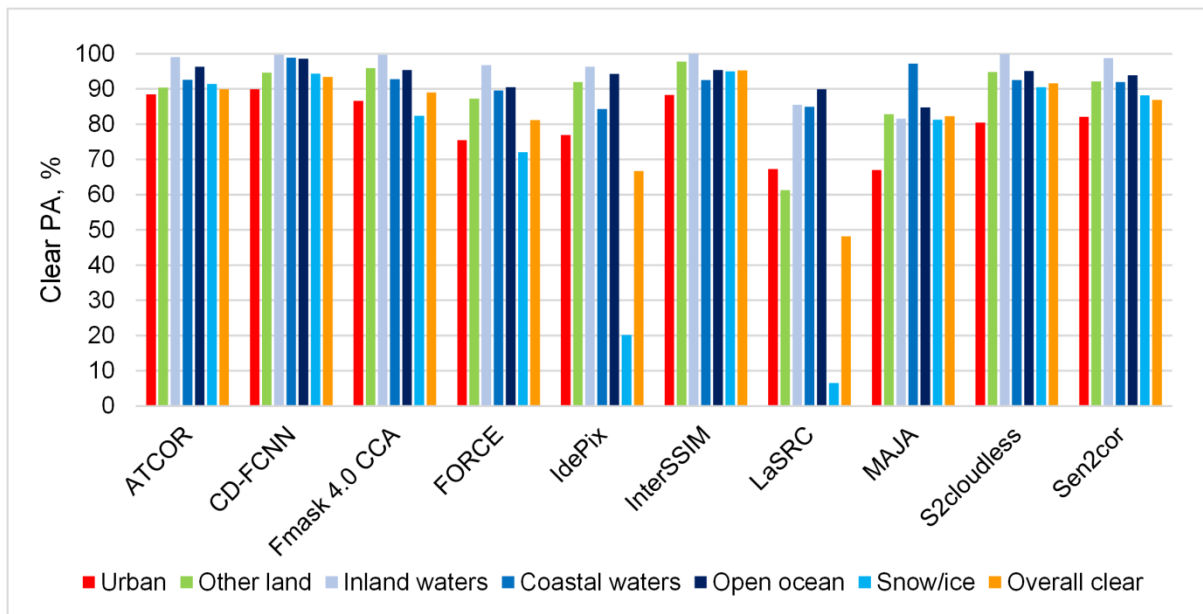
705

706

707

Figure 13. Examples of cloud masking by various algorithms over the Sentinel-2 scene S2A_MSIL1C_20170629T103021_N0205_R108_T31TFJ_20170629T103020.

708 Figure 14 shows performance of algorithms on clear pixels depending on the major
 709 land cover classes (proportion>4%) from the PixBox S2 data. LaSRC, IdePix and FORCE
 710 showed the worst performance for the clear snow pixels, which was expected given
 711 limitations of these algorithms. Excluding snow, overall performance of algorithms was
 712 uniform throughout the land clover classes. All algorithms showed worst performance for the
 713 urban area given the presence of bright targets. Even approaches utilizing the Sentinel-2
 714 multi-band parallax (e.g., FORCE, Frantz et al., 2018) over-detected clouds in the urban
 715 areas.



716
 717 Figure 14. Performance of algorithms in terms of clear producer’s accuracy over the non-
 718 cloudy regions depending on the land cover types in the PixBox S2 dataset.

719

720 **3.2 Performance of cloud masking algorithms for Landsat 8**

721 **3.2.1 GSFC L8 reference dataset**

722 This dataset included six Landsat 8 scenes and all algorithms showed high
 723 performance (Table 11). Fmask showed the highest values of performance metrics. Two
 724 algorithms achieved 100% cloud UA, meaning no cloud over-detection in this dataset.

725

726 Table 11. Performance metrics of cloud masking algorithms for the GSFC L8 dataset. All

727 algorithms processed six reference scenes (with 49.4% fraction of cloud in reference data).

Processor	Cloud			
	OA	BOA	PA	UA
ATCOR	97.3	97.3	94.8	99.8
CD-FCNN	97.3	97.3	94.6	100.0
Fmask 4.0 CCA	98.7	98.7	97.3	100.0
FORCE	98.2	98.1	96.5	99.7
<u>LaSRC</u>	96.5	96.5	94.8	98.0
<i>Average</i>	<i>97.6</i>	<i>97.6</i>	<i>95.6</i>	<i>99.5</i>
<i>Standard deviation</i>	<i>0.8</i>	<i>0.8</i>	<i>1.1</i>	<i>0.7</i>

728

729 *3.2.2 L8Biome reference dataset*

730 Table 12 provides a summary of performance metrics for the L8Biome dataset.

731 Results in this table should not be used directly for intercomparing algorithms for the

732 following reasons: (i) ATCOR processed only 86 images out of 96 images, since images in

733 polar regions were removed due to Sun elevation lower than 30°; (ii) LaSRC processed 80

734 images, since snow/ice scenes were not considered; (iii) all algorithms, except ATCOR, had

735 on average 2.4% pixels not classified—those pixels are on the boundary of the Landsat 8

736 scene, and do not have valid values for all spectral bands. In addition, since CD-FCNN was

737 trained on the L8Biome and the L8SPARCS datasets (80% and 20%, respectively), the CD-

738 FCNN results on this dataset are omitted in order to avoid overoptimistic (overfitted)

739 detection results. Fmask partially used L8Biome data to find optimal thresholds for some of

740 the rules, namely weight of cirrus cloud probability, spectral-contextual snow index, and

741 morphology-based post-processing (Qiu et al., 2019; personal communication, Zhe Zhu and

742 Shi Qiu, University of Connecticut, November 2021). Since the foundation of the Fmask

743 algorithm was developed well before the L8Biome dataset release, we still included Fmask

744 4.0 for the inter-comparison, though with caveats.

745

746 Table 12. Performance metrics of cloud masking algorithms for the L8Biome dataset.

747 ATCOR and LaSRC processed 86 (48.3% of clouds in reference data) and 80 (49.4%)

748 scenes, respectively, while Fmask and FORCE processed all 96 scenes (47.9%).

Processor	Cloud			
	OA	BOA	PA	UA
ATCOR* (86/96)	86.8	86.7	83.2	88.8
<u>Fmask 4.0 CCA</u>	90.0	90.2	93.6	86.6
FORCE	84.9	85.3	96.0	77.7
LaSRC* (80/96)	90.9	90.9	92.7	89.2
<i>Average</i>	<i>88.1</i>	<i>88.3</i>	<i>91.4</i>	<i>85.6</i>
<i>Standard deviation</i>	<i>2.4</i>	<i>2.3</i>	<i>4.9</i>	<i>4.7</i>

749

750 Table 13 provides a correct intercomparison between algorithms since the amount of

751 reference scenes and pixels used was the same. The average BOA was $90.0 \pm 1.4\%$ and

752 $91.5 \pm 1.8\%$ for all types of clouds and without thin clouds, respectively. Removing thin

753 clouds from the reference increases BOA and Cloud-PA accuracies by $+1.5 \pm 0.7\%$ and

754 $+3.0 \pm 1.4\%$, respectively.

755

756 Table 13. Performance metrics of cloud masking algorithms for the L8Biome dataset using

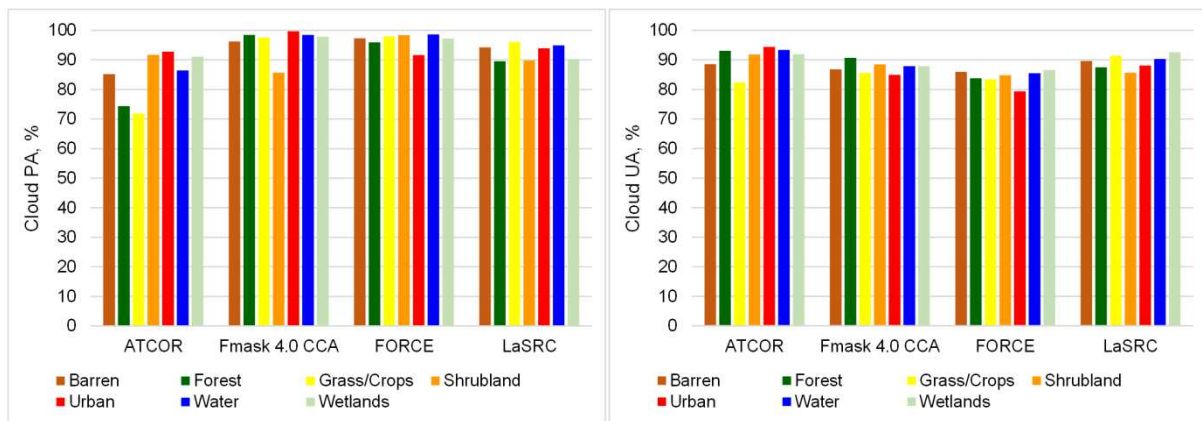
757 the same set of 80 Landsat 8 scenes. Fraction of cloud reference pixels for all types of clouds

758 and without thin clouds was 49.4% and 42.6%, respectively.

Processor	All types of clouds				Without thin clouds			
	Cloud				Cloud			
	OA	BOA	PA	UA	OA	BOA	PA	UA
ATCOR	88.2	88.2	84.6	90.9	89.6	89.2	86.8	88.6
<u>Fmask 4.0 CCA</u>	91.3	91.4	96.2	87.4	92.1	93.1	99.7	84.6
FORCE	89.4	89.5	96.8	84.2	89.0	90.2	98.1	80.4
LaSRC	90.9	90.9	92.7	89.2	92.8	93.5	97.8	86.9
<i>Average</i>	<i>89.9</i>	<i>90.0</i>	<i>92.6</i>	<i>87.9</i>	<i>90.9</i>	<i>91.5</i>	<i>95.6</i>	<i>85.1</i>
<i>Standard deviation</i>	<i>1.2</i>	<i>1.3</i>	<i>4.9</i>	<i>2.5</i>	<i>1.6</i>	<i>1.8</i>	<i>5.1</i>	<i>3.1</i>

759

760 Analysis of algorithms performance by biomes showed little variability (Figure 15).
 761 Exceptions are ATCOR which showed lower cloud PA values over forest and grass/cropland
 762 biomes, and Fmask which lower cloud PA values over shrubland. It is worth noting though
 763 that those are generic land cover classes and don't enable analysis of the dynamic state of the
 764 land cover class during the scene overpass. For example, a cropland can be characterized by
 765 multiple physical stages during the year, such as bare land (e.g., fallow or after ploughing),
 766 sparse vegetation (during crop emergence), dense vegetation (during peak), snow (during the
 767 winter period). Therefore, per-land cover performance of algorithms should be taken
 768 cautiously.
 769



770
 771 Figure 15. Performance of the Landsat 8 cloud detection algorithms for the L8Biome dataset
 772 depending on the biomes. The same set of 80 Landsat 8 scenes was used to calculate PA and
 773 UA accuracy values.
 774

775 3.2.3 PixBox L8 reference dataset

776 Table 14 shows the algorithm performance for the PixBox dataset. Fmask and
 777 ATCOR yielded the best performance in terms of BOA (87.9% and 86.3%, respectively),
 778 however PA/UA values exhibited a different behavior: for Fmask, PA and UA were mostly
 779 balanced (82.5% and 81.8%), while for ATCOR omission error (26.7%) was much higher

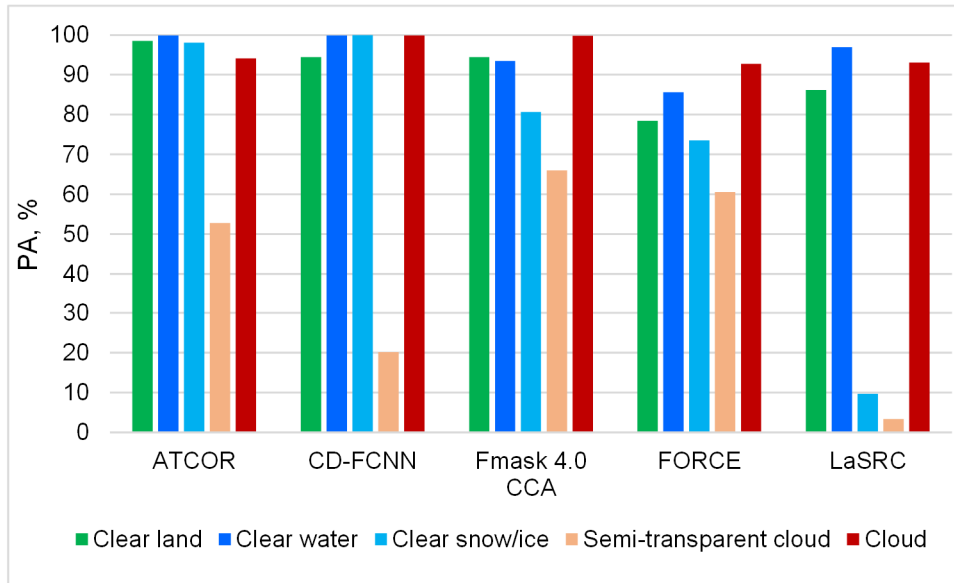
780 than commission error (2.8%). Overall, performance over the PixBox dataset was lower than
781 for L8Biome and GSFC, as the case with PixBox S2. Performance metrics substantially
782 improved when semi-transparent clouds were removed from the analysis. For all algorithms
783 cloud PA increased on average by $28.1\pm 13.9\%$ reaching $95.9\pm 3.6\%$. While there was an
784 overall agreement between algorithms on detecting opaque clouds from the PixBox L8
785 dataset (with average PA $95.9\pm 3.6\%$) all algorithms failed to detect semi-transparent clouds
786 (average PA was $40.6\pm 27.4\%$) (Figure 16). It's worth noting that all algorithms showed
787 equally good performance for clear land and water classes. ATCOR and CD-FCNN were also
788 successful in discriminating clouds from snow, while Fmask and FORCE showed
789 intermediate results. LaSRC failed to identify clouds over snow, as expected from the
790 algorithm's design.

791

792 Table 14. Performance metrics of cloud masking algorithms for the PixBox dataset. All
793 algorithms processed all 11 Landsat 8 reference scenes. Fraction of cloud reference pixels
794 was 27.4% for all types of clouds and 15.8%, when removing semi-transparent clouds.

Processor	<i>All types of clouds</i>				<i>Without semi-transparent clouds</i>			
	OA	BOA	Cloud		OA	BOA	Cloud	
			PA	UA			PA	UA
ATCOR	92.1	86.3	73.3	97.2	98.4	96.7	94.1	95.6
CD-FCNN	87.2	78.2	59	89.4	97.8	98.7	99.9	87.4
Fmask 4.0 CCA	90.4	87.9	82.5	81.8	94.3	96.6	99.8	72.6
FORCE	80.3	79.1	76.5	61.3	83.5	87.2	92.8	48.7
LaSRC	76.8	67.8	47.8	59.5	88.5	90.4	93.1	58.6
<i>Average</i>	83.7	78.2	66.5	73.0	92.5	93.9	95.9	72.6
<i>Standard deviation</i>	5.4	7.1	13.8	12.9	6.4	4.8	3.6	19.5

795



796

797 Figure 16. PA values for various types of classes in the PixBox L8 dataset.

798

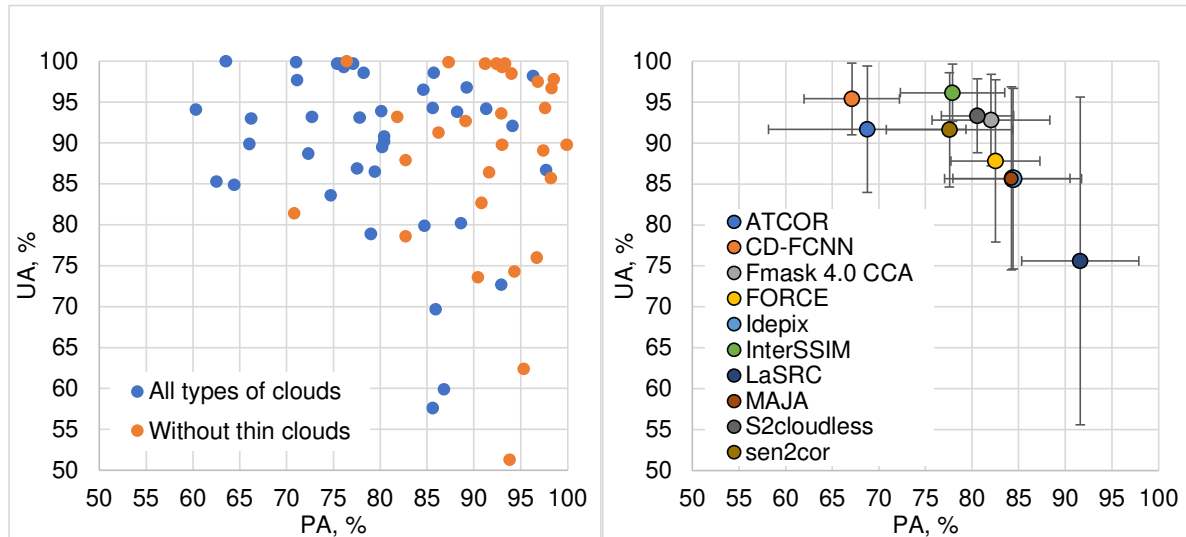
799 4 Discussion

800 4.1 Algorithm intercomparison

801 Figure 17 shows the distribution of cloud PA and UA values for Sentinel-2 cloud
 802 masking algorithms. Overall, cloud PA/UA values are located in the areas defined by lines
 803 $PA > 80\%$ or $UA > 80\%$. While individual values are located in the area of $PA > 90\%$ and
 804 $UA > 90\%$ (Figure 17, left), suggesting a very good balance of commission and omission
 805 errors, however that is not the case for averaged values across all reference datasets (Figure
 806 17, right). No algorithm yielded the $PA > 90\%$ and $UA > 90\%$ performance when averaging
 807 over reference datasets. Five algorithms (Fmask, FORCE, Idepix, MAJA and S2cloudless)
 808 yielded the average performance of cloud $PA > 80\%$ and $UA > 80\%$, providing some balance
 809 (within $\sim 10\%$) between commission and omission errors. Four algorithms (ATCOR, CD-
 810 FCNN, InterSSIM and sen2cor) yielded performance with cloud $UA > 90\%$ (cloud
 811 conservative), meaning these algorithms committed less clouds over clear regions, however
 812 at the expense of missing clouds. LaSRC yielded the cloud $PA > 90\%$ performance (non-cloud
 813 conservative), detecting most of the clouds, however, at the expense of masking out also

814 valid non-cloudy observations, and with a large standard deviation in UA across the datasets
 815 (potentially, owing to various rules defining the cloud and the use of conservative threshold
 816 for the cirrus band).

817



818

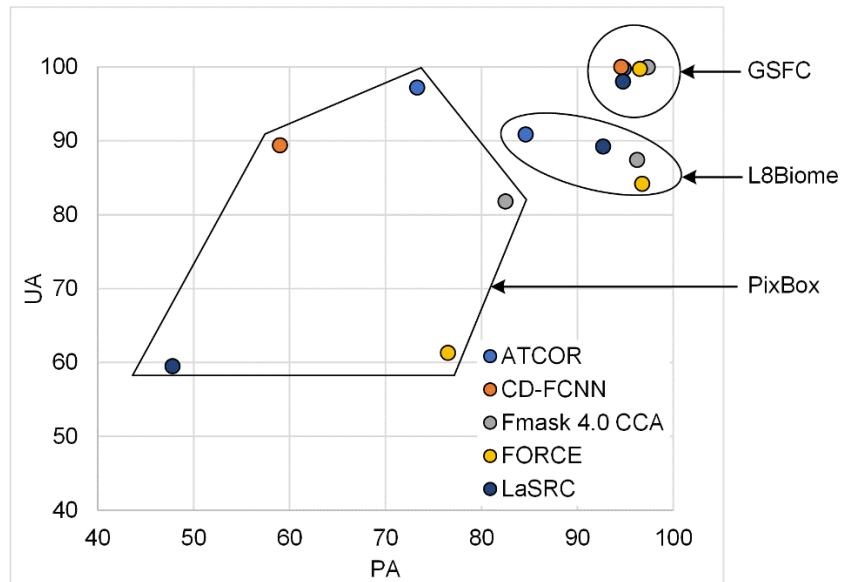
819 Figure 17. Distribution of cloud PA and UA over all Sentinel-2 cloud masking algorithms
 820 and reference datasets (left) and algorithms' average values along with the standard deviation
 821 over four reference datasets (right). Averaging was performed using PA and UA values from
 822 Table 6, Table 7, Table 8 and Table 9 for all cloud types.

823

824 Since only three datasets were used for Landsat 8, we did not perform the averaging
 825 (Figure 18). Three distinct clusters corresponding to the three reference datasets were evident
 826 with varying performance. The highest performance was for the GSFC dataset with only six
 827 Landsat 8 scenes over the same area, which probably is not fully representative of the
 828 performance of the algorithms. GSFC L8 had mostly thick and well-identifiable clouds that
 829 algorithms were able to classify successfully. L8Biome yielded the second highest
 830 performance with PA/UA values distributed over PA>90% (Fmask, FORCE and LaSRC) and
 831 UA>90% (ATCOR). Performance for the PixBox dataset was the lowest with algorithms

832 scattered in the cloud PA/UA space. Fmask yielded PA>80% and UA>80% for PixBox;
 833 ATCOR and CD-FCNN yielded UA>90%; while FORCE and LaSRC yielded both cloud PA
 834 and UA less than 80%.

835



836

837 Figure 18. Distribution of cloud PA and UA over all Landsat 8 clouds masking algorithms
 838 and reference datasets.

839

840 A summary of strengths and weaknesses of cloud algorithms known at the design
 841 stage and further identified/elaborated during the CMIX are presented in Table 15.

842

843 Table 15. Summary of algorithms parameters to control cloud commission/omission errors
 844 along with strengths and weaknesses.

Processor	Parameter	Strengths	Weaknesses
ATCOR	Cloud buffer size (default size is 7 px). Increase will lead to higher cloud PA.	<ul style="list-style-type: none"> – Water vapor map (S2) is used to reduce cloud pixel commission error – Elevation-dependent cirrus masking 	<ul style="list-style-type: none"> – Conservative cloud mask – Cloud buffer too small – Thin cirrus threshold of $\rho(\text{TOA})=0.01$ underestimates thin cirrus
CD-FCNN	A posteriori cloud probability (default value is 0.5). Decrease will lead to	<ul style="list-style-type: none"> – Single architecture to provide global cloud masks for both Landsat-8 	<ul style="list-style-type: none"> – Model can underperform compared to customized algorithms for Sentinel-2

	higher cloud PA (cloud-free conservative). Increase will lead to higher cloud UA (cloud conservative).	and Sentinel-2 images – No ancillary data required – Mitigation of training data requirements: transfer learning from Landsat-8 to Sentinel-2 – General approach directly learnt from available data	– Model performance is fully constrained by the quality of training data – Presence of errors in thin clouds, cloud borders, urban areas, and snow. – It does not provide shadow detection. – It does not provide cloud type classes (e.g. cirrus, thin or thick clouds).
Fmask 4.0 CCA	Cloud dilation (default is 3 px), cloud probability threshold (CPT), and potential false positive cloud (PFPC) extension and erosion. The CPT default value is 17.5% for Landsat 8, and 20% for Sentinel 2. Increase will reduce the number of potential cloud pixels. The PFPC parameters affect how the potential cloud mask is reduced to the final cloud mask. Changing its values will affect the algorithm's performance over bright targets.	– Generic algorithm – Applicable over land and water – Good performance over bright targets (urban, ice/snow)	– Performance decreases when thermal band is not used
FORCE	Cloud probability (default 22.5%). Increase will reduce the number of potential cloud pixels. Clouds were buffered by 300 m. Higher values will increase cloud commission but reduce commission.	– Rigorous cloud mask with emphasis on reducing cloud commission for safe usage in time series applications – Parallax effect is used to reduce bright false positives in Sentinel-2 imagery – Multiple flags can be set, e.g. snow and cloud	– Rigorous cloud mask with emphasis on reducing cloud commission with potential drawbacks for single-scene analysis – Parallax effect may occasionally introduce false positives in bright areas due to micro-vibrations on sensor – Snow and cloud often not mutually exclusively
IdePix	The CLOUD_AMBIGUOUS flag is currently quite probe to clear commission of urban and other very bright surfaces. Cloud buffer was not used, as it would increase cloud commission error.	– Mono-temporal approach – Detects thin clouds quite well – Allows user defined cloud dilation	– Snow detection could be better (bug in code during CMIX) – Commission error of bright (mostly urban) surfaces
s2cloudless	Cloud probability (default is 0.4). Lower values will lead to higher cloud PA (cloud-free conservative). Post-processing:	– Fast single-observation cloud masking – Works on any resolution and even on aggregated values (objects)	– Prone to errors on very bright areas – No spatial context is taken into account – No cloud shadow detection

	convolution (22 px) and dilation (11 px). The convolution smoothens the masks, reducing the amount of salt-and-pepper effect, while the dilation of masks closes small openings and increases the cloud masks on the outside.	<ul style="list-style-type: none"> – Provides pseudo-probability that user can tweak to get better cloud masks for her use-case 	
InterSSIM	Similar to s2cloudless. Number of prior satellite observations. Increase will lead to better performance, especially bright targets, but increase the usage of computational and storage resources.	<ul style="list-style-type: none"> – Using spatio-temporal context results in lower rate of false positive detections (particularly over consistently bright areas) – Provides pseudo-probability that user can tweak to get better cloud masks for her use-case 	<ul style="list-style-type: none"> – Resource intensive calculation – Higher rate of cirrus misclassifications – Higher rate of misclassifications over large waterbodies – No cloud shadow detection
LaSRC	Threshold for residuals from aerosol retrievals (default is 0.05). Increase will lead to higher cloud UA (cloud conservative).	<ul style="list-style-type: none"> – Simple, interpretable criteria Easily transferable – Conservative and tune to keep best high-quality data rather than questionable (low-quality) 	<ul style="list-style-type: none"> – Might confuse bad retrievals of aerosol with clouds (high aerosol, urban area) – Not suitable over snow cover region
MAJA	Four major parameters: <ul style="list-style-type: none"> – Multi-temporal: threshold on increase of surface reflectance in the blue. – Correlation: each neighborhood of a cloud is correlated with previous observations. If the correlation is high, it is not a cloud. – High clouds: threshold for the reflectance of the cirrus band, that depends on the squared altitude of the pixel to account for the fact that mountains may peak above the water vapor layer. – Buffer: all pixels close to a cloud within a buffer of 240 m are classified as clouds, which is rather conservative, and avoids omissions due to the parallax between spectral bands or to fuzzy limits of the cloud. 	<ul style="list-style-type: none"> – Multi-temporal criterion to better detect low clouds that brings a much better separation between cloud / non clouds – Moderate threshold for the cirrus bands, as the multi-temporal threshold already detects clouds which have a significant impact on reflectances – Large buffer (240m), possible thanks to the very low level of cloud commission errors before dilation 	<ul style="list-style-type: none"> – Some very rapid changes of vegetation could be interpreted as clouds – Multi-temporal algorithm is less efficient in places where the cloudiness is extremely high – Working at 120 m resolution (240 m resolution during CMIX, but it has been upgraded since), may cause omissions of very small clouds – The buffer will include some cloud free pixels (but they are in fact are affected by large adjacency effects)
sen2cor	The parameters used to run Sen2Cor version 2.8 for	<ul style="list-style-type: none"> – Cloud mask at “moderate” resolution (20 	<ul style="list-style-type: none"> – Potential cloud omissions on cloud edges/boundaries

	CMIX were the default parameters used in Sentinel-2 operational ground segment and available in L2A_CAL_SC_GIPP.xml. No cloud mask dilation is applied and cloud boundaries can be omitted.	m) – Robustness. Used operationally in all types of meteorological conditions and solar geometries – Processing time (<5 min for a full Sentinel-2 tile)	– Potential cloud omissions for cloud over water – Potential cloud commissions for bright buildings in urban area or bright surfaces
--	---	--	---

845

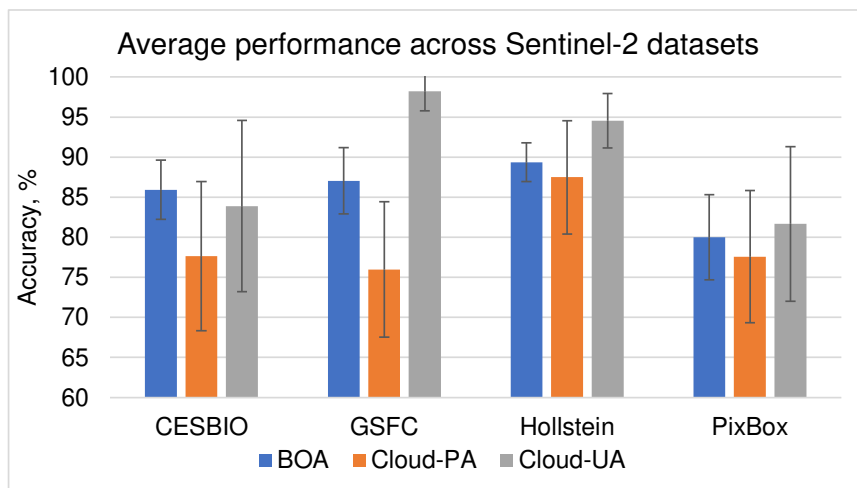
846 *4.2 Dependence of the performance on the reference datasets*

847 Performance of cloud masking algorithms for Sentinel-2 varied depending on the
 848 reference dataset (Figure 19): average BOA was $80.0 \pm 5.3\%$ (PixBox) to $89.4 \pm 2.4\%$
 849 (Hollstein). Performance of algorithms was the worst for the PixBox dataset compared to
 850 datasets. This can be explained by the following. PixBox dataset was sampled in such a way,
 851 so non-challenging (e.g., opaque thick clouds) and challenging (e.g., semi-transparent clouds,
 852 cloud boundaries) cases are equally present in the dataset. At the same time, other datasets
 853 were aimed at labelling the full images (L8Biome, CESBIO) or provide homogeneous
 854 polygons (Hollstein, GSFC), where the weight of challenging cases would be lower than for
 855 PixBox. In this regard, the question is about whether to weight samples according to the area
 856 or not. Both characteristics (based on equal allocation and area proportions) can be valuable
 857 to describe separability of classes by a given algorithm (model accuracy) and to estimate
 858 probability of a pixel being mapped correctly (map accuracy) (Blickensdörfer et al., 2022;
 859 Congalton, 1991).

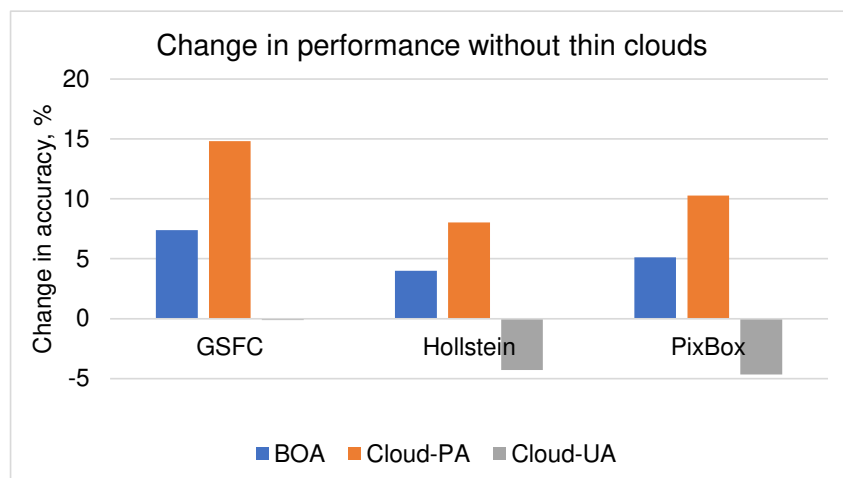
860 Across the four reference datasets algorithms showed better performance in terms of
 861 cloud UA, which was consistently higher than cloud PA. Removing thin/semi-transparent
 862 clouds from the reference datasets improves performance of algorithms (Figure 20), though at
 863 the expense of cloud UA. This happens because thin clouds have higher uncertainties and
 864 therefore are more challenging to the algorithms in contrast to thick clouds. When thin clouds
 865 removed from reference datasets the proportion of correctly detected classes increases and

866 therefore cloud PA increases. At the same time, cloud UA can experience both increase or
 867 decrease depending on the proportion of thin clouds and algorithm's performance on thin
 868 clouds.

869 The issue of thin/semi-transparent cloud detection has a significant impact on the
 870 subsequent shadow detection. Figure 21 shows an example of a cloud with different levels of
 871 transparency depending on wavelength used and its shadow. While the cloud is semi-
 872 transparent in the false color composite (SWIR-NIR-red), its shadow is clearly visible and
 873 impacts the reflectance.



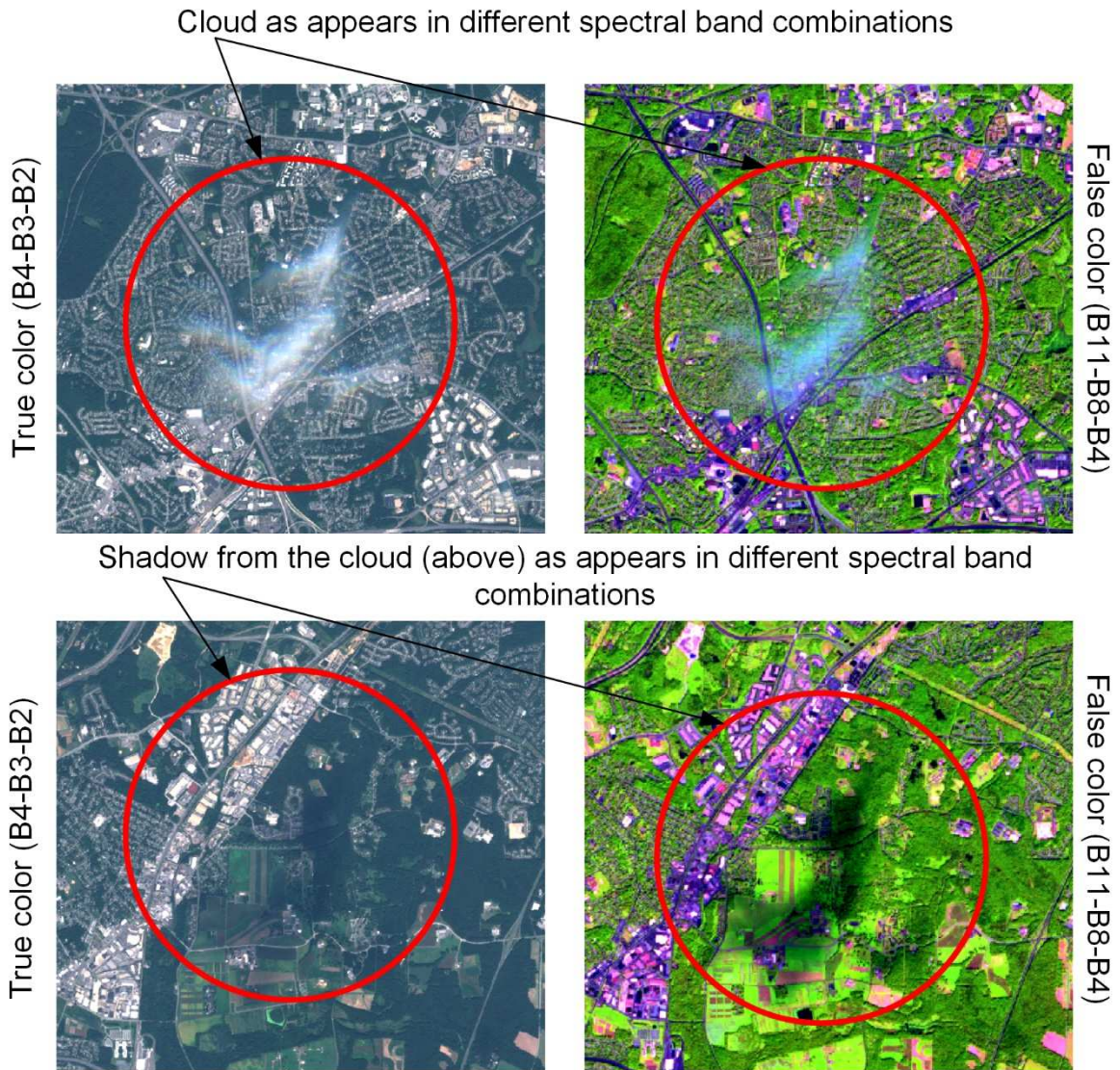
874
 875 Figure 19. Average performance of algorithms for Sentinel-2 for four cloud reference
 876 datasets.



877

878 Figure 20. Change in performance of Sentinel-2 cloud masking algorithms, when thin/semi-
879 transparent clouds removed from the reference datasets.

880



881

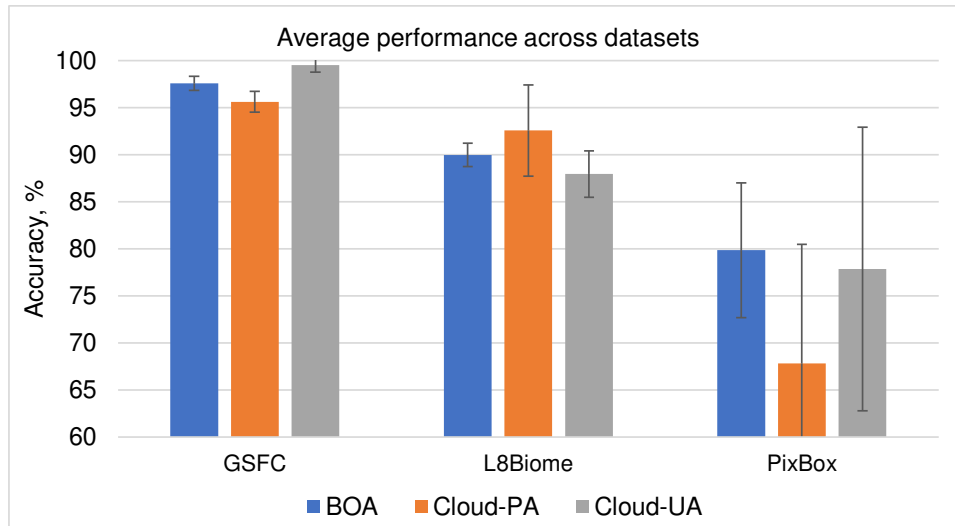
882 Figure 21. Example of thin/semi-transparent cloud in various band combinations (true color
883 and false color in top-of-atmosphere reflectance) along with the shadow from that cloud
884 (Sentinel-2 scene, L1C_T18SUJ_A011777_20170923T160124).

885

886 Figure 22 shows averaged BOA values across multiple Landsat 8 algorithms. As with
887 Sentinel-2, the performance varied across datasets yielding BOA of $97.6\pm 0.8\%$, $90.0\pm 1.3\%$

888 and $79.8 \pm 7.1\%$ for GSFC, L8Biome and PixBox, respectively. As with Sentinel-2, cloud PA
889 was higher than cloud UA for GSFC and PixBox datasets, but not for L8Biome.

890



891

892 Figure 22. Average performance of algorithms for Landsat 8 for three cloud reference
893 datasets.

894

895 In terms of various land cover classes, it is difficult to draw conclusions since only
896 generic “static” information on land cover was available for some of the datasets. We did not
897 observe any substantial differences in algorithm’s performance over various land cover
898 classes, except for urban areas in the PixBox S2 data, which is expected. Sentinel-2 does not
899 have a thermal band and, therefore, detection of clouds over bright targets in urban areas
900 remains a challenging task. The use of multi-spectral parallax (Skakun et al., 2017) only
901 partially addresses this problem (Frantz et al., 2018).

902

903 **5 Recommendations**

904 Results and lessons learned from CMIX-I provide a good foundation for future
905 activities for improving practices related to the development and validation of cloud masking
906 algorithms for passive optical satellite imagery.

907 The first area for improvement should aim at initially providing an agreed upon
908 definition of “cloud” (Mejia et al., 2016; Stubenrauch et al., 2013) that is passed beforehand
909 to intercomparison participants and validation dataset originators. Ideally this would be an
910 objective (quantitative) definition of clouds, which would include a numerical metric. As
911 results from CMIX-I showed, existing validation datasets varied in how a cloud was defined
912 through mostly photointerpretation, and it influenced the performance of the algorithms. For
913 example, one potential metric to define the cloud would be the cloud optical thickness.
914 However, this poses the questions at which wavelength the thickness should be defined, what
915 threshold to apply, and how it could be estimated for sizeable quantity of images. For
916 example, Mejia et al. (2016) use a radiative transfer model to estimate cloud optical depth (τ_c)
917 from ground-based sky images and define thick clouds with $\tau_c > 30$, thin clouds with $\tau_c \sim 1$, and
918 clear sky with $\tau_c \sim 0$ (all in the visible spectrum). While there was a consensus between
919 algorithms and developers in defining thick non-transparent clouds, there was a disagreement
920 (sometimes by design and depending on the intended applications) in transparent (semi-
921 transparent) clouds, such as cirrus, stratus and cloud edges. Also, the effect of those clouds
922 can vary with wavelengths, which adds complexity to the analysis.

923 Based on the cloud definition, the second area for improvement would include
924 generation of new reference/validation datasets. The strengths and weaknesses of existing
925 cloud reference datasets were thoroughly analyzed and discussed within this study, and new
926 datasets should substantially address those weaknesses. A special attention should be paid to
927 ensure a balanced statistical distribution of surface and cloud types, as well as the need to
928 cover a wide range of environmental conditions, in order to thoroughly test the performance
929 of the algorithms at global scale. Some of the recommendations include:

- 930 • Consistently implementing the cloud definition, and adding cloud shadows to the
931 analysis. Recommended practices for labelling clouds should be developed and

932 implemented for new datasets, whether through visual interpretation or ground
933 measurements or ancillary data (e.g. geostationary satellites). Cloud shadows should
934 be also part of the analysis, since an inaccurate cloud shadow mask can lead to
935 substantial artifacts in the downstream products.

936 • Defining a proper dilation of cloud masks to be applied, taking account the effect of
937 parallax between spectral bands, smooth variation of clouds at their borders, and
938 adjacency effects. The dilation could then be applied to the reference datasets and to
939 the algorithm results.

940 • Increasing the number of sites collecting ground-based imagery of the sky and use
941 them in coordination with Aeronet measurements. Some limitations of the use of
942 ground-based sky imagery include radiance contrast which could yield better
943 detection of thin clouds; furthermore, the geometrical matching between sky-camera
944 and satellite pixel may introduce some errors, which are related to the cloud height.

945 • Acquire multiple datasets (time-series) over the same area to analyze consistent errors
946 in cloud detection. This would enable temporal metrics to be exploited when
947 assessing the efficiency of cloud masks.

948 The third set of activities should focus on expanding the analysis framework, which
949 would include:

950 • A sample-based approach versus an area-based approach, when comparing reference
951 cloud mask with a predicted one. The problem with an area-based approach is that
952 more weight would be given to large clouds (which cover the larger area), whereas
953 smaller clouds might have a small impact on the performance metrics. At the same
954 time, sampled-based approaches can also miss some specific land cover features
955 (unless a stratification scheme can be constructed with strata describing those
956 features), and often do not address the boundaries of the clouds or more broadly

957 segmentation aspects. Area-based approaches are likewise necessary to study the
958 effects of cloud dilation. Therefore, both approaches should be considered.

959 • Temporal analysis of cloud masks over the same area. Originally planned for CMIX-I,
960 the idea of using temporal metrics was abandoned, since no reference data (except
961 GSFC, which were assisted with sky imagery and Aeronet measurements) was
962 available for these purposes. As undetected clouds add noise on time-series, it is
963 possible to evaluate the noise on time-series and compute the contribution of different
964 cloud masks to this noise.

965 • Application-based approach to cloud validation. One way to analyze efficiency of the
966 cloud/shadow masks is to “validate” them indirectly within the downstream products.
967 An example could include a generic land cover mapping workflow, when the same set
968 of satellite data will be processed by various cloud detection algorithms and used as
969 input to the classification algorithm. The derived land cover maps will be validated
970 using the same validation data and intercompared.

971 And finally, CMIX-I was limited to Landsat 8 and Sentinel-2 data. Future activities
972 could include adding hyper-spectral data (such as PRISMA or DESIS), coarse resolution data
973 (such as MODIS, VIIRS, Sentinel-3), and commercial very high spatial resolution satellites,
974 such as Planet or hyperspectral sensors.

975

976 **6 Conclusion**

977 The Cloud Mask Intercomparison eXercise (CMIX) was a community-wide effort to
978 intercompare the state-of-the-art and commonly-used cloud masking algorithms, with a focus
979 on moderate spatial resolution data acquired by Landsat 8 and Sentinel-2 missions. Ten
980 algorithms developed by nine teams from fourteen organizations representing universities,
981 industry and space agencies were evaluated within CMIX using existing cloud reference data.

982 Overall, the performance of algorithms varied depending on the reference dataset, which can
983 be attributed to differences in which reference datasets were generated. Average overall
984 accuracy (across algorithms) varied $80.0\pm 5.3\%$ to $89.4\pm 2.4\%$ for Sentinel-2, and $79.8\pm 7.1\%$
985 to $97.6\pm 0.8\%$ for Landsat 8, depending on the reference dataset. An overall accuracy of 90%
986 yields twice less errors than an overall accuracy of 80%. The study highlighted algorithms
987 that provided a balance between commission and omission errors, as well as algorithms
988 which are cloud conservative (high UA) and non-cloud (clear) conservative (high PA). With
989 repetitive observations like those of Sentinel-2, it seems reasonable to favor cloud
990 conservative approaches, with maybe the exception of very cloudy regions where every cloud
991 free observation is critical. When thin/semi-transparent clouds were not considered in the
992 reference datasets algorithms' performance generally improved: overall accuracy values
993 increased from +1.5% to 7.4%. It should be noted though that these clouds are commonly
994 occurring and are often present in optical imagery. We concluded the paper with
995 recommendations for further activities, which include provision of a quantitative definition
996 for clouds (targeting moderate spatial resolution imagery by Landsat 8 and Sentinel-2),
997 generation of new reference datasets, and expansion of the analysis framework (for example,
998 multi-temporal analysis and application-driven validation). Such intercomparison studies will
999 hopefully help the community to improve the algorithms and move towards standardization
1000 of cloud masking. Given the importance of cloud masking in optical imagery we encourage
1001 CEOS to continue the CMIX activities.

1002

1003

1004 **Acknowledgment**

1005 We would like to thank to Chris Justice (University of Maryland) for helpful
1006 comments on an earlier draft of paper and Gasmine Myers (University of Maryland) for

1007 proof-reading the paper. L.G.C., D.L.P. and G.M.G. (University of Valencia) were supported
1008 for this work by the Spanish Ministry of Science and Innovation (project PID2019-
1009 109026RB- I00, ERDF) and the European Social Fund. S.S., J.C.R. (University of Maryland)
1010 and E.V. (NASA GSFC) were supported by NASA grants 80NSSC19K1592,
1011 80NSSC19M0222 and 80NSSC21M0080.
1012

1013 **Appendix A.**

1014

1015 Table A1. Cloud and non-cloud classes that were used from the original reference datasets.

Dataset	Cloud	Non-cloud
CESBIO	Low clouds, high clouds	Shadow, land, water, snow
GSFC	Cloud, thin cloud	Clear, cloud shadow
Hollstein	Cloud, cirrus	Clear, water, shadow, snow
L8Biome	Thin cloud, thick cloud	Shadow, clear
PixBox S2	Opaque, thick semi-transparent cloud, average density semi-transparent cloud, semi-transparent cloud, thin semi-transparent cloud, fog, haze	Clear
PixBox L8	Cloud, semi-transparent cloud	Clear land, clear snow/ice, clear water, mixed snow_ice/water

1016

1017

1018 **References**

- 1019 Baetens, L., Desjardins, C., Hagolle, O., 2019. Validation of Copernicus Sentinel-2 Cloud
1020 Masks Obtained from MAJA, Sen2Cor, and FMask Processors Using Reference
1021 Cloud Masks Generated with a Supervised Active Learning Procedure. *Remote Sens.*
1022 11, 433.
- 1023 Blickensdörfer, L., Schwieder, M., Pflugmacher, D., Nendel, C., Erasmi, S., Hostert, P.,
1024 2022. Mapping of crop types and crop sequences with combined time series of
1025 Sentinel-1, Sentinel-2 and Landsat 8 data for Germany. *Remote Sens. Environ.* 269,
1026 112831.
- 1027 Brodersen, K.H., Ong, C.S., Stephan, K.E., Buhmann, J.M., 2010. The balanced accuracy and
1028 its posterior distribution, in: *Proc. 2010 20th International Conference on Pattern*
1029 *Recognition. IEEE*, pp. 3121–3124.
- 1030 Chai, D., Newsam, S., Zhang, H. K., Qiu, Y., Huang, J., 2019. Cloud and cloud shadow
1031 detection in Landsat imagery based on deep convolutional neural networks. *Remote*
1032 *Sens. Environ.* 225, 307–316.
- 1033 Congalton, R.G., 1991. A review of assessing the accuracy of classifications of remotely
1034 sensed data. *Remote Sens. Environ.* 37(1), 35-46.
- 1035 Doxani, G., Vermote, E., Roger, J.C., Gascon, F., Adriaensen, S., Frantz, D., Hagolle, O.,
1036 Hollstein, A., Kirches, G., Li, F., Louis, J., 2018. Atmospheric correction inter-
1037 comparison exercise. *Remote Sens.* 10(2), 352.
- 1038 EO Research Team, 2020. Cloud Masks at Your Service. [https://medium.com/sentinel-](https://medium.com/sentinel-hub/cloud-masks-at-your-service-6e5b2cb2ce8a)
1039 [hub/cloud-masks-at-your-service-6e5b2cb2ce8a](https://medium.com/sentinel-hub/cloud-masks-at-your-service-6e5b2cb2ce8a) (accessed 11 July 2021).
- 1040 Foga, S., Scaramuzza, P.L., Guo, S., Zhu, Z., Dilley Jr, R.D., Beckmann, T., Schmidt, G.L.,
1041 Dwyer, J.L., Hughes, M.J., Laue, B., 2017. Cloud detection algorithm comparison and

1042 validation for operational Landsat data products. *Remote Sens. Environ.* 194, 379–
1043 390.

1044 [Dataset] Francis, A., Mrziglod, J., Sidiropoulos, P., Muller, J.-P., 2020. Sentinel-2 Cloud
1045 Mask Catalogue (Version 1.0). Zenodo. <https://doi.org/10.5281/zenodo.4172871>

1046 Frantz, D., 2019. FORCE—Landsat + Sentinel-2 Analysis Ready Data and Beyond. *Remote*
1047 *Sens.* 11, 1124.

1048 Frantz, D., Haß, E., Uhl, A., Stoffels, J., Hill, J., 2018. Improvement of the Fmask algorithm
1049 for Sentinel-2 images: Separating clouds from bright surfaces based on parallax
1050 effects. *Remote Sens. Environ.* 215, 471–481.

1051 Frantz, D., Röder, A., Stellmes, M., Hill, J., 2016. An Operational Radiometric Landsat
1052 Preprocessing Framework for Large-Area Time Series Applications. *IEEE Trans.*
1053 *Geosc. Remote Sens.* 54 (7), 3928–3943.

1054 Frantz, D., Röder, A., Udelhoven, T., Schmidt, M., 2015. Enhancing the detectability of
1055 clouds and their shadows in multitemporal dryland Landsat imagery: Extending
1056 Fmask. *IEEE Geosci. Remote Sens. Lett.* 12(6), 1242–1246.

1057 Gascon, F., Bouzinac, C., Thépaut, O., Jung, M., Francesconi, B., Louis, J., Lonjou, V.,
1058 Lafrance, B., Massera, S., Gaudel-Vacaresse, A., Languille, F., 2017. Copernicus
1059 Sentinel-2A calibration and products validation status. *Remote Sens.* 9(6), 584.

1060 Hagolle, O., Huc, M., Pascual, D.V., Dedieu, G., 2010. A multi-temporal method for cloud
1061 detection, applied to FORMOSAT-2, VEN μ S, LANDSAT and SENTINEL-2 images.
1062 *Remote Sens. Environ.* 114(8), 1747–1755.

1063 [Dataset] Hagolle, O., Huc, M., Desjardins, C., Auer, S., Richter, R., 2017. MAJA Algorithm
1064 Theoretical Basis Document (Version 1.0). Zenodo.
1065 <http://doi.org/10.5281/zenodo.1209633>

1066 Holben, B.N., Eck, T.F., Slutsker, I.A., Tanre, D., Buis, J.P., Setzer, A., Vermote, E., et al.
1067 1998. AERONET—A federated instrument network and data archive for aerosol
1068 characterization. *Remote Sens. Environ.* 66(1), 1–16.

1069 Hollingsworth, B.V., Chen, L., Reichenbach, S.E., Irish, R.R., 1996. November). Automated
1070 cloud cover assessment for Landsat TM images. In: *Imaging Spectrometry II*, Vol.
1071 2819. International Society for Optics and Photonics, pp. 170–179.

1072 Hollstein, A., Segl, K., Guanter, L., Brell, M., Enesco, M., 2016. Ready-to-use methods for
1073 the detection of clouds, cirrus, snow, shadow, water and clear sky pixels in Sentinel-2
1074 MSI images. *Remote Sens.* 8(8), 666.

1075 Hughes, M.J., Hayes, D.J., 2014. Automated detection of cloud and cloud shadow in single-
1076 date Landsat imagery using neural networks and spatial post-processing. *Remote*
1077 *Sens.* 6 (6), 4907–4926.

1078 Irish, R.R., Barker, J.L., Goward, S.N., Arvidson, T., 2006. Characterization of the Landsat-7
1079 ETM+ automated cloud-cover assessment (ACCA) algorithm. *Photogramm. Eng.*
1080 *Remote Sens.* 72 (10), 1179–1188.

1081 Jeppesen, J.H., Jacobsen, R.H., Inceoglu, F., Toftegaard, T.S., 2019. A cloud detection
1082 algorithm for satellite imagery based on deep learning. *Remote Sens. Environ.* 229,
1083 247–259.

1084 Kotchenova, S.Y., Vermote, E.F., Matarrese, R., Klemm Jr, F.J., 2006. Validation of a vector
1085 version of the 6S radiative transfer code for atmospheric correction of satellite data.
1086 Part I: Path radiance. *Appl. Opt.* 45(26), 6762–6774.

1087 Lamarche, C., Santoro, M., Bontemps, S., d’Andrimont, R., Radoux, J., Giustarini, L.,
1088 Brockmann, C., Wevers, J., Defourny, P., Arino, O., 2017. Compilation and
1089 validation of SAR and optical data products for a complete and global map of

1090 inland/ocean water tailored to the climate modeling community. *Remote Sens.* 9(1),
1091 36.

1092 López-Puigdollers, D., Mateo-García, G., Gómez-Chova, L., 2021. Benchmarking Deep
1093 Learning Models for Cloud Detection in Landsat-8 and Sentinel-2 Images. *Remote*
1094 *Sens.* 13(5), 992.

1095 Louis, J., Debaecker, V., Pflug, B., Main-Knorn, M., Bieniarz, J., Mueller-Wilm, U., Cadau,
1096 E., Gascon, F., 2016. Sentinel-2 sen2cor: L2a processor for users, in: *Proceedings*
1097 *Living Planet Symposium 2016*. Spacebooks Online, pp. 1–8.

1098 Louis, J., 2021. Sentinel-2 Level-2A Algorithm Theoretical Basis Document.
1099 [https://sentinels.copernicus.eu/documents/247904/446933/Sentinel-2-Level-2A-](https://sentinels.copernicus.eu/documents/247904/446933/Sentinel-2-Level-2A-Algorithm-Theoretical-Basis-Documents-ATBD.pdf)
1100 [Algorithm-Theoretical-Basis-Documents-ATBD.pdf](https://sentinels.copernicus.eu/documents/247904/446933/Sentinel-2-Level-2A-Algorithm-Theoretical-Basis-Documents-ATBD.pdf) (accessed 9 July 2021)

1101 Mateo-García, G., Laparra, V., López-Puigdollers, D., Gómez-Chova, L., 2020. Transferring
1102 deep learning models for cloud detection between Landsat-8 and Proba-V. *ISPRS J.*
1103 *Photogramm. Remote Sens.* 160, 1–17.

1104 Mejia, F.A., Kurtz, B., Murray, K., Hinkelman, L.M., Sengupta, M., Xie, Y., Kleissl, J.,
1105 2016. Coupling sky images with radiative transfer models: a new method to estimate
1106 cloud optical depth. *Atm. Meas. Techn.* 9(8), 4151–4165.

1107 Pahlevan, N., Mangin, A., Balasubramanian, S.V., Smith, B., Alikas, K., Arai, K., Barbosa,
1108 C., et al. 2021. ACIX-Aqua: A global assessment of atmospheric correction methods
1109 for Landsat-8 and Sentinel-2 over lakes, rivers, and coastal waters. *Remote Sens.*
1110 *Environ.* 258, 112366.

1111 [Dataset] Paperin, M., Wevers, J., Stelzer, K., Brockmann, C., 2021a. PixBox Sentinel-2
1112 pixel collection for CMIX (Version 1.0). Zenodo.
1113 <https://doi.org/10.5281/zenodo.5036991>

1114 [Dataset] Paperin, M., Stelzer, K., Lebreton, C., Brockmann, C., Wevers, J., 2021b. PixBox
1115 Landsat 8 pixel collection for CMIX (Version 1.0). Zenodo.
1116 <https://doi.org/10.5281/zenodo.5040271>

1117 Puc, J., Žust, L., 2019. On cloud detection with multi-temporal data.
1118 [https://medium.com/sentinel-hub/on-cloud-detection-with-multi-temporal-data-](https://medium.com/sentinel-hub/on-cloud-detection-with-multi-temporal-data-f64f9b8d59e5)
1119 [f64f9b8d59e5](https://medium.com/sentinel-hub/on-cloud-detection-with-multi-temporal-data-f64f9b8d59e5) (accessed 09 June 2021)

1120 Qiu, S., Zhu, Z., He, B., 2019. Fmask 4.0: Improved cloud and cloud shadow detection in
1121 Landsats 4–8 and Sentinel-2 imagery. *Remote Sens. Environ.* 231, 111205.

1122 Richter, R., Schläpfer, D., 2019a. ATCOR-3 User Guide. Version 9.3.0. [https://www.rese-](https://www.rese-apps.com/pdf/atcor3_manual.pdf)
1123 [apps.com/pdf/atcor3_manual.pdf](https://www.rese-apps.com/pdf/atcor3_manual.pdf) (accessed 09 June 2021)

1124 Richter, R., Schläpfer, D., 2019b. ATCOR Theoretical Background Document. Version 1.0.
1125 https://www.rese-apps.com/pdf/atcor_ATBD.pdf (accessed 03 July 2021).

1126 Scaramuzza, P.L., Bouchard, M.A., Dwyer, J.L., 2012. Development of the Landsat data
1127 continuity mission cloud-cover assessment algorithms. *IEEE Trans. Geosci, Remote*
1128 *Sens.* 50(4), 1140–1154.

1129 Segal-Rozenhaimer, M., Li, A., Das, K., Chirayath, V., 2020. Cloud detection algorithm for
1130 multi-modal satellite imagery using convolutional neural-networks (CNN). *Remote*
1131 *Sens. Environ.* 237, 111446.

1132 Skakun, S., Vermote, E.F., Artigas, A.E.S., Rountree, W.H., Roger, J.-C., 2021. An
1133 experimental sky-image-derived cloud validation dataset for Sentinel-2 and Landsat 8
1134 satellites over NASA GSFC. *Int. J. Appl. Earth Observ. Geoinform.* 95, 102253.

1135 Skakun, S., Vermote, E., Roger, J.-C., Justice, C., 2017. Multispectral misregistration of
1136 Sentinel-2A images: Analysis and implications for potential applications. *IEEE*
1137 *Geosci. Remote Sens. Lett.* 14(12), 2408–2412.

1138 Skakun, S., Vermote, E.F., Roger, J.-C., Justice, C.O., Masek, J.G., 2019. Validation of the
1139 LaSRC cloud detection algorithm for Landsat 8 images. *IEEE J. Sel. Topics Appl.*
1140 *Earth Observ. Remote Sens.* 12(7), 2439–2446.

1141 Stubenrauch, C.J., Rossow, W.B., Kinne, S., Ackerman, S., Cesana, G., Chepfer, H., Di
1142 Girolamo, L., Getzewich, B., Guignard, A., Heidinger, A., Maddux, B.C., *et al.*, 2013.
1143 Assessment of global cloud datasets from satellites: Project and database initiated by
1144 the GEWEX radiation panel. *Bull. American Meteorol. Soc.* 94(7), 1031–1049.

1145 Tarrío, K., Tang, X., Masek, J.G., Claverie, M., Ju, J., Qiu, S., Zhu, Z., Woodcock, C.E.,
1146 2020. Comparison of cloud detection algorithms for Sentinel-2 imagery. *Sci. Remote*
1147 *Sens.* 2, 100010.

1148 [Dataset] U.S. Geological Survey, 2016. L8 SPARCS Cloud Validation Masks. U.S.
1149 Geological Survey data release. USGS. doi:10.5066/F7FB5146

1150 Vermote, E., Justice, C., Claverie, M., Franch, B., 2016. Preliminary analysis of the
1151 performance of the Landsat 8/OLI land surface reflectance product. *Remote Sens.*
1152 *Environ.* 185, 46–56.

1153 Vermote, E.F., Tanré, D., Deuze, J.L., Herman, M., Morcette, J.J., 1997. Second simulation
1154 of the satellite signal in the solar spectrum, 6S: An overview. *IEEE Trans Geosci.*
1155 *Remote Sens.* 35(3), 675–686.

1156 Wevers, J., Müller, D., Scholze, J., Kirches, G., Quast, R., Brockmann, C., 2021. IdePix for
1157 Sentinel-2 MSI Algorithm Theoretical Basis Document (Version 1.0). Zenodo.
1158 <https://doi.org/10.5281/zenodo.5788067>

1159 Wieland, M., Li, Y., Martinis, S., 2019. Multi-sensor cloud and cloud shadow segmentation
1160 with a convolutional neural network. *Remote Sens. Environ.* 230, 111203.

1161 Wulder, M.A., Loveland, T.R., Roy, D.P., Crawford, C.J., Masek, J.G., Woodcock, C.E.,
1162 Allen, R.G., et al., 2019. Current status of Landsat program, science, and applications.
1163 Remote Sens. Environ. 225, 127–147.

1164 Xie, F., Shi, M., Shi, Z., Yin, J., Zhao, D., 2017. Multilevel cloud detection in remote sensing
1165 images based on deep learning. IEEE J. Sel. Topics Appl. Earth Observ. Remote
1166 Sens. 10(8), 3631–3640.

1167 Zhu, Z., Wang, S., Woodcock, C.E., 2015. Improvement and expansion of the Fmask
1168 algorithm: Cloud, cloud shadow, and snow detection for Landsats 4–7, 8, and Sentinel
1169 2 images. Remote Sens. Environ. 159, 269–277.

1170 Zhu, Z., Woodcock, C.E., 2012. Object-based cloud and cloud shadow detection in Landsat
1171 imagery. Remote Sens. Environ. 118, 83–94.

1172 Zhu, Z., Woodcock, C.E., 2014. Automated cloud, cloud shadow, and snow detection in
1173 multitemporal Landsat data: An algorithm designed specifically for monitoring land
1174 cover change. Remote Sens. Environ. 152, 217–234.

1175 Zupanc, A., 2017. Improving Cloud Detection with Machine Learning.
1176 [https://medium.com/sentinel-hub/improving-cloud-detection-with-machine-learning-](https://medium.com/sentinel-hub/improving-cloud-detection-with-machine-learning-c09dc5d7cf13)
1177 [c09dc5d7cf13](https://medium.com/sentinel-hub/improving-cloud-detection-with-machine-learning-c09dc5d7cf13) (accessed 09 June 2021)

1178

1179 **List of Figure Captions**

1180

1181 Figure 1. Geographical distribution of the Landsat 8 and Sentinel-2 scenes in the reference
1182 datasets used in CMIX.

1183

1184 Figure 2. Distribution of labeled pixels in the CESBIO dataset.

1185

1186 Figure 3. Distribution of labeled pixels in the GSFC S2 dataset (left) and land cover classes
1187 (right).

1188

1189 Figure 4. Distribution of labeled pixels in the Hollstein dataset.

1190

1191 Figure 5. Distribution of labeled pixels in the L8Biome dataset.

1192

1193 Figure 6. Distribution of labeled pixels and land cover classes in the PixBox dataset.

1194

1195 Figure 7. Part of the L8Biome scene (LC81570452014213LGN00) with some thin clouds not
1196 labelled. Thin clouds are shown in orange, and thick clouds in maroon.

1197

1198 Figure 8. Examples of labeled data in the three datasets: CESBIO (fully labeled images);
1199 GSFC (polygons avoiding uncertain areas, such cloud boundaries); PixBox (sample-based
1200 approach).

1201

1202 Figure 9. Comparison of BOA values and distribution of PA/UA for the CESBIO reference
1203 dataset.

1204

1205 Figure 10. Comparison of BOA values and distribution of PA/UA (for all clouds) for the
1206 GSFC S2 reference dataset.

1207

1208 Figure 11. Comparison of BOA values and distribution of PA/UA (for all clouds) for the
1209 Hollstein reference dataset.

1210

1211 Figure 12. Comparison of BOA values and distribution of PA/UA (for all clouds) for the
1212 PixBox S2 reference dataset.

1213

1214 Figure 13. Examples of cloud masking by various algorithms over the Sentinel-2 scene
1215 S2A_MSIL1C_20170629T103021_N0205_R108_T31TFJ_20170629T103020.

1216

1217 Figure 14. Performance of algorithms in terms of clear producer's accuracy over the non-
1218 cloudy regions depending on the land cover types in the PixBox S2 dataset.

1219

1220 Figure 15. Performance of the Landsat 8 cloud detection algorithms for the L8Biome dataset
1221 depending on the biomes. The same set of 80 Landsat 8 scenes was used to calculate PA and
1222 UA accuracy values.

1223

1224 Figure 16. PA values for various types of classes in the PixBox L8 dataset.

1225

1226 Figure 17. Distribution of cloud PA and UA over all Sentinel-2 cloud masking algorithms
1227 and reference datasets (left) and algorithms' average values along with the standard deviation

1228 over four reference datasets (right). Averaging was performed using PA and UA values from
1229 Table 6, Table 7, Table 8 and Table 9 for all cloud types.

1230

1231 Figure 18. Distribution of cloud PA and UA over all Landsat 8 clouds masking algorithms
1232 and reference datasets.

1233

1234 Figure 19. Average performance of algorithms for Sentinel-2 for four cloud reference
1235 datasets.

1236

1237 Figure 20. Change in performance of Sentinel-2 cloud masking algorithms, when thin/semi-
1238 transparent clouds removed from the reference datasets.

1239

1240 Figure 21. Example of thin/semi-transparent cloud in various band combinations (true color
1241 and false color in top-of-atmosphere reflectance) along with the shadow from that cloud
1242 (Sentinel-2 scene, L1C_T18SUJ_A011777_20170923T160124).

1243

1244 Figure 22. Average performance of algorithms for Landsat 8 for three cloud reference
1245 datasets.